

# Talend connect to AWS and Snowflake Tutorial

In this tutorial we will be examining connections from Talend to various AWS services and the Snowflake Data Warehouse package. In AWS we will use S3, SQS and EC2 services together with Lambda functions. To demonstrate the usage the following scenario will be modelled.

A retail company employs staff at several different sites which open 12 hours daily with a staffing FTE (full time equivalent) level of 5. Each store has a pool of 20 staff who work 6-hour shifts and to fulfil the FTE, any 10 of the 20 staff will be required each day. The timesheets for each day are produced by an external system that periodically sends a data file in JSON format to an Amazon S3 bucket owned by the company, who will then load the information from this file into their Snowflake database, for further analysis.

In the example a single timesheet file for each week will be produced for all stores in a simulated process, as is the creation of staff records. These simulations are useful for demonstrating Talend concepts.

## Requirements

To undertake this assignment, the following software and accounts are required.

- Talend Big Data Studio community edition version 7.3 or later.
- Amazon Web Services account, all examples are free tier eligible
- Snowflake account, trial version used in this tutorial.
- Java 8 JRE for installation onto AWS EC2 instance.

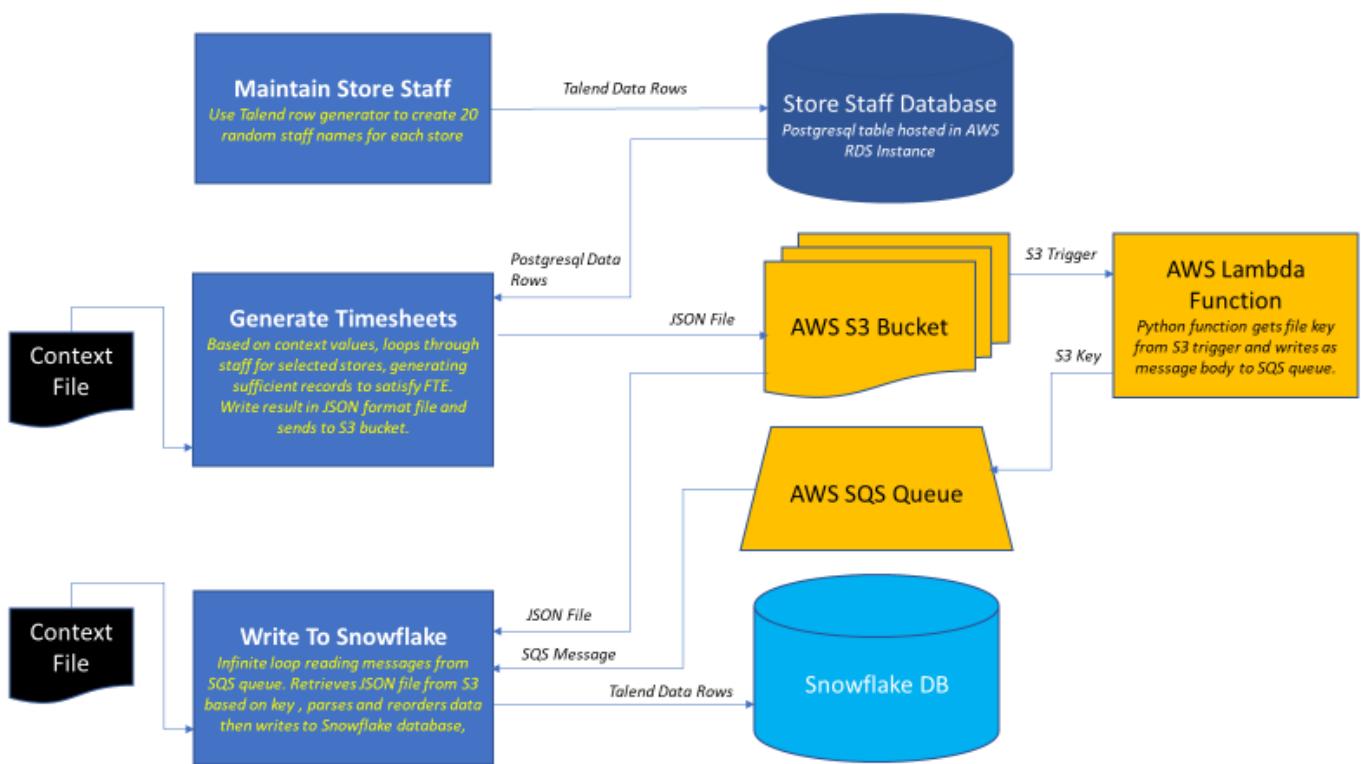
## Design Overview

Store staff records will be created by a Talend job and stored in a Postgresql table hosted via RDS an AWS. A second Talend job running on an EC2 instance will extract data from the staff database and use it to simulate the generation of timesheets. These will be written to a JSON format file and be copied to an S3 bucket. An AWS Lambda function written in Python responds to an S3 trigger event raised by placement of the file into the bucket and writes the file key to the body of a message, which is placed on an AWS SQS queue. A final Talend job also running on EC2 responds to the SQS message and retrieves the file key from the message body.

The key is used to download the JSON file from S3 and its contents are parsed and reordered before being written to a table in Snowflake. Both Talend jobs running on the EC2 instance use context files to load the initial variables

The use of the Lambda function allows serverless technology to be used to automate the snowflake load on receipt of a timesheet file.

The following illustration shows the design concept.



## Task List

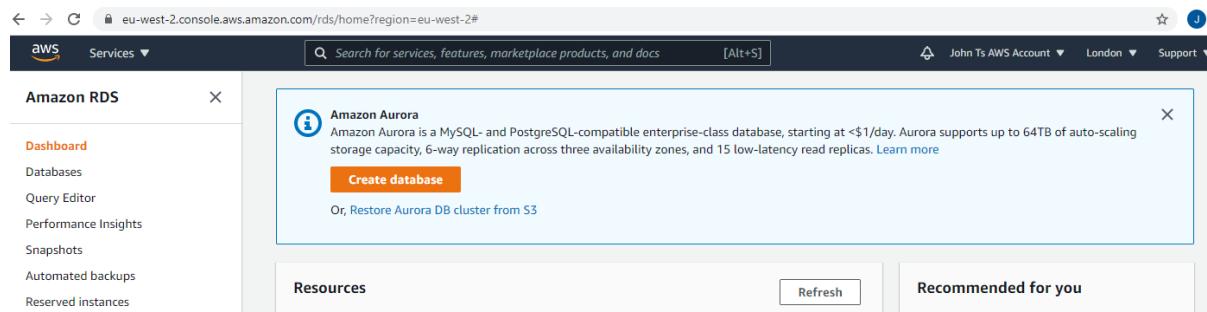
Before generating the Talend jobs for this project there are several configuration tasks to be undertaken. The following table provides a check list of these:

No	Task	Description
1	Create Postgresql instance on AWS RDS	Working in the AWS console for RDS, create a new Postgresql instance.
2	Add database and table to Postgres	Using the instance created previously create a database and table within it to hold details of store staff.
3	Setup Snowflake schema.	Create a new database in Snowflake and add a table to store staff hours.
4	Create AWS S3 bucket	Create a bucket in the S3 console for timesheet files to be loaded into.
5	Generate AWS Keys	Generate AWS key pair and download for use in connecting to S3 and SQS.
6	Create an EC2 Instance	Create AWS EC2 instance for Talend standalone jobs to run on.
7	Install Java 8 JRE onto EC2	Install the java runtime onto the EC2 instance created previously. This will allow Talend standalone jobs to run on the server.
8	Create S3 role for EC2	Create a role with access to S3 buckets and assign the role to the EC2 instance to allow applications access to S3.
9	Create AWS SQS queue	Create a FIFO queue in AWS which will be written to by Lambda and consumed by Talend.
10	Create AWS Lambda function	Create a serverless Lambda function in AWS triggered by S3 input that write file key from the trigger into an SQS message. Code will be written in Python, directly into AWS console.
11	Create job to generate store staff	Create Talend job to generate store staff.
12	Create job to generate timesheets	Talend job to create timesheets and write as a JSON file to the S3 bucket. This will trigger the Lambda, writing a message to SQS.
13	Create job to write to Snowflake	Job that loops indefinitely, reading SQS queue and retrieving data from S3. Output goes to Snowflake database.
14	Deploy jobs to EC2	Deploy the 2 Talend jobs that will be ran on the EC2 server.

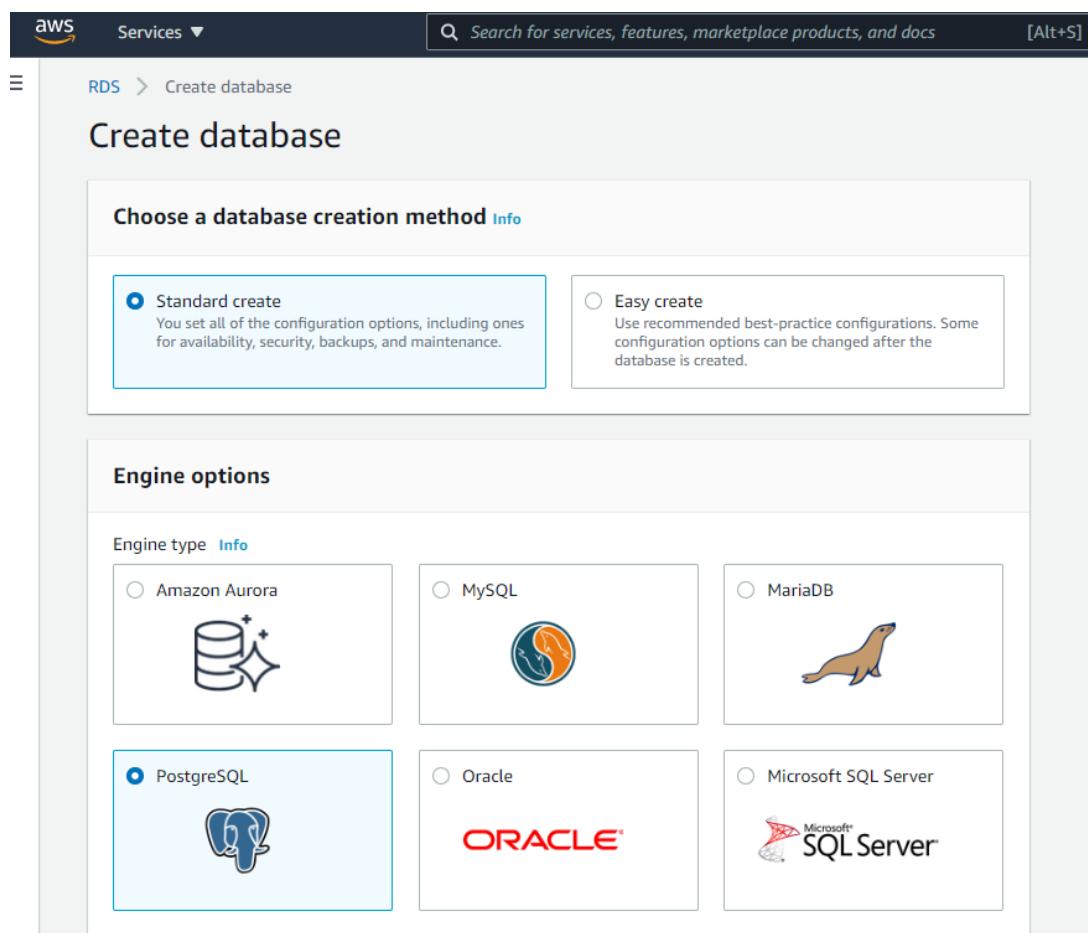
Each task will be examined in turn.

## 1. Create Postgresql instance on AWS RDS

To add a new PostgreSQL database into AWS RDS an instance must be created. Log onto AWS management console and enter the RDS option. Click on the **Create Database** button.



In the Create Database dialog select **PostgreSQL** from the Engine options box and check the box to use the **Free tier template**.



## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

Give the instance a name, in this case jt-dbpostgresretailpoc, accept all defaults and click **Create**. The instance will be created, and this can be checked by selecting **Databases** from the left-hand side menu and clicking on the **instance name**. Information regarding the instance including endpoint and port details will be displayed.

The screenshot shows the AWS RDS instance details page for 'jt-dbpostgresretailpoc'. The left sidebar shows navigation options like Dashboard, Databases (selected), Query Editor, Performance Insights, Snapshots, Automated backups, Reserved instances, Proxies, Subnet groups, Parameter groups, Option groups, Events, Event subscriptions, Recommendations, and Certificate update. The main content area has tabs for Summary, Connectivity & security (selected), Monitoring, Logs & events, Configuration, Maintenance & backups, and Tags. Under the 'Connectivity & security' tab, it displays the Endpoint (jt-dbpostgresretailpoc.c22aqflm1aig.eu-west-2.rds.amazonaws.com), Port (5432), Networking (Availability zone: eu-west-2c, VPC: vpc-39f4a551, Subnet group: default-vpc-39f4a551), and Security (VPC security groups: default (sg-87ee7ae1) active). The instance identifier is jt-dbpostgresretailpoc, CPU usage is 4.00%, Status is Available, Class is db.t2.micro, and Region & AZ is eu-west-2c.

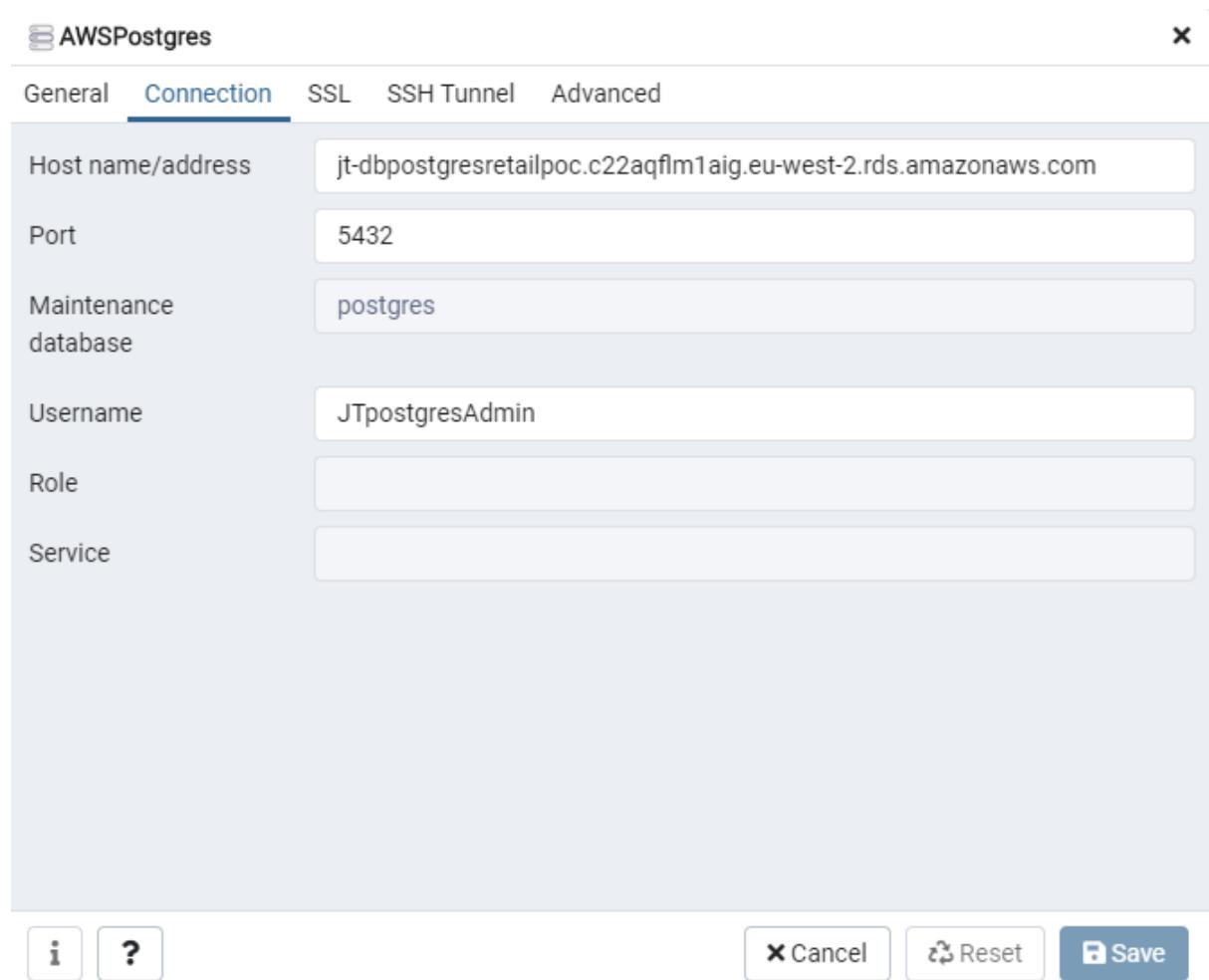
## 2. Add database and table to Postgres

The previous step has launched a PostgreSQL instance in AWS, but we now need to create a database and table within the instance. This is achieved by connecting to the instance from outside AWS by using the PGAdmin tool, the standard IDE for PostgreSQL.

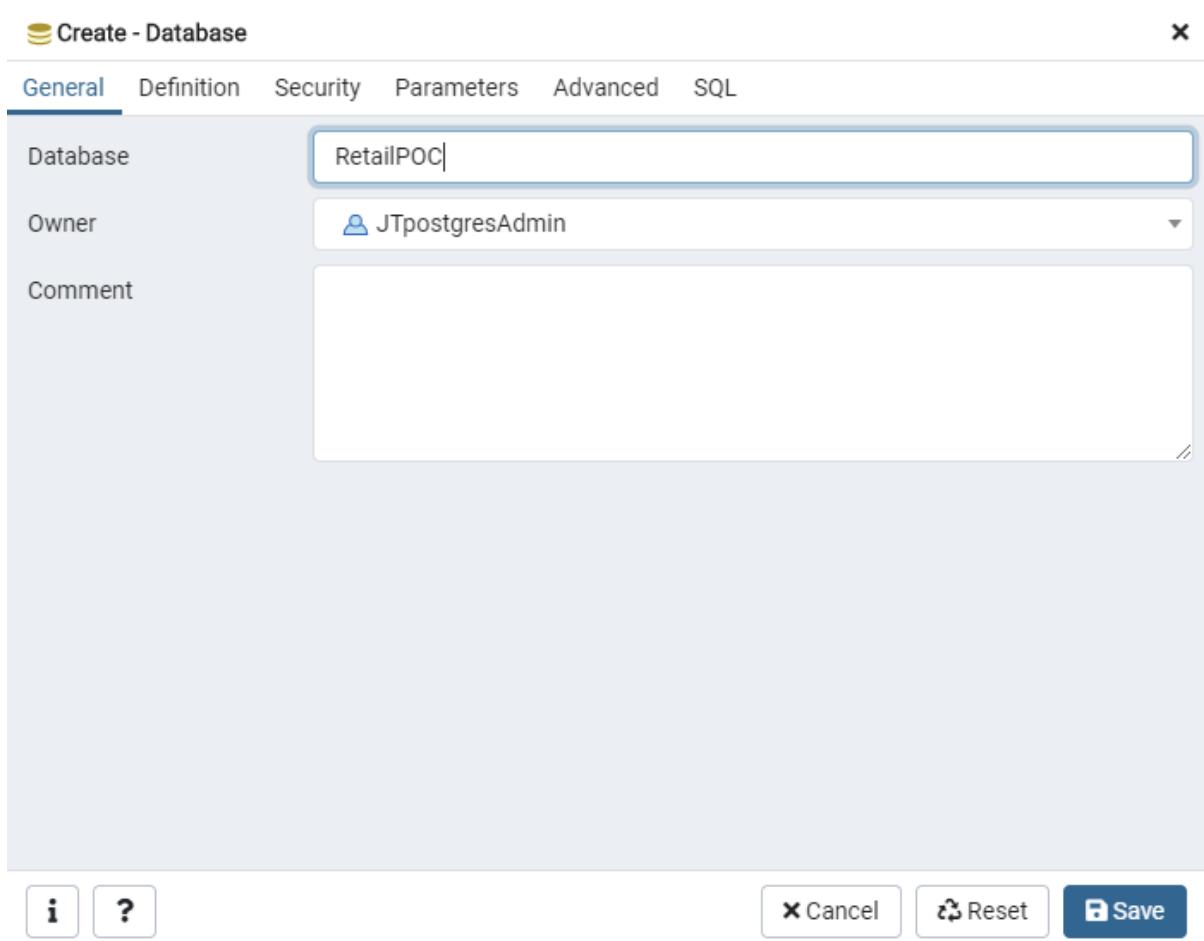
The screenshot shows the PGAdmin 'Create - Server' dialog. The 'General' tab is selected. The 'Name' field is set to 'AWSPostgres'. The 'Server group' dropdown is set to 'Servers'. The 'Background' and 'Foreground' checkboxes are unchecked. The 'Connect now?' checkbox is checked. The 'Comments' text area is empty. At the bottom are buttons for 'Cancel', 'Reset', and 'Save'.

In the PGAdmin browser pain right click on **Servers** and select **Create > Server**. In the first tab of the dialog box enter the server name AWSPostgres.

Select **Connection** from the menu and enter the end point as defined in the AWS RDS instance screen from the previous section, confirming that the port number is also correct. The maintenance database can be left at the default value of postgres and the username/password should be the ones defined in AWS.



Working again in the Browser pane, expand the AWSPostgres server and right click on the **Databases** node. Select **Create > Database** from the menu and name it RetailPOC prior to saving.



Right click the **RetailPOC** database in the Browser and select **Query Tool** from the menu.

Copy the following SQL statement into the query window then execute it:

```
CREATE TABLE public.tbstorestaff
(
    employee_no integer NOT NULL,
    store_code integer,
    first_name text COLLATE pg_catalog."default",
    last_name text COLLATE pg_catalog."default",
    CONSTRAINT tbstorestaff_pkey PRIMARY KEY (employee_no)
)
```

The table structure should now be visible in the Browser hierarchy as below:

Browser

The screenshot shows a hierarchical tree view of a PostgreSQL database named 'RetailPOC'. At the top level, there are several categories: Casts, Catalogs (2), Event Triggers, Extensions, Foreign Data Wrappers, Languages, and Schemas (1). Under 'Schemas (1)', there is a single entry for 'public'. Below 'public', there are numerous objects: Collations, Domains, FTS Configurations, FTS Dictionaries, FTS Parsers, FTS Templates, Foreign Tables, Functions, Materialized Views, Procedures, Sequences, and Tables (1). Under 'Tables (1)', there is one entry for 'tbstorestaff'. This table has four columns: employee\_no, store\_code, first\_name, and last\_name. It also has one constraint: tbstorestaff\_pkey.

The PostgreSQL table is now ready to receive data.

### 3. Setup Snowflake schema.

Log in to the Snowflake account and select the Databases option from the ribbon menu. Click on **Create** and enter the name of the database RETAIL\_POC in the name field and press **Finish**.

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

The screenshot shows the Snowflake web console interface. At the top, there's a navigation bar with icons for Databases, Shares, Data Marketplace, Warehouses, Worksheets, and History. Below the navigation bar, a banner says "Enjoy your free trial! Visit our documentation to learn more about using Snowflake or c...". The main area is titled "Databases" and shows four databases: RETAIL\_POC, SNO, DEMO, and UTIL. There are buttons for "Create...", "Clone...", "Drop...", and "Transfer Ownership". A search bar shows "4 databases". A modal dialog box is open, titled "Create Database", with the "Name \*" field set to "RETAIL\_POC". It also has a "Comment" field and a "Show SQL" button. At the bottom of the dialog are "Cancel" and "Finish" buttons.

The database will now be shown in the list as below.

The screenshot shows the same Snowflake web console interface as before, but now it displays the newly created database "RETAIL\_POC" in the list. The table header includes "Database", "Origin", "Creation Time", "Owner", and "Comment". The "RETAIL\_POC" entry shows "4/13/2021, 1:53 AM" as the creation time and "SYSADMIN" as the owner.

Click the **Snowflake icon** at the left-hand end of the ribbon to open a query window. Select the RETAIL\_POC database and the PUBLIC schema, then paste the following code into the query window:

```
CREATE OR REPLACE TABLE "RETAIL_POC"."PUBLIC"."TBSTAFFHOURS" (
    STORE_CODE NUMBER(38,0),
    DAY_NO NUMBER(38,0),
    HOURS_WORKED NUMBER(38,0),
    EMPLOYEE_NO NUMBER(38,0),
    LAST_NAME VARCHAR(16777216),
    WEEK_NO NUMBER(38,0),
    FIRST_NAME VARCHAR(16777216)
);
```

Run the query which will create the table. Selecting the Databases option from the menu will now display the TBSTAFFHOURS table. Clicking on the table name should display the definition.

The screenshot shows the Snowflake interface. At the top, there are navigation icons for Databases, Shares, Data Marketplace, Warehouses, Worksheets, and History, along with a Preview button. Below this, the path 'Databases > RETAIL\_POC > TBSTAFFHOURS (PUBLIC)' is displayed. A sub-menu bar includes Tables, Views, Schemas, Stages, File Formats, Sequences, and Pipes, with 'Tables' being the active tab. A 'Load Table' button is visible. The main area displays the table schema in a grid:

Column Name	Ordinal ▲	Type	Nullable	Default	Comment
STORE_CODE	1	NUMBER(38,0)	true	NULL	
DAY_NO	2	NUMBER(38,0)	true	NULL	
HOURS_WORKED	3	NUMBER(38,0)	true	NULL	
EMPLOYEE_NO	4	NUMBER(38,0)	true	NULL	
LAST_NAME	5	VARCHAR(16777216)	true	NULL	
WEEK_NO	6	NUMBER(38,0)	true	NULL	
FIRST_NAME	7	VARCHAR(16777216)	true	NULL	

Snowflake is now ready to receive the test data. Since we will be accessing the table directly using our Snowflake account there is no need to worry about permissions. In the real world it would be necessary to set up the correct access.

## 4. Create AWS S3 bucket

An S3 bucket is created to store the JSON files simulating employee hours worked sent from an external application. Since access will be controlled by a key pair, public access is not required and the setup will be relatively simple, predominately using default values for security.

Within AWS navigate to the S3 console and click **Create Bucket**. Give the bucket a unique name and everything else can be left as default. My bucket name is emeraldmill.sales but any unique name is fine. Press the **Create Bucket** button once complete.

The screenshot shows the 'Create bucket' wizard in the AWS S3 console. The top navigation bar includes the AWS logo, Services dropdown, search bar ('Search for services, features, marketplace products, and docs'), and a keyboard shortcut '[Alt+S]'. The current page is 'Amazon S3 > Create bucket'. The 'Create bucket' section has a sub-header 'Buckets are containers for data stored in S3. Learn more' with a link icon. The 'General configuration' section contains fields for 'Bucket name' (set to 'emeraldmill.sales') and 'AWS Region' (set to 'EU (London) eu-west-2'). Below these, there is an optional 'Copy settings from existing bucket - optional' section with a 'Choose bucket' button. The final section is 'Block Public Access settings for this bucket', which contains a note about granting public access through ACLs, bucket policies, or access point policies, and a link to 'See rules for bucket naming'.

The bucket will be created and will appear in your list of buckets.

Name	AWS Region
emeraldmill.sales	EU (London) eu-west-2
emeraldmilldestinationbucket	EU (London) eu-west-2
emeraldmilllockingbucket	EU (London) eu-west-2

Clicking on the name will take you to the property pages and display the list of objects currently stored in the bucket. Initially this will be empty, but the illustration below shows the bucket after some files have been loaded.

Name	Type	Last modified	Size	Storage class
fldr1/	Folder	-	-	-
RetailPOC_04052021_124249	-	May 4, 2021, 13:42:50 (UTC+01:00)	845.5 KB	Standard
RetailPOC_04052021_182305	-	May 4, 2021, 19:23:06 (UTC+01:00)	845.2 KB	Standard
RetailPOC_04052021_182508	-	May 4, 2021, 19:25:09 (UTC+01:00)	845.5 KB	Standard
RetailPOC_05052021_003458	-	May 5, 2021, 01:34:59 (UTC+01:00)	845.4 KB	Standard

## 5. Generate AWS Keys

Authentication for the S3 bucket in this POC will use public key cryptography. A key pair generated within AWS will grant root user access to AWS services. In practice this would be a serious security weakness and a user with the minimum clearance to perform necessary tasks should be created. A key pair against this more restricted user could then be generated limiting access to just necessary areas.

To generate an access key, go to the AWS IAM dashboard and click on the **My access key** link on the right-hand side of the screen.

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

The screenshot shows the AWS IAM dashboard. On the left, a sidebar menu includes 'Identity and Access Management (IAM)' with 'Dashboard' and 'Access management' sections, and 'Access reports' with 'Access analyzer' and 'Archive rules'. The main content area displays 'IAM dashboard' statistics: 'Users: 0', 'Roles: 5', 'User groups: 0', 'Identity providers: 0', and 'Customer managed policies: 1'. It also lists 'Security alerts' and 'Best practices'.

This takes you to the security credentials page. Expand the **Access Keys** section and click on **Create New Access Key**.

### Your Security Credentials

Use this page to manage the credentials for your AWS account. To manage credentials for AWS Identity and Access Management (IAM) users, use the [IAM Console](#).

To learn more about the types of AWS credentials and how they're used, see [AWS Security Credentials](#) in AWS General Reference.

The screenshot shows the 'Your Security Credentials' page. A sidebar on the left lists 'Password', 'Multi-factor authentication (MFA)', and 'Access keys (access key ID and secret access key)'. The 'Access keys' section is expanded, showing a table with one row:

Created	Access Key ID	Last Used	Last Used Region	Last Used Service	Status	Actions
Apr 14th 2021	AKIAYYXALVCPRF4W3MXO	2021-05-05 19:39 UTC+0100	eu-west-2	s3	Active	<a href="#">Make Inactive</a>   <a href="#">Delete</a>

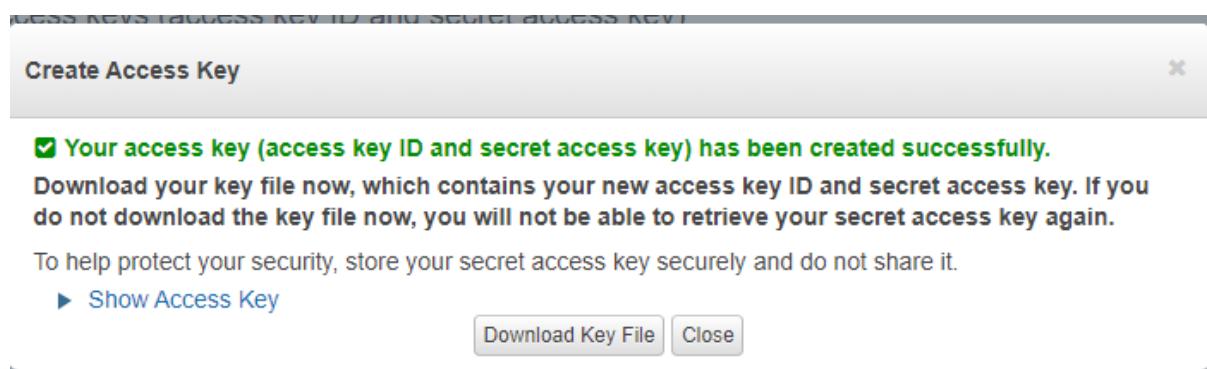
A 'Create New Access Key' button is located below the table. A note in a callout box states: 'Root user access keys provide unrestricted access to your entire AWS account. If you need long-term access keys, we recommend creating a new IAM user with limited permissions and generating access keys for that user instead. [Learn more](#)'.

The sidebar also lists 'CloudFront key pairs' and 'X.509 certificate'.



*Another way of getting to this screen is from the drop-down account menu on the ribbon at the top of all console screens. Select the My Security Credentials option.*

The key pair will be created and can be downloaded in a text file. Keep this file in a secure place so the key values can be used later.



*It is common practice to have several keys created with all but one set de-activated. The active key is then changed on a cyclic basis to form an additional layer of security. This is known as key rotation.*

## 6. Create an EC2 Instance

Two of the Talend jobs will run as standalone jobs on a dedicated server. Creation of timesheet files will run on demand whereas the consumption of these files is a process that loops continuously once initiated. The server will be an AWS EC2 instance and the only pre-requisite will be installation of Java 8. There is no requirement for Talend Studio to be installed as all dependencies are included in the standalone jobs which build in a similar way to a fat jar.

From the AWS EC2 console click on the **Launch Instance** button to create a new virtual server.

The first stage is to choose the Amazon Machine Image (AMI) on which the new server will be based. Click the Free tier only option on the left-hand pane and scroll down to the **Microsoft Windows Server 2019 Base** option and select.

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

**Step 1: Choose an Amazon Machine Image (AMI)**

1. Choose AMI   2. Choose Instance Type   3. Configure Instance   4. Add Storage   5. Add Tags   6. Configure Security Group   7. Review   Cancel and Exit

Free tier only ⓘ

	<b>Red Hat Enterprise Linux 8 (HVM), SSD Volume Type</b> - ami-06178cf087598769c (64-bit x86) / ami-025e95bc52b79028e (64-bit Arm)	<input type="button" value="Select"/>	<input checked="" type="radio"/> 64-bit (x86) <input type="radio"/> 64-bit (Arm)
	<b>SUSE Linux Enterprise Server 15 SP2 (HVM), SSD Volume Type</b> - ami-0d7db5fc4b5075b0d (64-bit x86) / ami-0fd4500e38324e55 (64-bit Arm)	<input type="button" value="Select"/>	<input checked="" type="radio"/> 64-bit (x86) <input type="radio"/> 64-bit (Arm)
	<b>Ubuntu Server 20.04 LTS (HVM), SSD Volume Type</b> - ami-0194c3e07668a7e36 (64-bit x86) / ami-0960f1036d6edacf5 (64-bit Arm)	<input type="button" value="Select"/>	<input checked="" type="radio"/> 64-bit (x86) <input type="radio"/> 64-bit (Arm)
	<b>Microsoft Windows Server 2019 Base</b> - ami-0ae15c1544cd06ac8	<input type="button" value="Select"/>	<input checked="" type="radio"/> 64-bit (x86)

The instance type for Free Tier eligibility will already be selected so click **Review and Launch** to initiate the instance.

**Step 2: Choose an Instance Type**

Amazon EC2 provides a wide selection of instance types optimized to fit different use cases. Instances are virtual servers that can run applications. They have varying combinations of CPU, memory, storage, and networking capacity, and give you the flexibility to choose the appropriate mix of resources for your applications. [Learn more](#) about instance types and how they can meet your computing needs.

Filter by: All instance families ▾ Current generation ▾ Show/Hide Columns

Currently selected: t2.micro (- ECUs, 1 vCPUs, 2.5 GHz, -, 1 GiB memory, EBS only)

	Family	Type	vCPUs ⓘ	Memory (GiB)	Instance Storage (GB) ⓘ	EBS-Optimized Available ⓘ	Network Performance ⓘ	IPv6 Support ⓘ
<input type="checkbox"/>	t2	t2.nano	1	0.5	EBS only	-	Low to Moderate	Yes
<input checked="" type="checkbox"/>	t2	<b>t2.micro</b> <span style="background-color: #00AEEF; color: white; padding: 2px 5px;">Free tier eligible</span>	1	1	EBS only	-	Low to Moderate	Yes
<input type="checkbox"/>	t2	t2.small	1	2	EBS only	-	Low to Moderate	Yes
<input type="checkbox"/>	t2	t2.medium	2	4	EBS only	-	Low to Moderate	Yes
<input type="checkbox"/>	t2	t2.large	2	8	EBS only	-	Low to Moderate	Yes
<input type="checkbox"/>	t2	t2.xlarge	4	16	EBS only	-	Moderate	Yes
<input type="checkbox"/>	t2	t2.2xlarge	8	32	EBS only	-	Moderate	Yes
<input type="checkbox"/>	t3	t3.nano	2	0.5	EBS only	Yes	Up to 5 Gigabit	Yes

Cancel Previous  Next: Configure Instance Details

The details of the instance will be displayed. For this proof of concept, the defaults are all fine so click **Launch** to start the server.

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

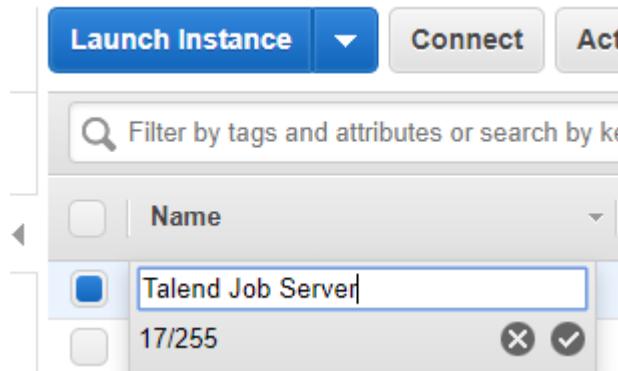
The screenshot shows the AWS Launch Wizard Step 7: Review Instance Launch. At the top, there's a navigation bar with the AWS logo, services dropdown, search bar, account information (John Ts AWS Account, London, Support), and a progress bar from 1. Choose AMI to 7. Review. The main content area is titled "Step 7: Review Instance Launch". It includes a note: "Please review your instance launch details. You can go back to edit changes for each section. Click **Launch** to assign a key pair to your instance and complete the launch process." A warning box says: "⚠ Improve your instances' security. Your security group, launch-wizard-2, is open to the world. Your instances may be accessible from any IP address. We recommend that you update your security group rules to allow access from known IP addresses only. You can also open additional ports in your security group to facilitate access to the application or service you're running, e.g., HTTP (80) for web servers. [Edit security groups](#)". Below this are sections for AMI Details (Microsoft Windows Server 2019 Base with Containers - ami-0a0b4d8b0e9aef6bc, Free tier eligible), Instance Type (t2.micro), and Security Groups (Security group name: launch-wizard-2, Description: launch-wizard-2 created 2021-05-07T12:48:34.329+01:00). At the bottom right are "Cancel", "Previous", and a large blue "Launch" button.

To launch the server, you will be asked to either use an existing key pair or you can create new ones. These are different from the key pairs generated earlier and relate specifically to EC2 instances. Using these allows you to access the instance via remote desktop (RDP). When a key pair file is generated keep it in a known place so it can be accessed when required.

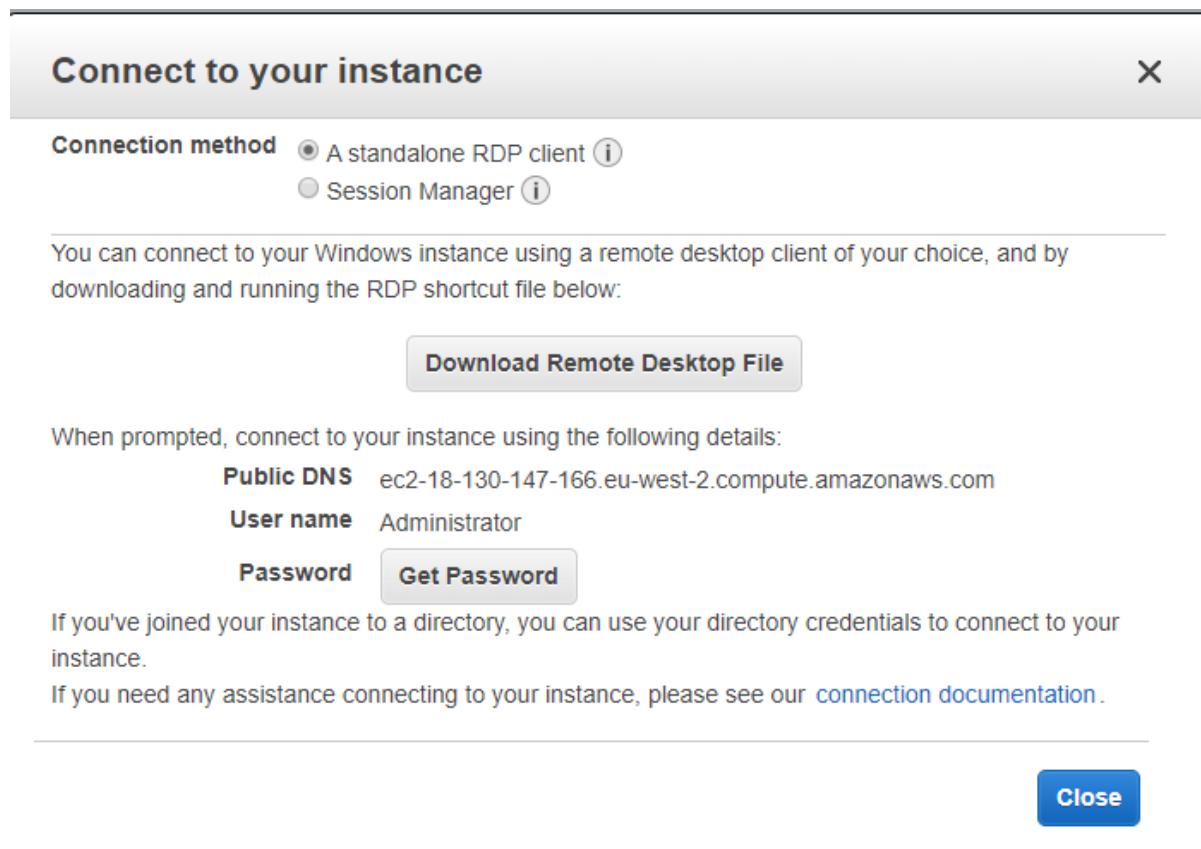
A modal dialog box titled "Select an existing key pair or create a new key pair". It contains text explaining that a key pair consists of a public key (AWS stores) and a private key file (you store) for secure connection. For Windows AMIs, the private key file is required for RDP. For Linux AMIs, it's for SSH. A note states: "Note: The selected key pair will be added to the set of keys authorized for this instance. Learn more about [removing existing key pairs from a public AMI](#) ."

Below this are two dropdown menus: "Choose an existing key pair" (set to "Select a key pair") and "EC2KeyPair". There's also a checkbox: "I acknowledge that I have access to the selected private key file (EC2KeyPair.pem), and that without this file, I won't be able to log into my instance." At the bottom are "Cancel" and a large blue "Launch Instances" button.

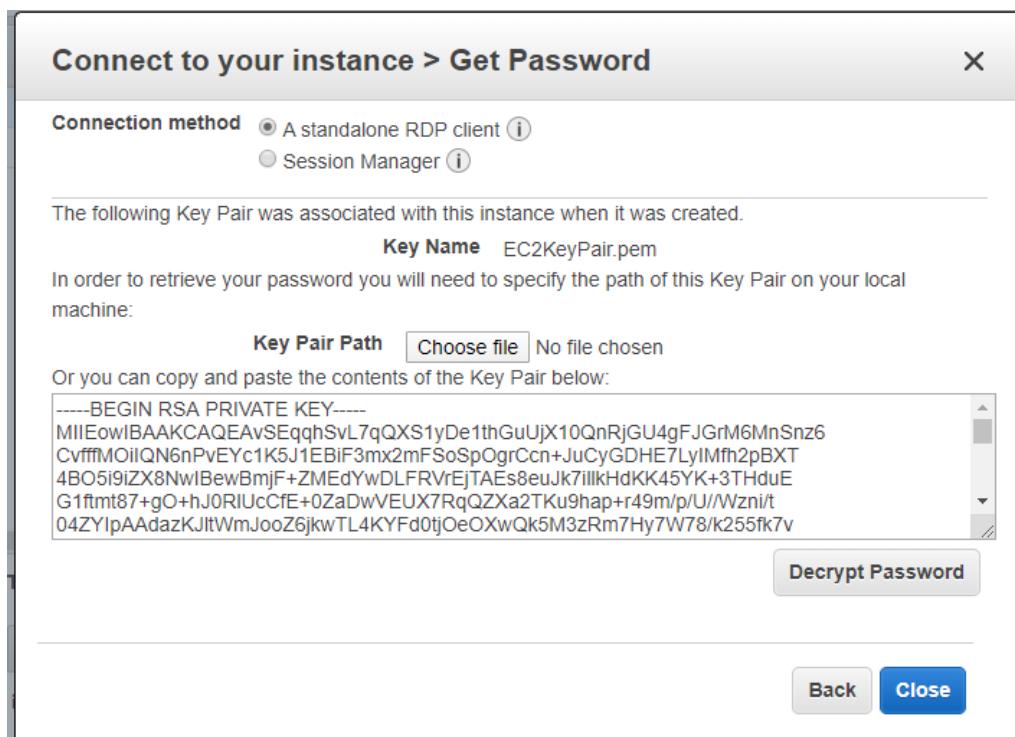
Check the acknowledgement box and launch the instance. The instance will now start and be displayed in the Instances list. The next thing to do is to give it a name and the field can be edited directly in the **Name** field.



To connect to the server from your local machine click on the **Connect** button with the server selected in the list. A dialog box will present 2 options: **Download Remote Desktop File** and **Get Password**.



Download the RDP file first and save in a suitable location then click on **Get Password**. Specify the key pair by clicking **Choose file**.



Click on **Decrypt Password** and you will return to the previous screen but with the password now visible. Store this password in a safe place as it will be needed any time you RDP onto this server.

Looking at the RDP file generated for the instance, full address and username are specified.

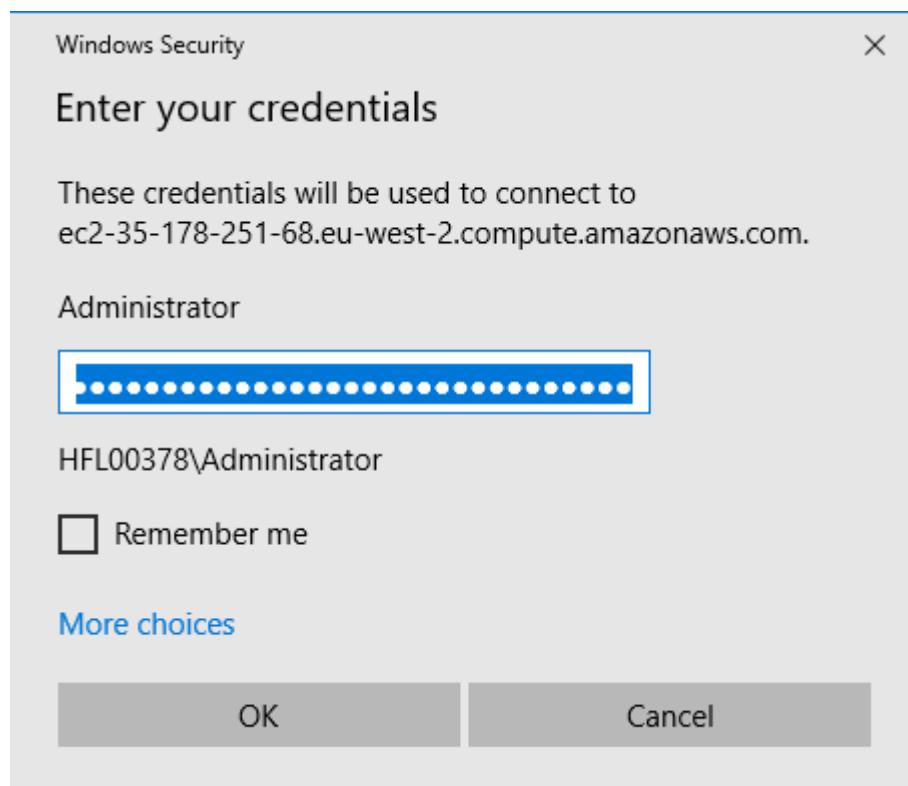
```
1 auto connect:i:1
2 full address:s:ec2-35-178-251-68.eu-west-2.compute.amazonaws.com
3 username:s:Administrator
```

length : 108 lines Ln : 5 Col : 1 Pos : 109 Unix (LF) UTF-8 INS

The address matches the public IPV4 DNS shown in the AWS console.

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status	Alarm Status	Public DNS (IPv4)
Talend Job Server	i-0fcccf0ca2aa73e20	t2.micro	eu-west-2a	running	2...	None	ec2-35-178-251-68.eu-west-2.compute.amazonaws.com

To connect to the server, double-click the RDP file to call the desktop client. Enter the password created and saved previously then press **OK**.



A remote desktop session to the server should now be established and any configurations can be made in the same way as if it was a local pc. The EC2 server is now ready to have Java installed.





*The public IPV4 address only applies to the running instance. If it is stopped and restarted later, a common practice to save resource, the new running instance will have a different address. The old RDP file will no longer create a connection so either create a new RDP file every time an instance is restarted or it is much simpler to just update the DNS in the existing file with the new one from the AWS console. In practice only the four octal byte values will change, the region and availability zone domains remain constant so it may be easier to just change the numbers in the address. The password remains constant for the life of the instance and is not affected by restarts.*

## 7. Install Java 8 JRE onto EC2

For the Talend jobs to run standalone on the EC2 instance it is necessary to install the Java runtime environment (JRE). Talend still recommends version 8 rather than 11.

Download the version 8 installer from the Oracle Java site.

Java Downloads for All Operating Systems  
Recommended Version 8 Update 291  
Release date April 20, 2021

**Important Oracle Java License Update**  
The Oracle Java License has changed for releases starting April 16, 2019.  
The new Oracle Technology Network License Agreement for Oracle Java SE is substantially different from prior Oracle Java licenses. The new license permits certain uses, such as personal use and development use, at no cost – but other uses authorized under prior Oracle Java licenses may no longer be available. Please review the terms carefully before downloading and using this product. An FAQ is available [here](#).  
Commercial license and support is available with a low cost [Java SE Subscription](#).  
Oracle also provides the latest OpenJDK release under the open source [GPL License](#) at [jdk.java.net](#).

Select the file according to your operating system from the list below to get the latest Java for your computer.

> [All Java Downloads](#)      > [Remove Older Versions](#)      > [What is Java?](#)

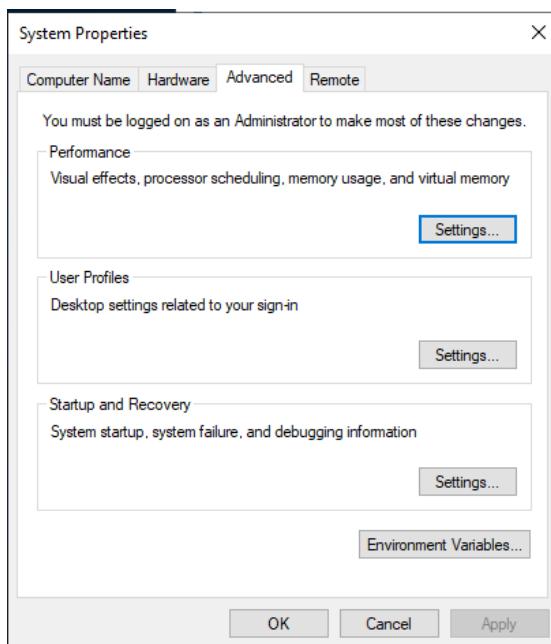
By downloading Java you acknowledge that you have read and accepted the terms of the [Oracle Technology Network License Agreement for Oracle Java SE](#)

Windows			
Which should I choose?			
	<a href="#">Windows Online</a> filesize: 1.98 MB	<a href="#">Instructions</a>	After installing Java, you may need to restart your
	<a href="#">Windows Offline</a>		

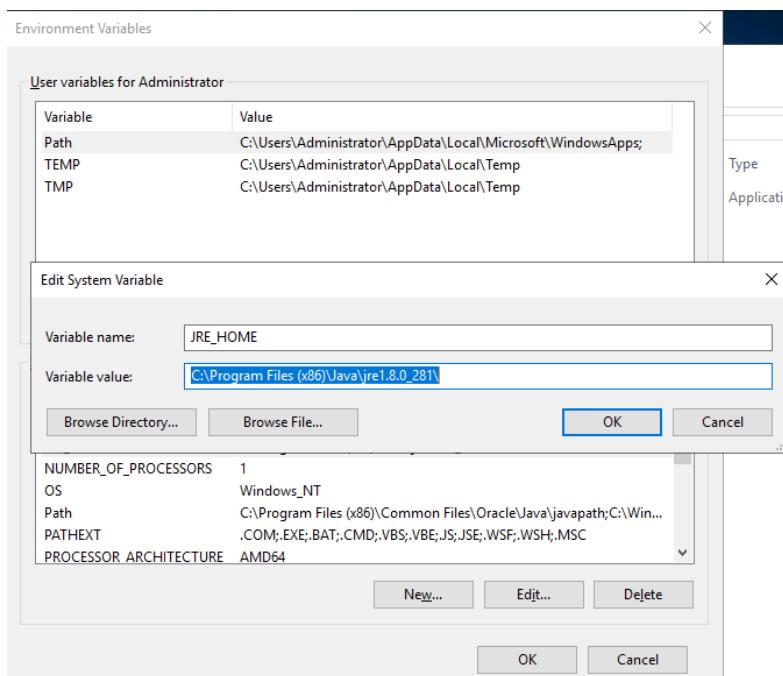
Click on the file to run the installer package and accept the defaults to install Java. This is a straightforward process but in the event of problems there are numerous on-line resources dedicated to the operation, so I won't repeat it here.

Once the installation is complete the next step is to modify the environment variables on the server to allow Talend to run communicate with the JRE correctly.

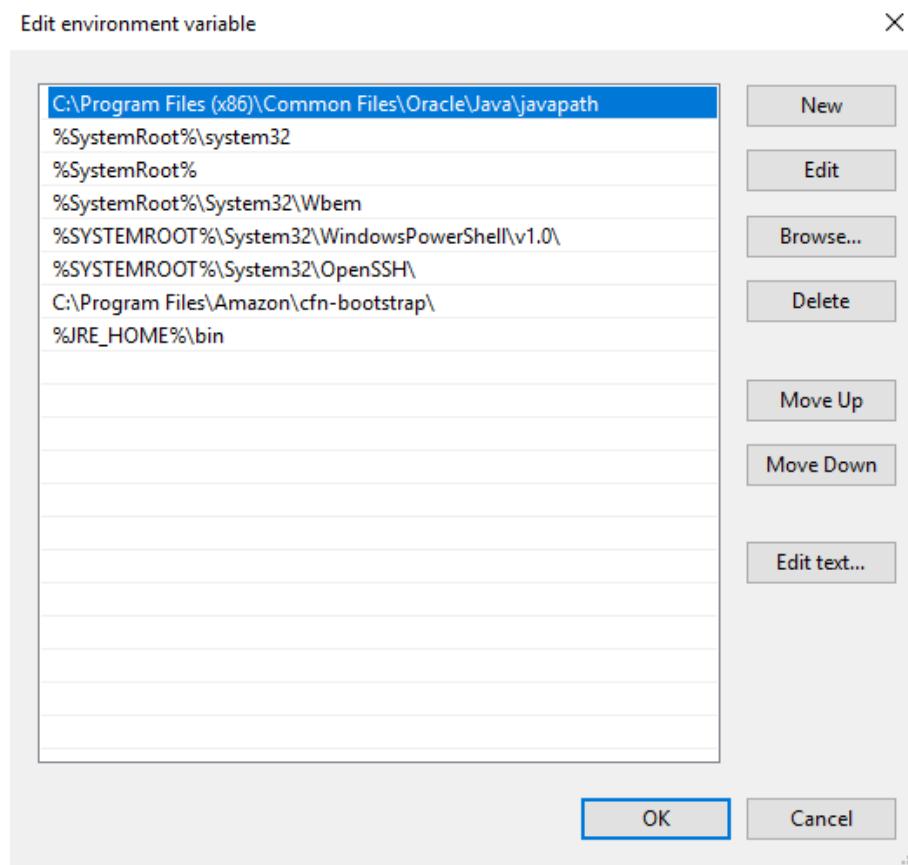
Working in the RDP session for the server, from the **System Properties** dialog click on **Environment Variables**.



In the Environment Variables Dialog click **New** and add a value for **JRE\_HOME** which points to the java location. This should be **C:\Program Files (x86)\Java\jre1.8.0\_281\** if the standard installation defaults were used.



Save the variable by pressing **OK** then click on the **PATH** variable in the System variables section. Add a reference to the **JRE\_HOME** value to the path.



The Java environment is now setup for use by Talend. Notice that only the JRE is required for runtime use, not the JDK. If you wished to install Talend Open Studio (TOS) on the server then the JDK would be required as well.

## 8. Create S3 role for EC2

By default, an EC2 instance has no access to S3 so any jobs running on that server the need to access a bucket will fail. To fix this a role implementing an S3 access policy can be assigned to the server. In normal operations this role should be restricted to the minimum required access but for the proof of concept a simple generalised policy will be used.

Working in the AWS IAM console select **Roles** from the right-hand pane and click **Create role**

Role name	Trusted entities	Last activity
AWSServiceRoleForRDS	AWS service: rds (Service-Linked role)	28 days
AWSServiceRoleForSupport	AWS service: support (Service-Linked role)	None
AWSServiceRoleForTrustedAdvisor	AWS service: trustedadvisor (Service-Linked ...)	None
myLambda-role-fioxmr7b	AWS service: lambda	3 days

In the Create role screen select AWS Service then EC2 from common use cases. Press **Next:Permissions** to proceed.

Create role

1 2 3 4

Select type of trusted entity

**AWS service** EC2, Lambda and others      **Another AWS account** Belonging to you or 3rd party      **Web identity** Cognito or any OpenID provider      **SAML 2.0 federation** Your corporate directory

Allows AWS services to perform actions on your behalf. [Learn more](#)

Choose a use case

**Common use cases**

**EC2**  
Allows EC2 instances to call AWS services on your behalf.

**Lambda**  
Allows Lambda functions to call AWS services on your behalf.

Or select a service to view its use cases

API Gateway	CodeBuild	EMR	IoT SiteWise	RDS
AWS Backup	CodeDeploy	EMR Containers	IoT Things Graph	Redshift
AWS Chatbot	CodeGuru	ElastiCache	KMS	Rekognition
AWS Marketplace	CodeStar Notifications	Elastic Beanstalk	Kinesis	RoboMaker
AWS Support	Comprehend	Elastic Container Registry	Lake Formation	S3

\* Required      Cancel      **Next: Permissions**

Type **S3** in the Filter policies box then check the box next to AmazonS3FullAccess and click **Next:Tags**

▼ Attach permissions policies

Choose one or more policies to attach to your new role.

**Create policy**

**Filter policies**  Showing 8 results

	Policy name	Used as
<input type="checkbox"/>	AmazonDMSRedshiftS3Role	None
<input type="checkbox"/>	AmazonS3FullAccess	Permissions policy (1)
<input type="checkbox"/>	AmazonS3OutpostsFullAccess	None
<input type="checkbox"/>	AmazonS3OutpostsReadOnlyAccess	None
<input type="checkbox"/>	AmazonS3ReadOnlyAccess	None
<input type="checkbox"/>	IVSRecordToS3	None
<input type="checkbox"/>	QuickSightAccessForS3StorageManagementAnalyticsReadOnly	None
<input type="checkbox"/>	S3StorageLensServiceRolePolicy	None

▶ Set permissions boundary

\* Required      Cancel      Previous      **Next: Tags**

Skip the tags screen and Review the role. Check the policy is shown then give it the name S3\_Access and click **Create**.

## Create role

1 2 3 4

### Review

Provide the required information below and review this role before you create it.

**Role name\*** S3\_Access

Use alphanumeric and '+,-,@-' characters. Maximum 64 characters.

**Role description**

Allows EC2 instances to call AWS services on your behalf.

Maximum 1000 characters. Use alphanumeric and '+,-,@-' characters.

**Trusted entities** AWS service: ec2.amazonaws.com

**Policies**  AmazonS3FullAccess

**Permissions boundary** Permissions boundary is not set

No tags were added.

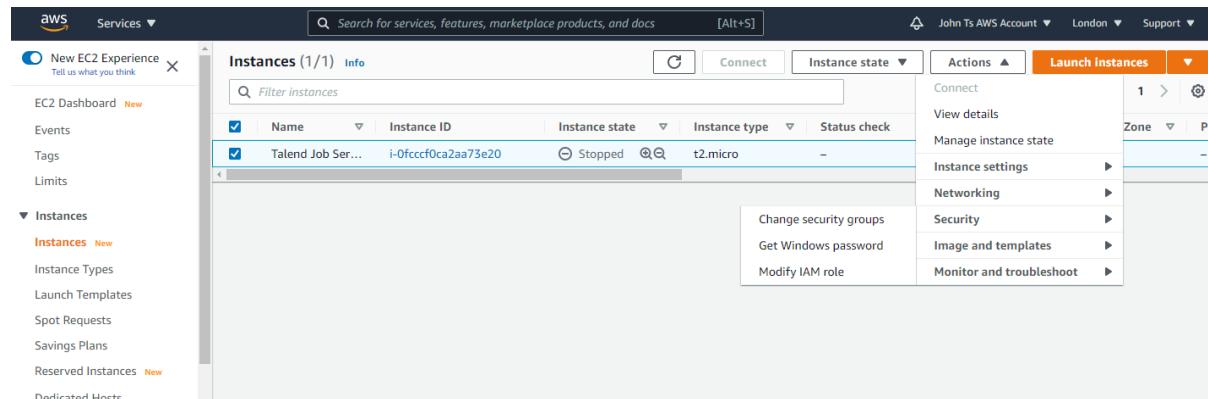
\* Required

Cancel

Previous

Create role

To attach the role to the EC2 instance, open the AWS EC2 console and select **Instances** from the right-hand pane. Select the instance by clicking the check box in the right-hand column, then click **Actions, Security and Modify IAM role**.



Name	Instance ID	Instance state	Instance type	Status check
Talend Job Ser...	i-0fccc0ca2aa73e20	Stopped	t2.micro	-

In the Modify IAM role screen, select the S3\_Access role as defined previously and click **Save** to assign the role.



*Note only 1 role can be applied to an instance and if other access requirements are necessary, their policies would need to be added to the role.*

The EC2 instance is now ready for use. Remember that it can be stopped when not in use and restarted as needed to save resource, but it will restart with a new ipv4 address, so the RDP file will have to be changed each time.

## 9. Create AWS SQS queue

A message queue is a useful way to transfer data between programs. Within Talend, ActiveMQ is often used as it is built into the product, however there are many other queue brokers and AWS has the Simple Queue Service or SQS.

Two types of queue exist: Standard where message ordering is not preserved and FIFO which guarantees first in first out delivery. We will create a FIFO queue for the proof of concept.

From the AWS management console select **Simple Queue Service** and click **Create Queue**. Check the FIFO radio button and name the queue RetailPOC.fifo. The rest of the config can be left as default for this scenario then click **Create Queue**.

Amazon SQS > Queues > Create queue

### Create queue

**Details**

Type  
Choose the queue type for your application or cloud infrastructure.

ⓘ You can't change the queue type after you create a queue.

Standard Info

At-least-once delivery, message ordering isn't preserved

- At-least once delivery
- Best-effort ordering

FIFO Info

First-in-first-out delivery, message ordering is preserved

- First-in-first-out delivery
- Exactly-once processing

Name  
  
A queue name is case-sensitive and can have up to 80 characters. You can use alphanumeric characters, hyphens (-), and underscores (\_).

**Configuration**  
Set the maximum message size, visibility to other consumers, and message retention. Info

After creation the queue will appear in the list which was also show any current messages available.

Amazon SQS > Queues

Queues (1)						
<span style="border: 1px solid #0072bc; padding: 2px;">C</span> <span>Edit</span> <span>Delete</span> <span>Send and receive messages</span> <span>Actions</span> <span style="background-color: #0072bc; color: white; border: 1px solid #0072bc; padding: 2px;">Create queue</span>						
<span>&lt;</span> <span>1</span> <span>&gt;</span> <span>②</span>						
	<input type="radio"/>	<a href="#">RetailPOC fifo</a>	FIFO	16/04/2021, 17:15:26 BST	0	0

Clicking the queue name will take you to the detail screen which contains full information including monitoring, plus the ability to send and receive messages to test the queue.

Amazon SQS > Queues > RetailPOC fifo

### RetailPOC fifo

Edit Delete Purge Send and receive messages

**Details** Info

Name	Type	ARN
<a href="#">RetailPOC fifo</a>	FIFO	<a href="#">arn:aws:sqs:eu-west-2:602839689375:RetailPOC fifo</a>
Encryption	URL	Dead-letter queue
-	<a href="https://sq.eu-west-2.amazonaws.com/602839689375/RetailPOC fifo">https://sq.eu-west-2.amazonaws.com/602839689375/RetailPOC fifo</a>	-

▶ More

Amazon SQS assigns a unique identifier called a queue URL to each new queue. The queue URL includes the owner account ID, the queue name, and the queue region. You provide the queue URL when you perform any action on a queue.

The name of a FIFO queue must end with the .fifo suffix. The suffix counts towards the 80-character queue name quota.

We recommend that you delete any queue that you are not using (and don't foresee using it in the near future).

You can delete a queue even when it isn't empty. If you want to delete the messages in a queue but not the queue itself, you can purge the queue.

When you purge a queue, the message deletion process takes up to 60 seconds. We recommend waiting for 60 seconds regardless of your queue's size.

To test the queue, press the **Send and receive messages** option and in the Send Message section add a test message and message group id and click **Send message**.



*As this is a FIFO queue it must also have a message group id. FIFO queue logic applies only per message group ID and all messages are sent and received in strict order. Enforcing the group id increases flexibility of the queue allowing it to be serviced by multiple clients using either the same or different group ids and utilising the FIFO capabilities as required. Standard queues do not use this value as it would have no logical purpose and is not displayed in the queue creation screen.*

Amazon SQS > Queues > RetailPOC fifo > Send and receive messages

Send and receive messages

Send messages to and receive messages from a queue.

**Send message** Info

Message body  
Enter the message to send to the queue.  
This is a test of the RetailPOC fifo queue

Message group ID  
The tag that specifies that a message belongs to a specific message group.  
TestGroup

Message deduplication ID - Optional  
The token used for deduplication of messages within the deduplication interval.  
Enter message deduplication id

► Message attributes - Optional Info

Looking at the Receive message section the message is now shown as available. Click the **Poll for messages** option to retrieve it.

Receive messages Info

Messages available: 1 Polling duration: 30 Maximum message count: 10 Polling progress: 60%  
1 receives/second

Messages (1)

ID	Sent	Size	Receive count
de10994a-8e4f-4a18-ac6a-3838e48c6624	13/05/2021, 10:42:03 BST	42 bytes	1

Clicking on the message ID will show the contents.

Message: de10994a-8e4f-4a18-ac6a-3838e48c6624

Details Body Attributes

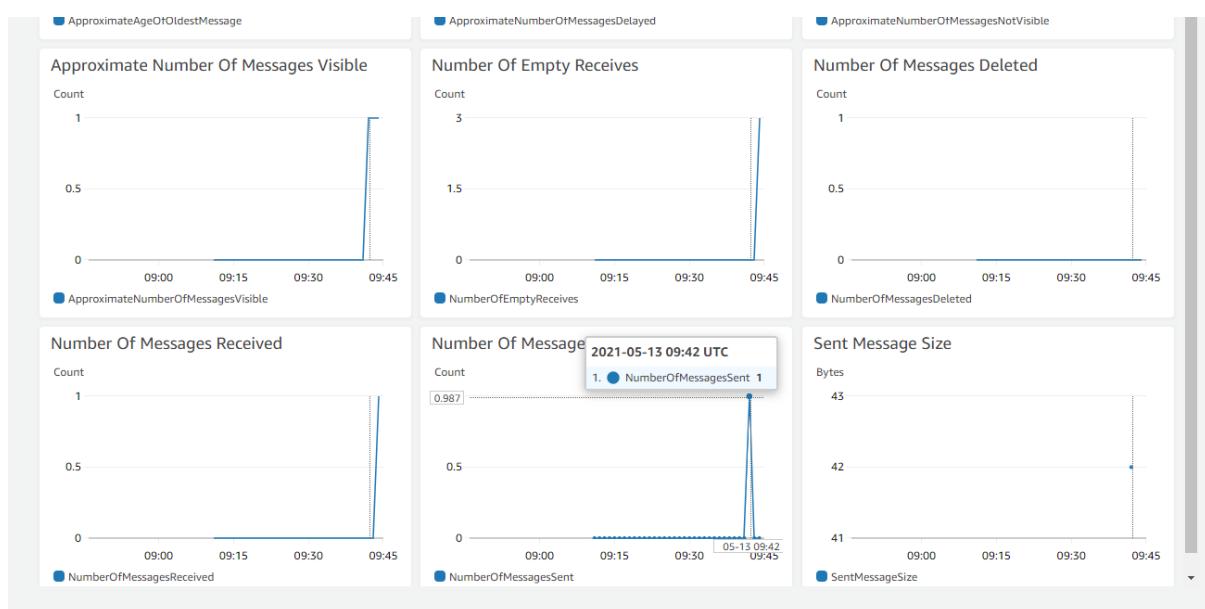
This is a test of the RetailPOC fifo queue

Done

Return to the detail screen and select the **Monitoring** tab. Set the window to 1 hour and the details of the messages sent and received will be displayed.



*This is the default setup, which is fine for the POC, but it can be customised or extended according to requirements.*



The SQS queue is now ready for use.

## 10. Create AWS Lambda function

The design goal of the system is for it to respond automatically to the arrival of a file in the S3 bucket and this is achieved by creating a serverless Lambda function triggered by the S3 arrival, which writes the filename to the body of a message stored on the SQS queue defined in the previous section. A Talend job can then monitor the queue and respond to any new messages.

Note Lambda functions are a component of serverless computing meaning that the function code is passed to the service which itself takes care of hosting and execution without any configuration from the user.

From the AWS console select Lambda then click on **Create function**.

The screenshot shows the AWS Lambda Functions list screen. It displays one function named "myLambda". The interface includes a search bar, filter options, and columns for Function name, Description, Package type, Runtime, Code size, and Last modified.

In the Create function screen select **Author from scratch**, name it myLambda and choose Python 3.7 for the runtime. All other defaults are fine in this case so click **Create function**



*Lambda offers several choices of language including C#, Java, Node.js, Python and Ruby. Some can be coded directly in the Lambda console whilst others such as C# use external tools such as visual studio. For this simple example Python will be used for clarity.*

Lambda > Functions > Create function

### Create function Info

Choose one of the following options to create your function.

- Author from scratch Info  
Start with a simple Hello World example.
- Use a blueprint Info  
Build a Lambda application from sample code and configuration presets for common use cases.
- Container image Info  
Select a container image to deploy for your function.
- Browse serverless app repository Info  
Deploy a sample Lambda application from the AWS Serverless Application Repository.

---

#### Basic information

Function name Info  
Enter a name that describes the purpose of your function.

Use only letters, numbers, hyphens, or underscores with no spaces.

Runtime Info  
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

Permissions Info

In the function overview click the **Add trigger** button

myLambda

▼ Function overview

+ Add trigger

Select **S3** from the list and the bucket name from the drop-down options. In this case it will be emeraldmill.sales. Other defaults can be left alone. Click the disclaimer at the bottom warning about recursive functions and press **Add** to create the trigger.

## Add trigger

### Trigger configuration



S3

aws storage



#### Bucket

Please select the S3 bucket that serves as the event source. The bucket must be in the same region as the function.

emeraldmill.sales



#### Event type

Select the events that you want to have trigger the Lambda function. You can optionally set up a prefix or suffix for an event. However, for each bucket, individual events cannot have multiple configurations with overlapping prefixes or suffixes that could match the same object key.

All object create events



#### Prefix - optional

Enter a single optional prefix to limit the notifications to objects with keys that start with matching characters.

e.g. images/

#### Suffix - optional

Enter a single optional suffix to limit the notifications to objects with keys that end with matching characters.

e.g. jpg

Lambda will add the necessary permissions for Amazon S3 to invoke your Lambda function from this trigger. [Learn more](#)

In the code window copy and paste the following code:

```
#####
Your module description
#####

import boto3
import json

def lambda_handler(event, context):

    sqs = boto3.resource('sqs')
    file_key = event['Records'][0]['s3']['object']['key']

    queue = sqs.get_queue_by_name(QueueName='RetailPOC fifo')

    response = queue.send_message(
        MessageBody=file_key,
        MessageGroupId='messageGroup1'
    )
```

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

Click on **Deploy** once complete to deploy the function in Lambda.

To test the function, select the **test** option from the menu and click **New event** and select **Amazon S3 Put** from the template. Name it S3PutTest and in the “key” element change the value to be “RetailPOCFile”. Save the changes and press **Test** to perform the evaluation.

```
10 "principalId": "EXAMPLE"
11 },
12 "requestParameters": {
13     "sourceIPAddress": "127.0.0.1"
14 },
15 "responseElements": {
16     "x-amz-request-id": "EXAMPLE123456789",
17     "x-amz-id-2": "EXAMPLE123456789abcdef0123456789abcdef"
18 },
19 "s3": {
20     "s3SchemaVersion": "1.0",
21     "configurationId": "testConfigRule",
22     "bucket": {
23         "name": "example-bucket",
24         "ownerIdentity": {
25             "principalId": "EXAMPLE"
26         },
27         "arn": "arn:aws:s3:::example-bucket"
28     },
29     "object": {
30         "key": "RetailPOCFile",
31         "size": 1024,
32         "eTag": "0123456789abcdef0123456789abcdef",
33         "versionId": "20210513T182802Z-0123456789abcdef"
34     }
35 }
```

To check the test has worked go the SQS console and select the **RetailPOC.fifo** queue. Click **Send and receive messages** and scroll to the Receive messages section and **Poll for messages**

Messages available	Polling duration	Maximum message count	Polling progress
1	30	10	67% 1 receives/second

Messages (1)			
Search messages		View details Delete	
ID	Sent	Size	Receive count
f9df0ae6-4af8-4272-aee4-3a4907cdbf1e	13/05/2021, 18:28:02 BST	13 bytes	5

A message should be shown in the queue. Click on the **ID** to view the message and select the body.

Message: f9df0ae6-4af8-4272-aee4-3a4907cdbf1e X

<a href="#">Details</a>	<b>Body</b>	<a href="#">Attributes</a>
<pre>RetailPOCFile</pre>		

**Done**

A message showing the key of the S3 object that we modified in the Lambda test is displayed, indicating that the function has responded correctly to the trigger and created the expected SQS message.



*Using the test mechanism, we have simulated the S3 event and no file has actually been deposited in the bucket however the AWS template code is designed such that the presence of an actual file would respond in exactly the same way.*

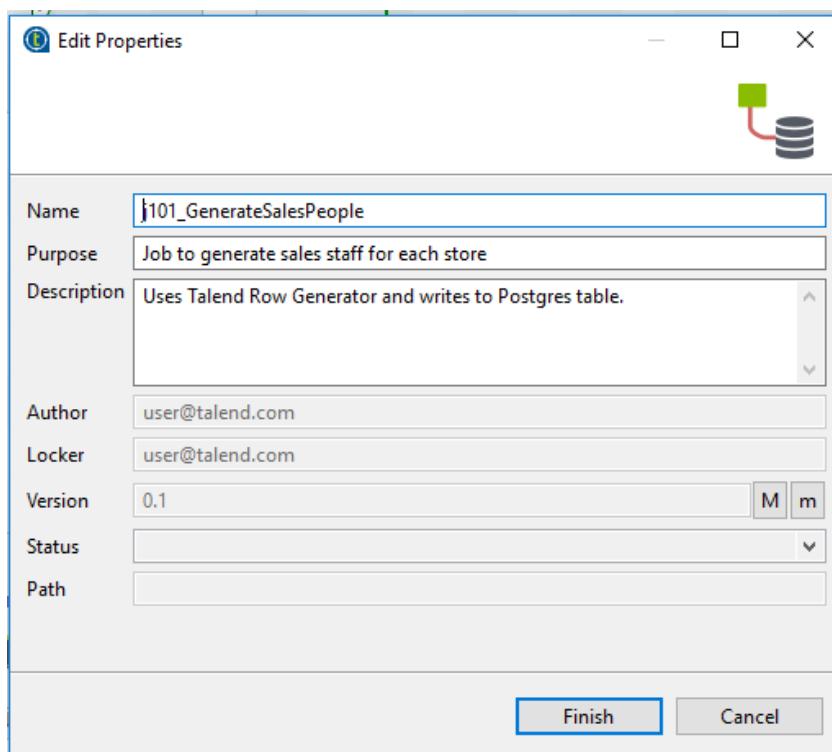
The Lambda function is now complete and will be called anytime a file is loaded into the S3 bucket.

## 11. Create job to generate store staff

The prerequisites are now complete so we will now look at creating the Talend jobs in Open Studio Big Data Edition. There will be three jobs to create the sales staff records, simulate the work hours and provide a time sheet file and finally to process that file and load the data into SnowflakeDB.

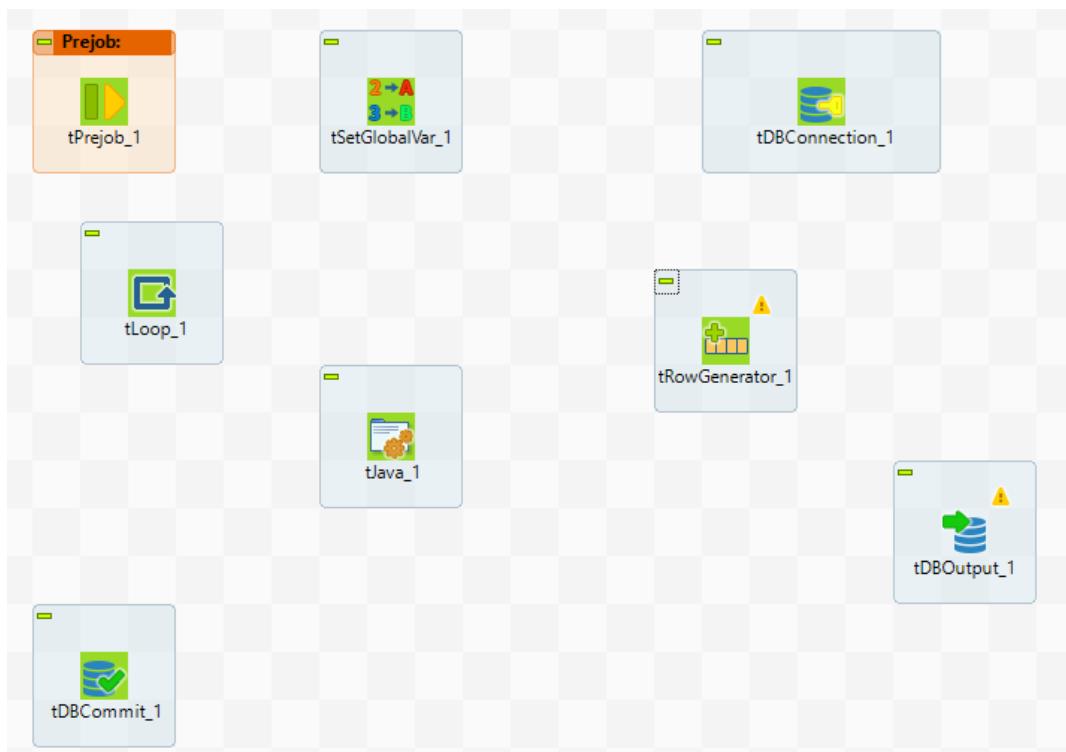
The first job will use the Talend Row generator to simulate sales staff records for each store. The details will be stored in a Postgresql database table.

In the studio repository create a new job and name it j101\_GenerateSalesPeople.



Add the following components to the design area:

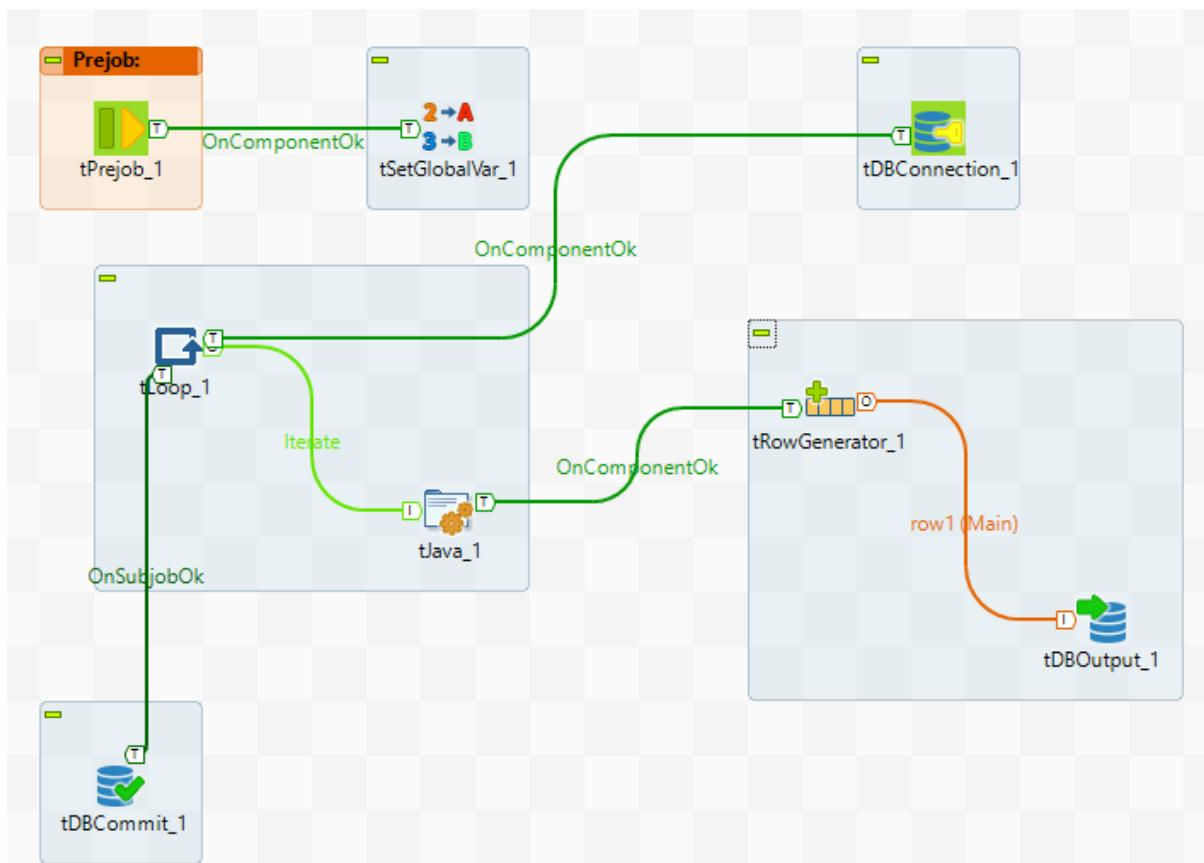
Tprejob	tLoop	tDBOutput
tSetGlobalVar	tJava	tDBCommit
tDBConnection	tRowGenerator	



Join the components as follows:

From Component	To Component	Join Type
<b>tPrejob</b>	tSetGlobalVar	OnComponentOk trigger
<b>tDBConnection</b>	tLoop	OnComponentOk trigger
<b>tLoop</b>	tJava	Iterate
<b>tLoop</b>	tDBCommit	OnSubjobOk trigger
<b>tJava</b>	tRowGenerator	OnComponentOk trigger
<b>tRowGenerator</b>	tDBOutput	Row Main

The job should now look like the diagram below:



Configure the components as follows:

**tPrejob** requires no configuration and just initiates a task each time the job starts and is guaranteed to run prior to the main job.

**tSetGlobalVar** initialises the variable for holding the current store code in the global cache. Click the component in the design area and select the **Component** tab in the bottom pane to show the editor.

Key	Value
"vStoreCode"	0

From the **Basic Settings** option, press the **green + icon** to create a new row and name it “vStoreCode”. Assign a default value of 0 then select the **View** option from the menu.

Add a user-friendly label to component by adding the expression “**<b>Initialise Globals</b><br>**” to the start of the Label format field which will show the label in bold together with the actual component name underneath.

**tDBConnection** is used to connect to the Postgresql instance which is running in AWS RDS as configured previously.

Select **PostgreSQL** as the database with version v9 or later. The Host value will be the URL shown in the AWS RDS console for the instance. Port should be the standard value 5432 unless a different setup has been used. The Database name will be RetailPOC and the public schema is used. Finally, Username and Password refer to the local user set up in RDS for the instance. A shared connection is not used in this case and since the job will just be ran in Studio a data source alias is not required either and both check boxes can remain unchecked.

Go into the **view** option and label the component “**AWS RDS Postgres Connection**”.

**tLoop** is used to loop through the 8 stores used in the example and on each iteration run a routine to generate 20 staff records and write them to the database table. On completion of the entire loop the records will be committed, closing the transaction.

**Loop Type:** For  
**From:** 1  
**To:** 8  
**Step:** 1  
 Values are increasing

As we have a defined number of iterations, a For loop should be chosen going from 1 to 8 in increments of 1. Label the loop “**Loop Through Stores**”.

**tJava** is used to assign the iteration value of tLoop to the global variable vStoreCode. Each call will override the previous value, so it is not necessary to re-initialise the variable each time.

**Code:**

```
globalMap.put("vStoreCode",((Integer)globalMap.get("tLoop_1_CURRENT_ITERATION")));
```

Label the component “**Set Global Value**”.

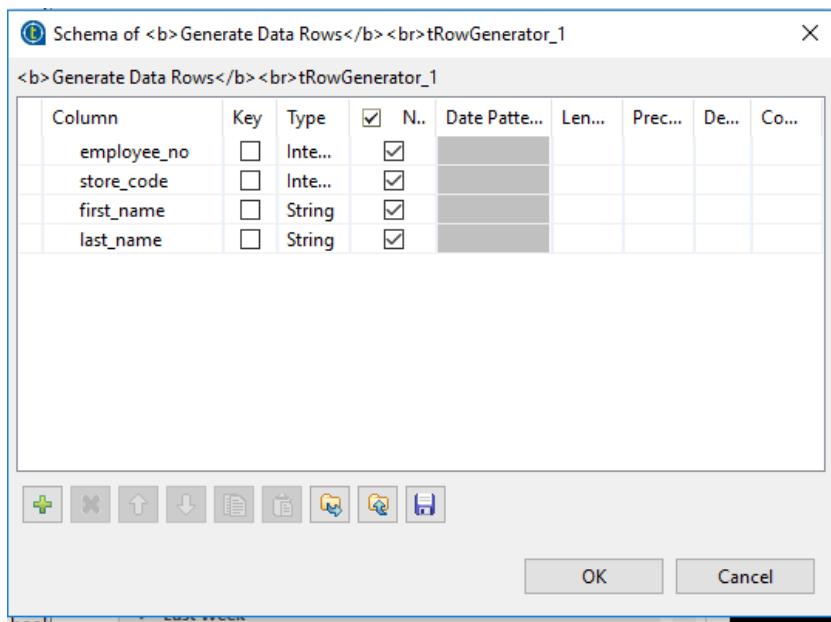
**tRowGenerator** is a Talend component that allows the generation of Random data in a readable format. For the proof of concept, we will generate some skeleton employee data with an employee number, store code, first name and last name. Store code and employee number will be generated as functions of the loop iteration and the names will be assigned from the row generator.



*Note the names will be all male. This is not overt sexism, rather that Talend uses a list of past USA presidents to generate random names and to date these have all been male.*

**Schema:** Built-In  
 Edit schema

Click the **Edit schema** ellipsis to call the editor.



Click the green + icon to add new rows of the following types:

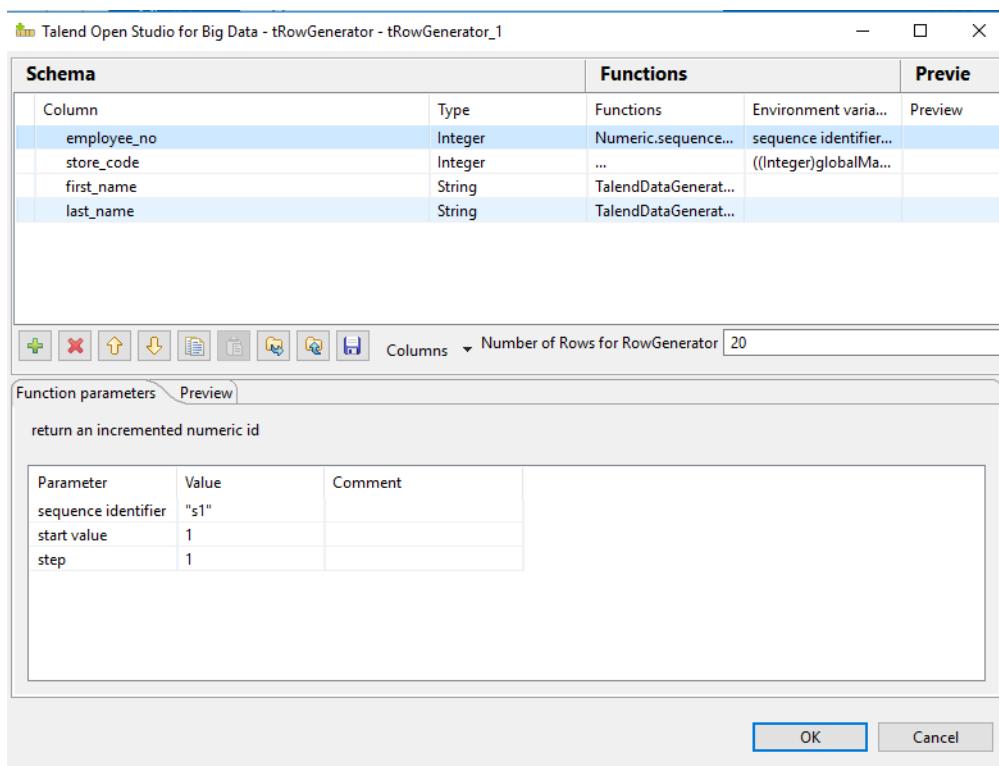
**employee\_no** Integer

**store\_code** Integer

**first\_name** String

**last\_name** String

Press **OK** to save the schema then click on **RowGenerator Editor**.



Insert the value 20 for the **number of rows to be generated**.

The schema fields will already be populated, click on each one in turn and perform the configuration steps.

Field	Actions																																				
<b>employee_no</b>	In the Functions drop down list select the Numeric Sequence option																																				
	<table border="1"> <thead> <tr> <th colspan="2">Functions</th> </tr> </thead> <tbody> <tr> <td>Type</td> <td>Numeric.sequence(String,int,int)</td> </tr> <tr> <td>Integer</td> <td>Mathematical.BITAND(int,int)</td> </tr> <tr> <td>Integer</td> <td>Mathematical.BITNOT(int)</td> </tr> <tr> <td>String</td> <td>Mathematical.BITOR(int,int)</td> </tr> <tr> <td>String</td> <td>Mathematical.BITXOR(int,int)</td> </tr> <tr> <td></td> <td>Mathematical.INT(String)</td> </tr> <tr> <td></td> <td>Mathematical.NUM(String)</td> </tr> <tr> <td></td> <td>Mathematical.SCMP(String,String)</td> </tr> <tr> <td></td> <td>Mathematical.SDIV(int,int)</td> </tr> <tr> <td></td> <td>Numeric.random(int,int)</td> </tr> <tr> <td></td> <td><b>Numeric.sequence(String,int,int)</b></td> </tr> <tr> <td>Columns</td> <td>Number of Rows</td> </tr> <tr> <td></td> <td>StringHandling.COUNT(String,String)</td> </tr> <tr> <td></td> <td>StringHandling.INDEX(String,String)</td> </tr> <tr> <td></td> <td>StringHandling.LEN(String)</td> </tr> <tr> <td></td> <td>TalendDate.compareDate(Date,Date)</td> </tr> </tbody> </table>			Functions		Type	Numeric.sequence(String,int,int)	Integer	Mathematical.BITAND(int,int)	Integer	Mathematical.BITNOT(int)	String	Mathematical.BITOR(int,int)	String	Mathematical.BITXOR(int,int)		Mathematical.INT(String)		Mathematical.NUM(String)		Mathematical.SCMP(String,String)		Mathematical.SDIV(int,int)		Numeric.random(int,int)		<b>Numeric.sequence(String,int,int)</b>	Columns	Number of Rows		StringHandling.COUNT(String,String)		StringHandling.INDEX(String,String)		StringHandling.LEN(String)		TalendDate.compareDate(Date,Date)
Functions																																					
Type	Numeric.sequence(String,int,int)																																				
Integer	Mathematical.BITAND(int,int)																																				
Integer	Mathematical.BITNOT(int)																																				
String	Mathematical.BITOR(int,int)																																				
String	Mathematical.BITXOR(int,int)																																				
	Mathematical.INT(String)																																				
	Mathematical.NUM(String)																																				
	Mathematical.SCMP(String,String)																																				
	Mathematical.SDIV(int,int)																																				
	Numeric.random(int,int)																																				
	<b>Numeric.sequence(String,int,int)</b>																																				
Columns	Number of Rows																																				
	StringHandling.COUNT(String,String)																																				
	StringHandling.INDEX(String,String)																																				
	StringHandling.LEN(String)																																				
	TalendDate.compareDate(Date,Date)																																				
	In the Function parameters tab enter the following values: Sequence identifier "s1" Start value 1 Step 1																																				
<b>store_code</b>	Select the ellipsis option from the top of the drop-down function list which indicates a custom value. In the Function parameters tab enter the value ((Integer)globalMap.get("vStoreCode")) to retrieve the current value of the Store Code global variable.																																				
<b>first_name</b>	Select TalendDataGenerator.getFirstname() from the Functions drop-down.																																				
<b>last_name</b>	Select TalendDataGenerator.getLastname() from the Functions drop-down.																																				

Click **OK** to save the configuration then label the component “**Generate Data Rows**”

**tDBOutput** writes the generated rows to the database table.

Job(j101\_GenerateSalesPeople 0.1) Contexts(j101\_GenerateSalesPeople) Component Run (Job j101\_GenerateSalesPeople)

**<b>Write to Postgres</b> <br>tDBOutput\_1(tDBOutput\_1)(PostgreSQL)**

**Basic settings**

Database: PostgreSQL

Use an existing connection

Component List: tDBConnection\_1 - <b>AWS RDS Postgres Connection</b> <br> tDBConnection\_1 \*

Table: "tbstorestaff"

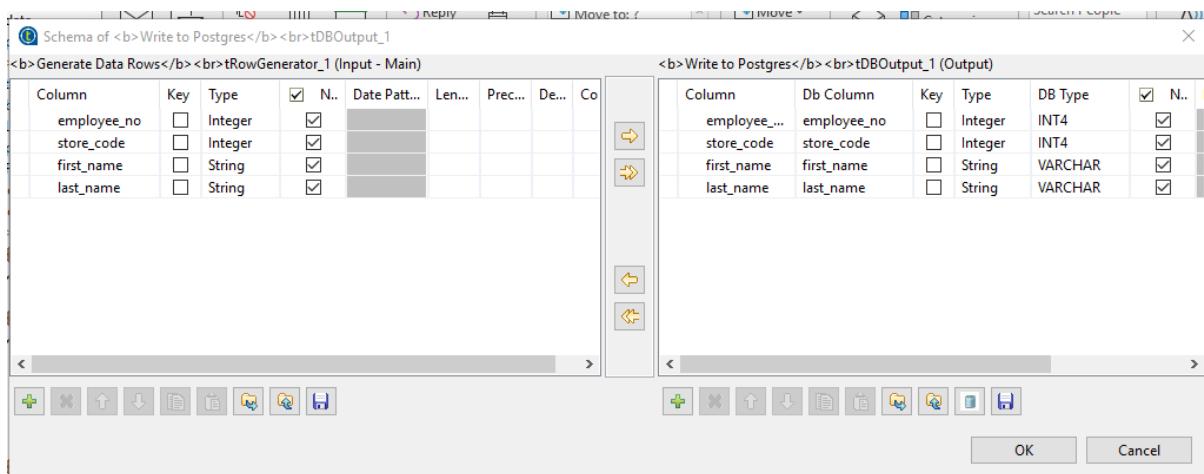
Action on table: Default

Action on data: Insert

Schema: Built-In

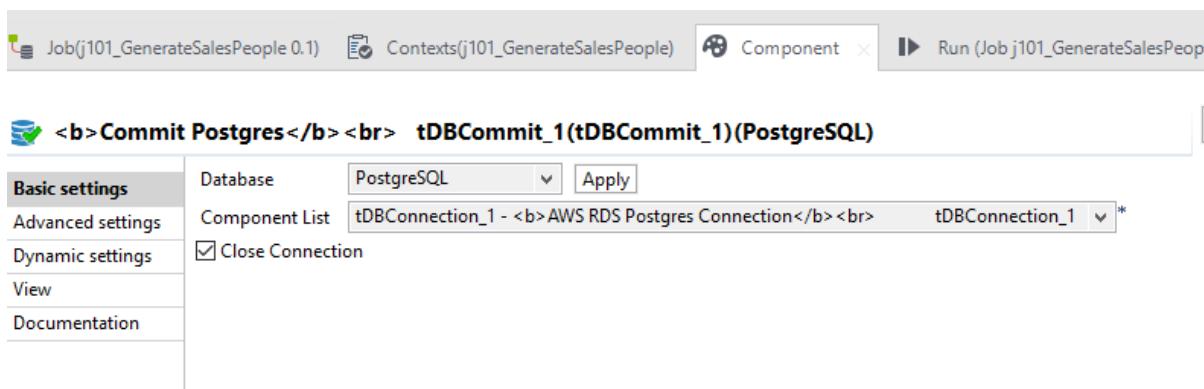
Die on error

Select **PostgreSQL** as the database type and check the **Use an existing connection** box. In the Component List select the DB connection defined previously. Enter “**tbstorestaff**” for the table and leave the Action on table setting as Default. Action on data should be set to **Insert**, then click **Sync columns** followed by **Edit schema** and verify the schema looks the same as below.



Click **OK** to exit the schema editor then label the component “**Write to Postgres**”

**tDBCommit** will finally commit the Postgres transactions and close the database connection.



Select **PostgreSQL** as the database type and choose the database connection component from the list. Check the **Close Connection** option and label the component “**Commit Postgres**”

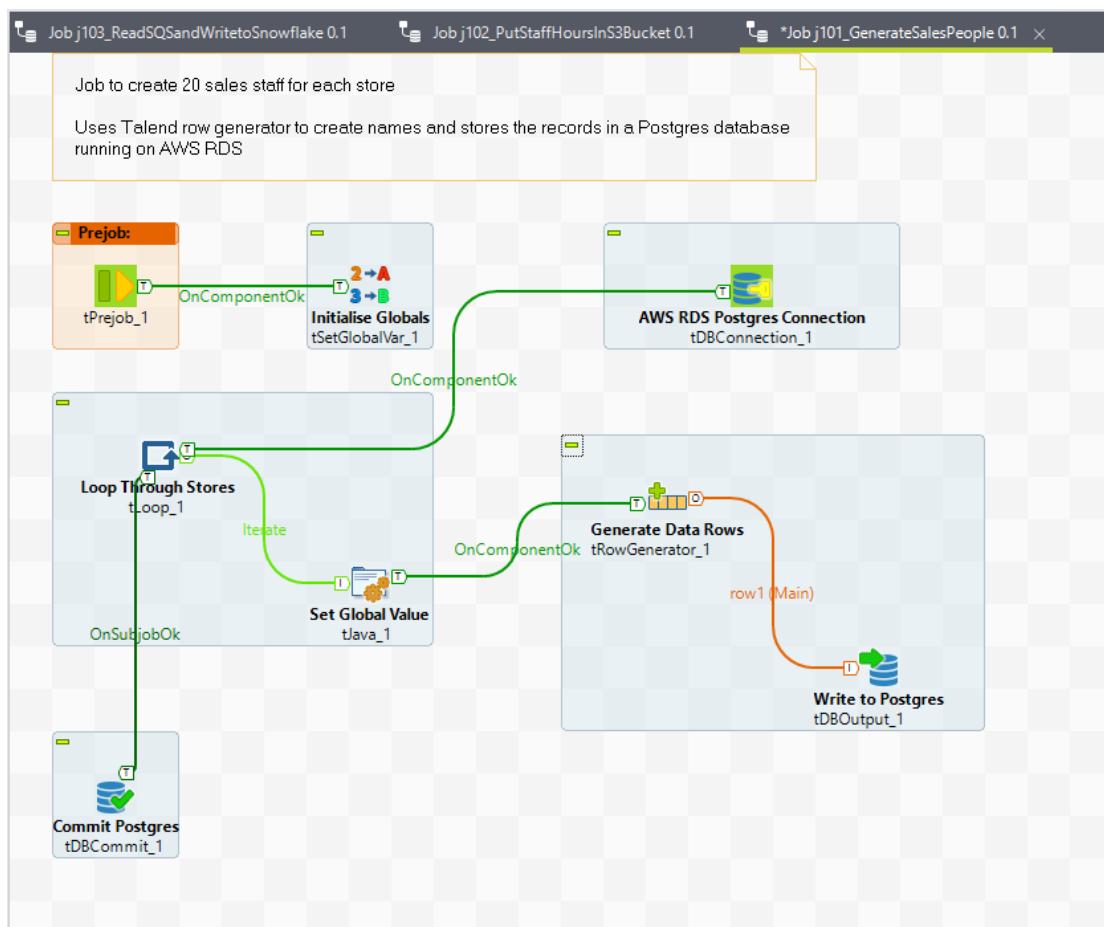


*Ensure that the connecting tLoop to tDBCommit is of the type OnSubjobOk in order that it fires after the loop has completed and all data is written. If an OnComponentOK trigger is inadvertently used instead it will call the commit immediately the loop is started and will close the database connection before data had been written, leading to an error when tDBOutput tries to access it.*

Complete the job by adding a note to the screen similar to:



The completed job should now look like this.



The job is now ready to run in Talend Studio but check that the database table is empty prior to execution or a primary key violation may occur.

After completion the results can be viewed in PGAdmin or any similar client.

The pgAdmin interface shows the following details:

- Servers:** AWSPostgres
- Databases:** RetailPOC
- Tables:** tbstorestaff
- Columns:** employee\_no, store\_code, first\_name, last\_name
- Query Editor:**

```

1 SELECT employee_no, store_code, first_name, last_name
2   FROM public.tbstorestaff;
    
```
- Data Output:**

	employee_no	store_code	first_name	last_name
153	153	8	Calvin	Quincy
154	154	8	Ronald	Harrison
155	155	8	Andrew	Grant
156	156	8	Woodrow	Van Buren
157	157	8	George	Cleveland
158	158	8	Harry	Taft
159	159	8	Warren	Harrison
160	160	8	Rutherford	Reagan

This is a rudimentary routine designed to demonstrate both the row generator and loading data into PostgreSQL hosted on AWS RDS. Flexibility could be enhanced by using context variables loaded via file or database to vary parameters such as number of store staff and total stores and update or insert used to prevent key errors on subsequent runs.

## 12. Create job to generate timesheets

The second job to simulate the creation of timesheets for the staff generated in the previous job and write the data as a JSON file to the S3 bucket defined earlier. This will trigger the AWS Lambda function which writes a message to AWS Simple Queue Service (SQS).

Create a new job in the repository called j102\_PutStaffHoursInS3Bucket.

Drop the following components onto the design area:

tPrejob	tLoop x 3	tFileDelete x 2
tFileInputDelimited x 2	tJava	tS3Connection
tContextLoad	tDBInput	tS3Put
tDBConnection	tFileOutputDelimited	
tSetGlobalVar	tFileOutputJSON	

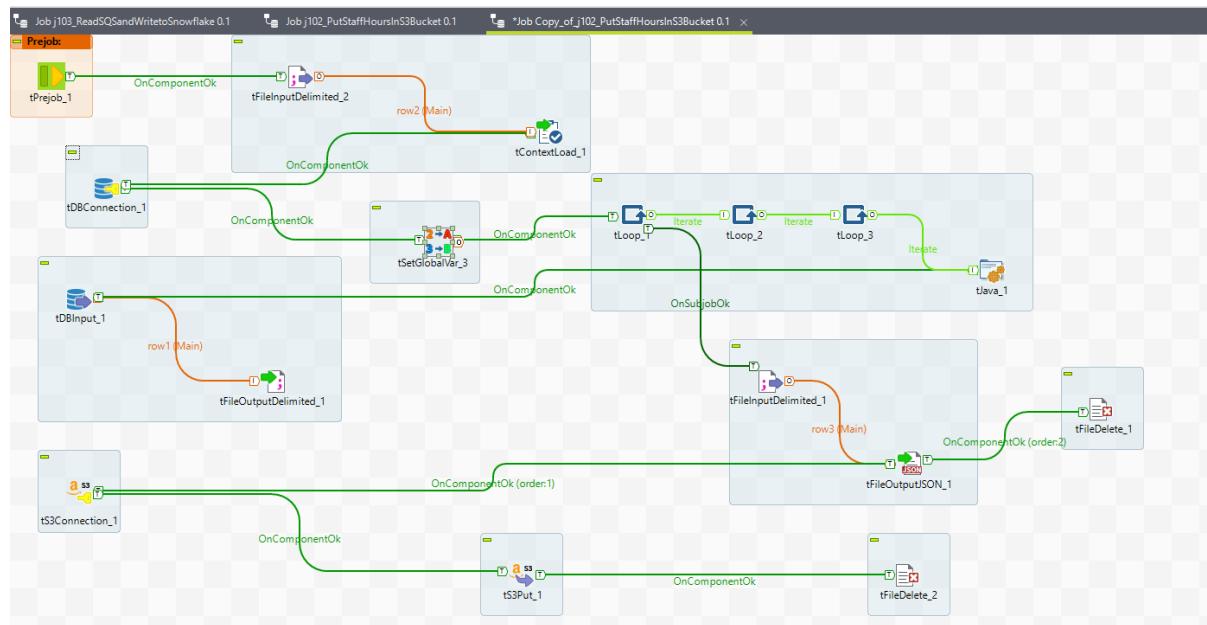
Arrange them in a similar layout to below.



Create the following connections between the components:

From Component	To Component	Join Type
<b>tPrejob</b>	1 <sup>st</sup> tFileInputDelimited	OnComponentOk trigger
1 <sup>st</sup> tFileInputDelimited	tContextLoad	Row Main
<b>tContextLoad</b>	tDBConnection	OnComponentOk trigger
<b>tDBConnection</b>	tSetGlobalVar	OnComponentOk trigger
<b>tSetGlobalVar</b>	1 <sup>st</sup> tLoop	OnComponentOk trigger
1 <sup>st</sup> tLoop	2 <sup>nd</sup> tLoop	Iterate
2 <sup>nd</sup> tLoop	3 <sup>rd</sup> tLoop	Iterate
3 <sup>rd</sup> tLoop	tJava	Iterate
<b>tJava</b>	tDBInput	OnComponentOk trigger
<b>tDBInput</b>	tFileOutputDelimited	Row Main
1 <sup>st</sup> tLoop	2 <sup>nd</sup> tFileInputDelimited	OnSubjobOk trigger
2 <sup>nd</sup> tFileInputDelimited	tFileOutputJSON	Row Main
<b>tFileOutputJSON</b>	1 <sup>st</sup> tFileDelete	OnComponentOk trigger
<b>tFileOutputJSON</b>	ts3Connection	OnComponentOk trigger
<b>ts3Connection</b>	ts3Put	OnComponentOk trigger
<b>ts3Put</b>	2 <sup>nd</sup> tFileDelete	OnComponentOk trigger

The resulting diagram should be similar to the following example:



Configure the components as follows:

**tPrejob** requires no configuration and just initiates a task each time the job starts and is guaranteed to run prior to the main job.

1<sup>st</sup> **tFileInputDelimited** reads in data used to populate the context variables

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

The screenshot shows the Talend Contexts editor interface. At the top, there are tabs for 'Job(j102\_PutStaffHoursInS3Buck...', 'Contexts(j102\_PutStaffHoursInS3...)', 'Component', and 'Run (Job j102\_PutStaffHoursInS3...)'. Below the tabs is a table with columns: Name, Type, Comment, Default, and Value. The 'Default' column contains dropdown menus, and the 'Value' column contains text input fields. The context variables listed are:

	Name	Type	Comment	Default	Value
1	BranchFirst	int   Integer		<input type="button" value="▼"/>	
2	BranchLast	int   Integer		<input type="button" value="▼"/>	
3	WeekFirst	int   Integer		<input type="button" value="▼"/>	
4	WeekLast	int   Integer		<input type="button" value="▼"/>	
5	DayFirst	int   Integer		<input type="button" value="▼"/>	
6	DayLast	int   Integer		<input type="button" value="▼"/>	
7	filename	String	C:\talend_files\Contexts\j102_PutStaffHoursInS3Bucket_Context.txt	<input type="button" value="▼"/>	

At the bottom left are icons for adding (+), deleting (x), up/down sorting (up/down arrows), and saving (disk). On the right, it says 'Default context environment' and 'Default' with a dropdown arrow.

Using a file in the following format:

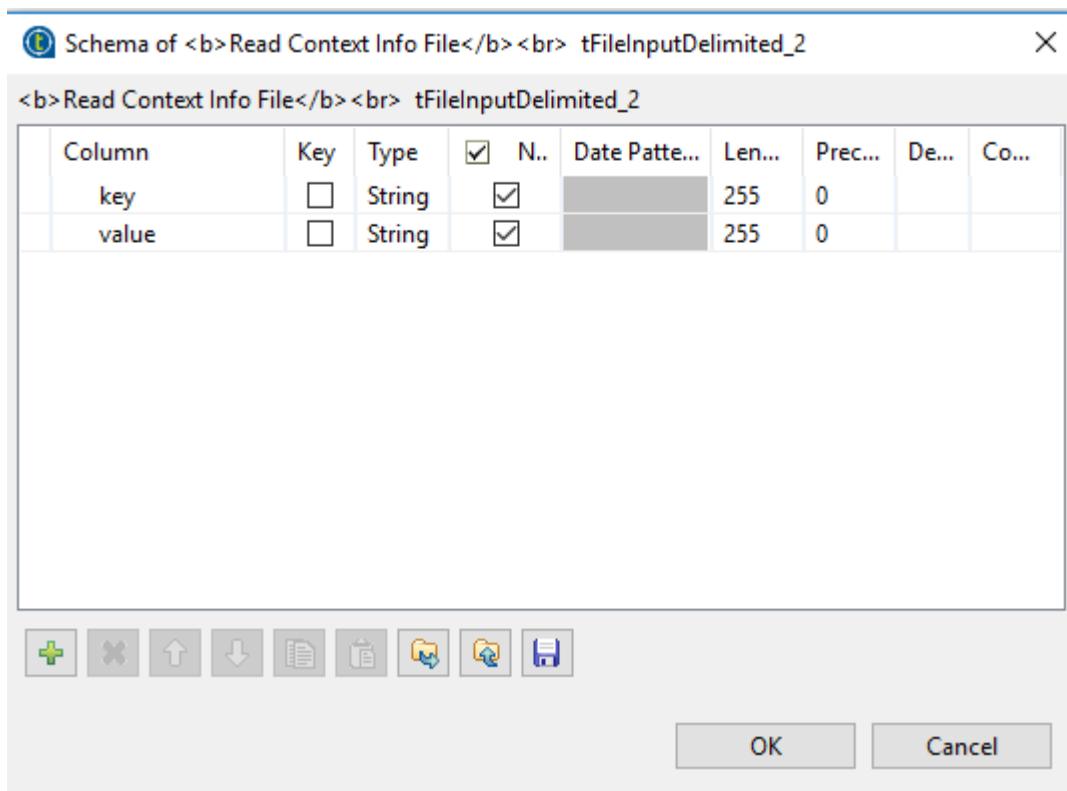
j102\_PutStaffHoursInS3Bucket\_Context.txt - Notepad

```
File Edit Format View Help
BranchFirst;1
BranchLast;8
WeekFirst;1
WeekLast;1
DayFirst;1
DayLast;7
```

Select the component in the design area and click the **Component** tab to display the editor.

The screenshot shows the Talend Component editor for the 'tFileInputDelimited\_2' component. At the top, there are tabs for 'Job(j102\_PutStaffHoursInS3Buck...', 'Contexts(j102\_PutStaffHoursInS3...)', 'Component', and 'Run (Job j102\_PutStaffHoursInS3...)'. Below the tabs is a configuration panel for the 'tFileInputDelimited\_2' component. The left sidebar has sections for 'Basic settings', 'Advanced settings', 'Dynamic settings', 'View', and 'Documentation'. The 'Basic settings' section is active. It includes fields for 'Property Type' (set to 'Built-In'), 'Schema' (set to 'Built-In'), 'File name/Stream' (set to 'context.filename'), 'Row Separator' (set to '\n'), 'Field Separator' (set to ','), 'Header' (set to 0), 'Footer' (set to 0), and 'Limit' (empty). There are also checkboxes for 'Skip empty rows' (checked), 'Uncompress as zip file' (unchecked), and 'Die on error' (unchecked). A note at the top states: "When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."

Use the context value of **filename** to identify the file and use the default values “\n” and “,” for row and field separators. There are zero header and footer lines and skip empty rows. Click on the **Edit Schema** ellipsis to check the definition.



The schema should consist of a key value pair both of string type. Press **OK** to close the editor.

Select the **View** option in the menu to display the label format.

Basic settings	Label format	<b>Read Context Info File</b>  _UNIQUE_NAME_
Advanced settings	Hint format	<b>_UNIQUE_NAME_</b>  _COMMENT_
Dynamic settings	Connection format	row
<b>View</b>		
Documentation		

Modify the label format from `_UNIQUE_NAME_` to `<b>Read Context Info File</b><br> _UNIQUE_NAME_`. This will give a user-friendly label to the component but also retain the component name to assist with monitoring. Labelling all subsequent components should follow a similar pattern.

**tContextLoad** takes the file data imported in the previous component and uses it to populate the context variables named in the file. No configuration is necessary on this component as the default values are sufficient.

Label the component “**Load File Data To Context**” using the method previously shown.

**tDBConnection** is used to connect to the Postgresql instance which is running in AWS RDS as configured previously.

Select **PostgreSQL** as the database with version v9 or later. The Host value will be the URL shown in the AWS RDS console for the instance. Port should be the standard value 5432 unless a different setup has been used. The Database name will be **RetailPOC** and the public schema is used. Finally, Username and Password refer to the local user set up in RDS for the instance. A shared connection is not used in this case and since the job will just be ran in Studio a data source alias is not required either and both check boxes can remain unchecked.

Go into the **view** option and label the component “**AWS RDS Postgres Connection**”.

**tSetGlobalVar** initialises the variable for holding the current week number, day number and store code in the global cache. Click the component in the design area and select the **Component** tab in the bottom pane to show the editor.

Key	Value
"vWeekNo"	...
"vDayNo"	...
"vStoreCode"	...

From the Basic Settings option, press the **green + icon** to create a new row and name it “**vWeekNo**”. Assign a default value of “”. Repeat the process for **vDayNo** and **vStoreCode** and add a user-friendly label “**Initialise Globals**” to the component.

Three nested loops are used to generate the store staff data.

The first **tLoop** component iterates through the defined stores.

 <b>Store Loop</b><br> tLoop\_1(tLoop\_1)

<input checked="" type="radio"/> Basic settings <input type="radio"/> Advanced settings <input type="radio"/> Dynamic settings <input type="radio"/> View <input type="radio"/> Documentation	<p>Loop Type  <input checked="" type="radio"/> For  <input type="radio"/> While</p> <p>From: context.BranchFirst</p> <p>To: context.BranchLast</p> <p>Step: 1</p> <p><input checked="" type="checkbox"/> Values are increasing</p>
---	--

Select **For** loop as the type use the branch context values to define **From** and **To** parameters with a step of 1. Once complete modify the label to be “**Store Loop**”

The second **tLoop** processes the week values.

 <b>Week Loop</b><br> tLoop\_2(tLoop\_2)

<input checked="" type="radio"/> Basic settings <input type="radio"/> Advanced settings <input type="radio"/> Dynamic settings <input type="radio"/> View <input type="radio"/> Documentation	<p>Loop Type  <input checked="" type="radio"/> For  <input type="radio"/> While</p> <p>From: context.WeekFirst</p> <p>To: context.WeekLast</p> <p>Step: 1</p> <p><input checked="" type="checkbox"/> Values are increasing</p>
---	--

Similar in definition to the previous component but use the week start and end values from the context. Label the component “**Week Loop**”.

The third **tLoop** traverses the day values.

 <b>Day Loop</b><br> tLoop\_3(tLoop\_3)

<input checked="" type="radio"/> Basic settings <input type="radio"/> Advanced settings <input type="radio"/> Dynamic settings <input type="radio"/> View <input type="radio"/> Documentation	<p>Loop Type  <input checked="" type="radio"/> For  <input type="radio"/> While</p> <p>From: context.DayFirst</p> <p>To: context.DayLast</p> <p>Step: 1</p> <p><input checked="" type="checkbox"/> Values are increasing</p>
---	--

Similar configuration again but select day values from the context. Label the component “**DayLoop**”.

**tJava** is used to set the current store, week and day values into the global variables for each iteration of the inner loop.

```
globalMap.put("vStoreCode", ((Integer)globalMap.get("tLoop_1_CURRENT_VALUE")).toString());
globalMap.put("vWeekNo", ((Integer)globalMap.get("tLoop_2_CURRENT_VALUE")).toString());
globalMap.put("vDayNo", ((Integer)globalMap.get("tLoop_3_CURRENT_VALUE")).toString());
```

Assign **vStoreCode** the current value of the outer loop, **vWeekNo** the current value of the middle loop and **vDayNo** the value of the inner loop as shown in the diagram.



*Make sure the CURRENT\_VALUE property is used and not CURRENT\_ITERATION. Although they will be the same if the loops all start from 1 this doesn't have to be the case. For example, if branches 3 to 4 only were specified in the context range then for the first iteration the value of CURRENT\_VALUE would be 3 but CURRENT\_ITERATION has a value of 1.*

Label the component “Set Globals”.

**tDBInput** retrieves data from the PostgreSQL database by running a query incorporating the global variables to generate simulated sales hours for each store on a daily basis.

Basic settings

Database: PostgreSQL

Use an existing connection:

Component List: tDBConnection\_1 - <b>AWS RDS Postgres Connection</b><br> tDBConnection\_1

Schema: Built-In

Table Name: tbstorestaff

Query Type: Built-In

Query:

```
"SELECT " + ((String)globalMap.get("vWeekNo")) + " AS WeekNo, "
+ ((String)globalMap.get("vDayNo")) + " AS DayNo,
store_code,
6 AS HoursWorked,
employee_no,
first_name,
last_name
FROM tbstorestaff
WHERE store_code = " + ((String)globalMap.get("vStoreCode"))
```

Select **PostgreSQL** as the database type and check the **Use an existing connection** option. Select the DB connection previously configured and the Table Name “**tbstorestaff**”. Select **Edit schema** to call the editor.

Column	Db Column	Key	Type	DB Type	<input checked="" type="checkbox"/> N..	Date Patter...	Length	Precis...
week_no	week_no	<input type="checkbox"/>	Integer	INT4	<input checked="" type="checkbox"/>			
day_no	day_no	<input type="checkbox"/>	Integer	INT4	<input checked="" type="checkbox"/>			
store_code	store_code	<input type="checkbox"/>	Integer	INT4	<input checked="" type="checkbox"/>			
hours_worked	hours_worked	<input type="checkbox"/>	Integer	INT4	<input checked="" type="checkbox"/>			
employee_no	employee_no	<input type="checkbox"/>	Integer	INT4	<input checked="" type="checkbox"/>			
first_name	first_name	<input type="checkbox"/>	String	VARC...	<input checked="" type="checkbox"/>			
last_name	last_name	<input type="checkbox"/>	String	VARC...	<input checked="" type="checkbox"/>			

Buttons at the bottom: +, X, Up, Down, Save, Undo, Redo, OK, Cancel.

Click the green plus icon to add each row and add the following entries to the schema.

Column	Type
week_no	Integer
day_no	Integer
store_code	Integer
hours_worked	Integer
employee_no	Integer
first_name	String
last_name	String

Press **OK** to close the schema editor and enter the following text into the Query field.

```
"SELECT " + ((String)globalMap.get("vWeekNo")) + " AS WeekNo, "
    + ((String)globalMap.get("vDayNo")) + " AS DayNo,
    store_code,
    6 AS HoursWorked,
    employee_no,
    first_name,
    last_name
FROM tbstorestaff
WHERE store_code =" + ((String)globalMap.get("vStoreCode"))
+ " ORDER BY RANDOM()
LIMIT 10"
```

The query models the rules of the proof of concept. Each store is open 12 hours a day and has a full time equivalent of 5, representing a total of 60 hours. Staff work in 6 hour shifts therefore a total of 10 staff from a pool of 20 is required for each store per day.



*Note the query is constructed using dynamic SQL to incorporate the global variables. There may be concerns about security vulnerability from SQL injection and if this was public facing such as a web site query I would agree. In this case it is only internal and doesn't present a risk but if concern remained it could be redesigned to use a parametrised stored procedure instead.*

Label the component “**Read Postgres Storestaff**”.

**tFileOutputDelimited** creates a temporary output file used to store the data retrieved from PostgreSQL. Ultimately the data will be written to a JSON format file however it is created using an iterative method requiring the output file to be regularly appended. JSON files in Talend don't have this option and would be overwritten with each iteration therefore the temporary delimited file which can be appended is used instead and this file used later to create the JSON document in a single pass.

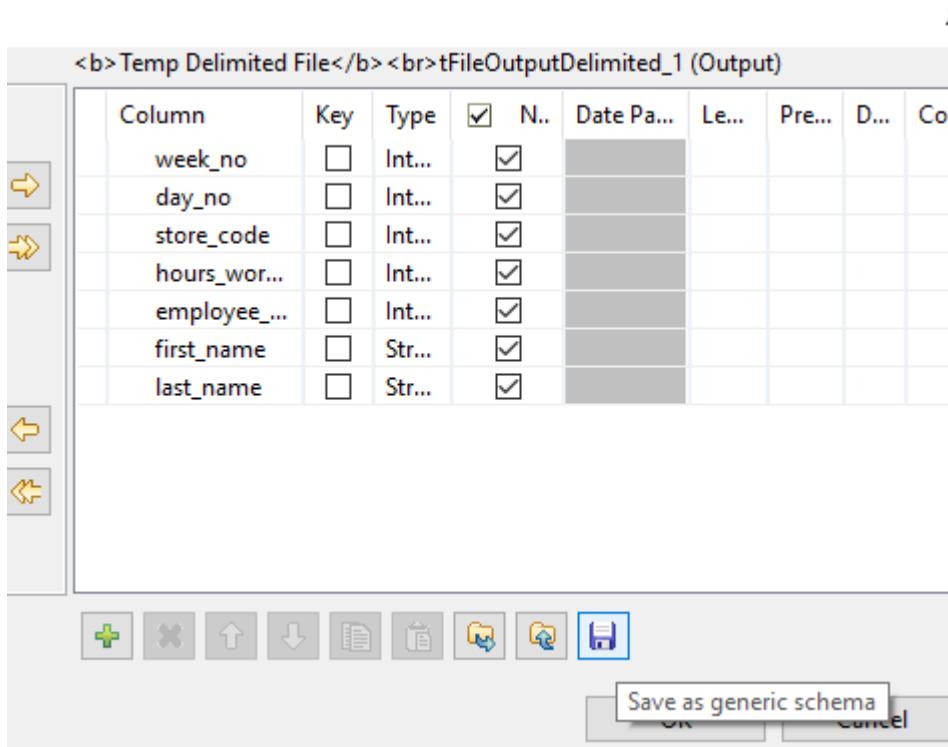
Enter the filename and leave the separators with the default values. Click on **Edit schema** to ensure it is the same as the tDBInput component. Ensure that the Append option is checked then label the component "**Temp Delimited File**".

The second **tFileInputDelimited** component is used to read the completed temporary file from the previous section. Initiation is by an **OnSubjobOk** trigger from the outer **tLoop** component. Using this type of trigger ensures that the loop iterations are fully complete including the writing of the temporary file, prior to it being consumed.

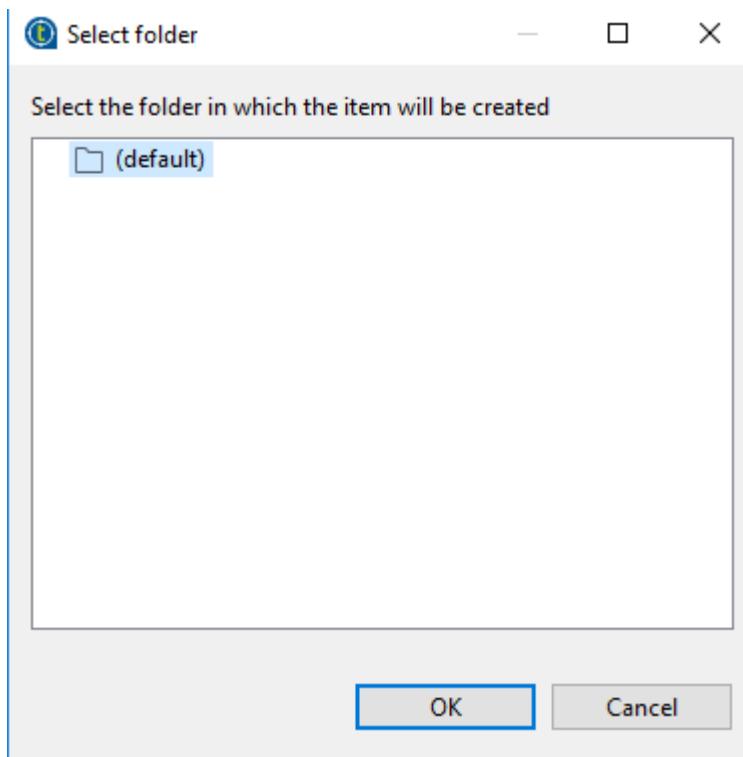
The schema will need to match the temporary file, and this can be achieved by clicking the ellipsis and manually entering the values. An alternative way that may be easier and avoid cumulative errors

is to save the schema from the **tFileOutputDelimited** component to the **repository** as a **generic schema**, which can then be used to set the schema on this component.

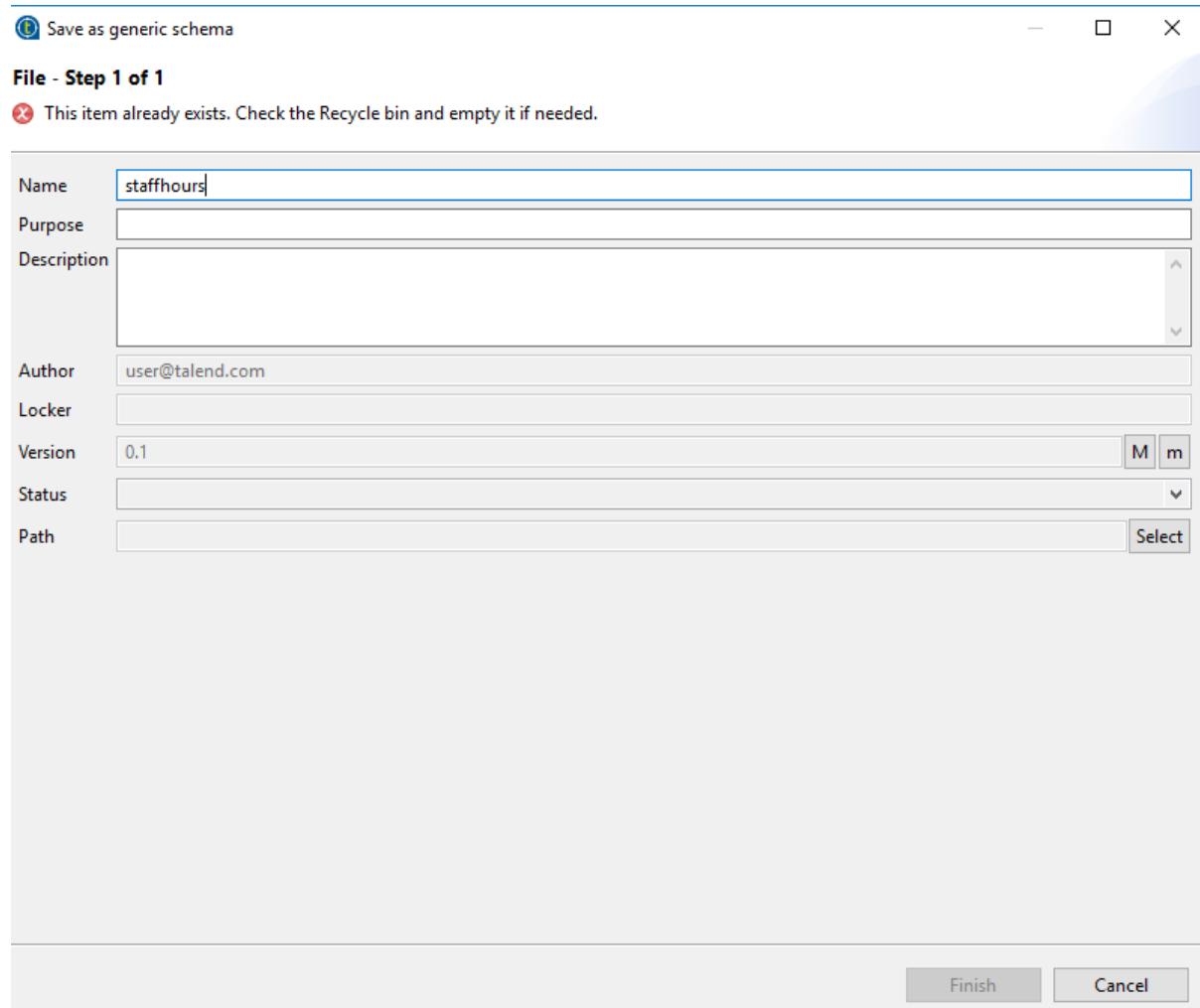
To do this return to the **tFileOutputDelimited** component and click the ellipsis to enter the schema editor.



Press the disk icon to Save as generic schema

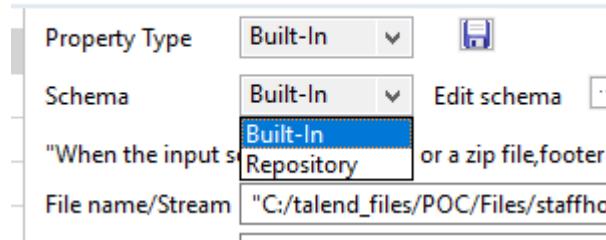


Click **OK** to accept the default location which equates to Generic schemas folder in the Metadata section of the repository.



Name the schema “**staffhours**” and click **OK** to complete, then return to the second **tFileInputDelimited** component.

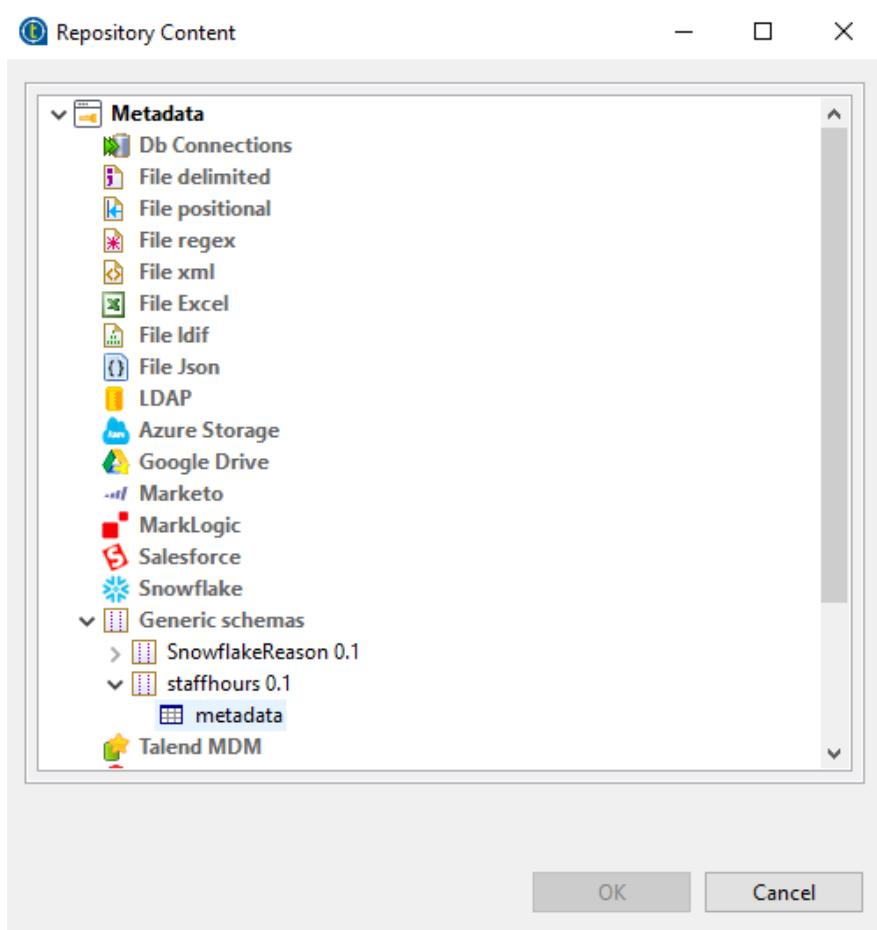
### | Temp File</b> <br> **tFileInputDelimited\_1**



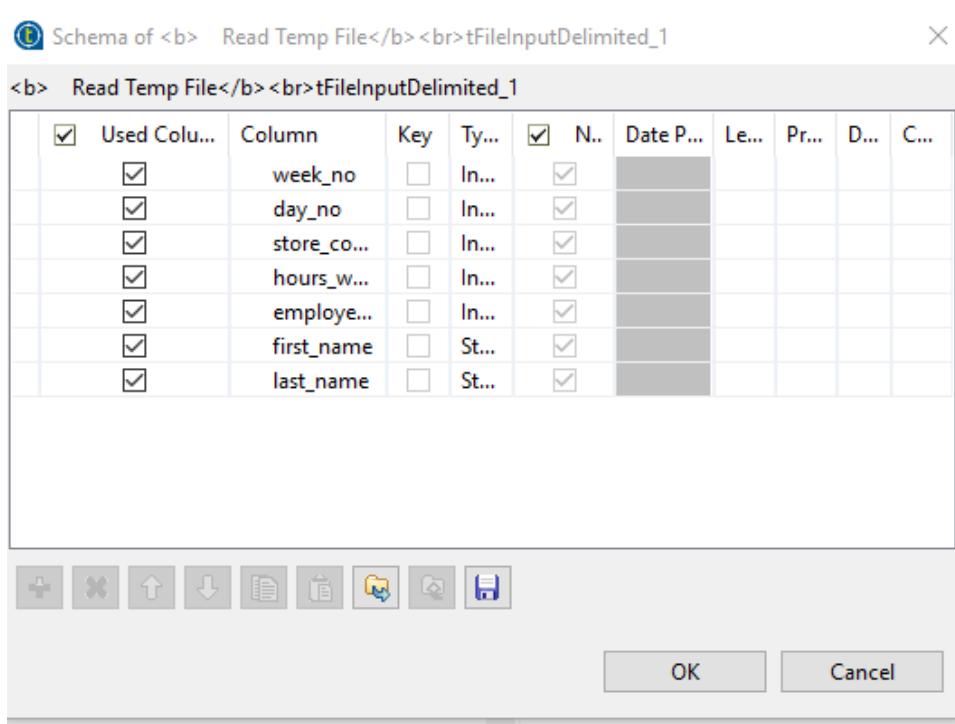
Select the **Repository** option from the Schema drop-down list which will add an edit control for the repository schema name.



Click the ellipsis next to the control to call the repository viewer.

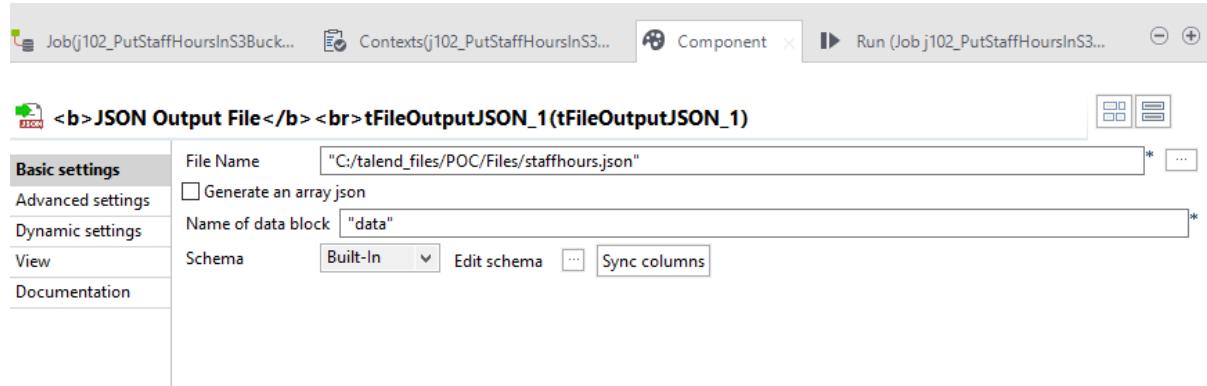


Expand the Generic schemas section. Select staffhours and click **OK** to complete. You can check the schema by clicking **Edit schema** and selecting the **View Schema** option.



File name should match the **tFileOutputDelimited** component and the separators will be the default values. There are no header or footer rows and empty rows can be skipped. Label the component “**Read Temp File**”.

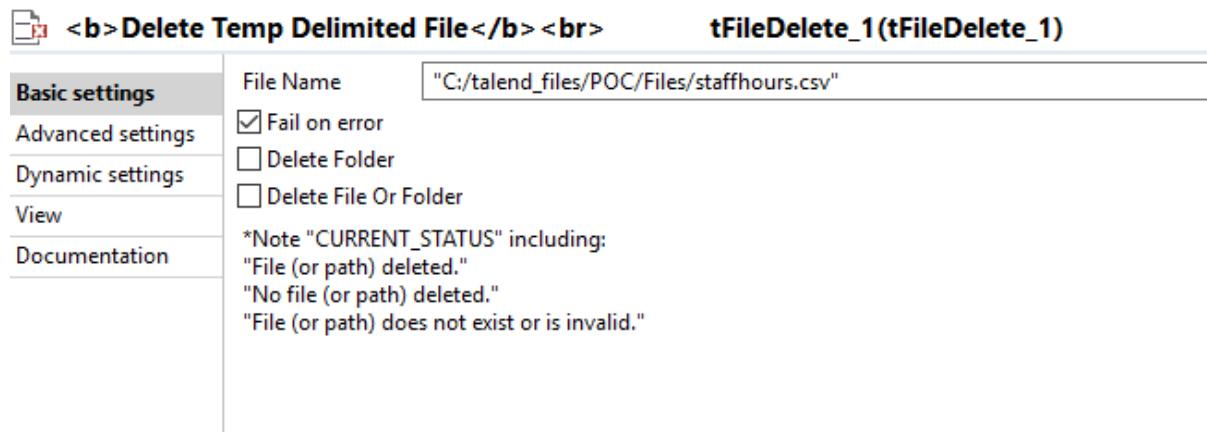
**tFileOutputJSON** will output the data read from the temporary file in JSON format. A very basic format is used for this POC with the fields of each record presented in a single data block with no nesting.



Inset the filename which will be an initial location prior to being stored in an S3 bucket. The name of the data block can be left as “**data**” and ensure the schema matches the input file by pressing **Sync columns**.

Label the component “**JSON Output File**”.

The first **tFileDelete** component removes the temporary delimited file after the JSON file has been created.



Enter the file name and label the component “**Delete Temp Delimited File**”.

Connection to the S3 bucket is initiated by the **tS3Connection** component.

 **AWS S3 Connection**   
 tS3Connection\_1(tS3Connection\_1)

<b>Basic settings</b>	Access Key "AKIAYYXALVCPRF4W3MXO"
<b>Advanced settings</b>	Secret Key *****
<b>Dynamic settings</b>	<input type="checkbox"/> Inherit credentials from AWS role
<b>View</b>	<input type="checkbox"/> Assume Role
<b>Documentation</b>	Region and Endpoint
	Region EU (London) *
	<input type="checkbox"/> Client-side Encrypt

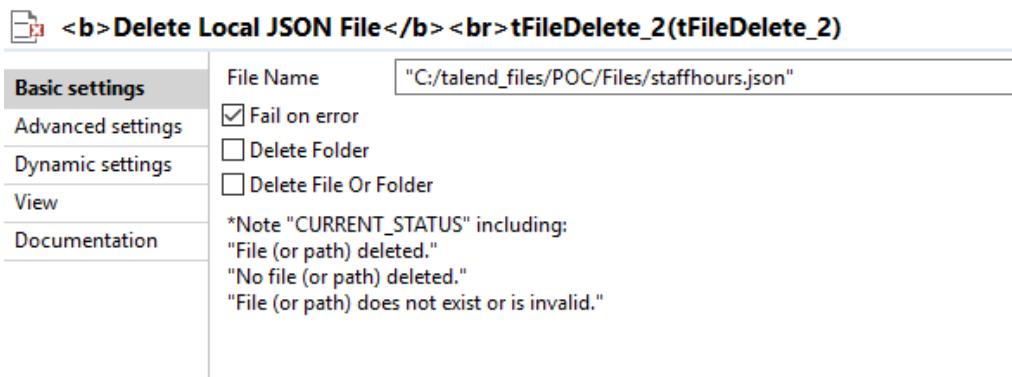
Use the keys generated when creating the bucket, as described earlier in this document and select the AWS Region which should match the bucket location. Label the component “**AWS S3 Connection**”

Having defined the connection, **tS3Put** is used to move the JSON file to the bucket.

 **Put JSON File in S3 Bucket**   
 tS3Put\_1(tS3Put\_1)

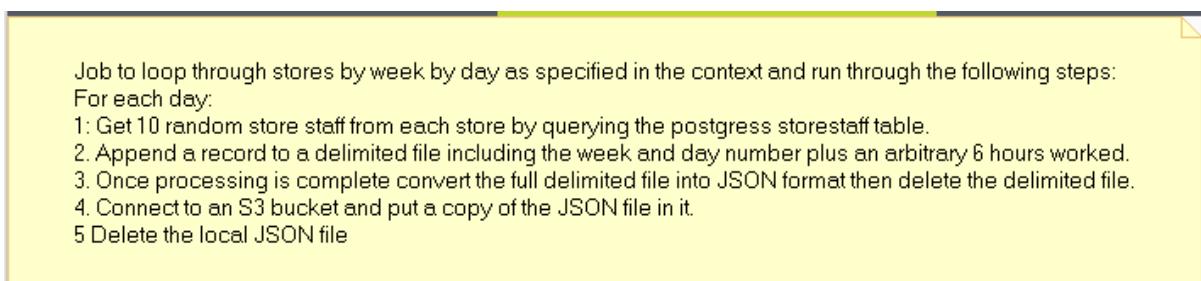
Check the Use an existing connection option and select the **S3 connection** from the list. Add the bucket name and the Key which is the S3 equivalent of file name. To prevent over writing a timestamp will be concatenated onto the name using Talend date functions. The key will take the form "RetailPOC\_" + TalendDate.formatDate("ddMMyyyy\_HH:mm:ss", TalendDate.getCurrentDate()). Enter the location of the JSON file and label the component “**Put JSON File in S3 Bucket**”.

The final action of the job is to remove the local JSON file after a copy has been stored in S3, using a second **tFileDelete** component.

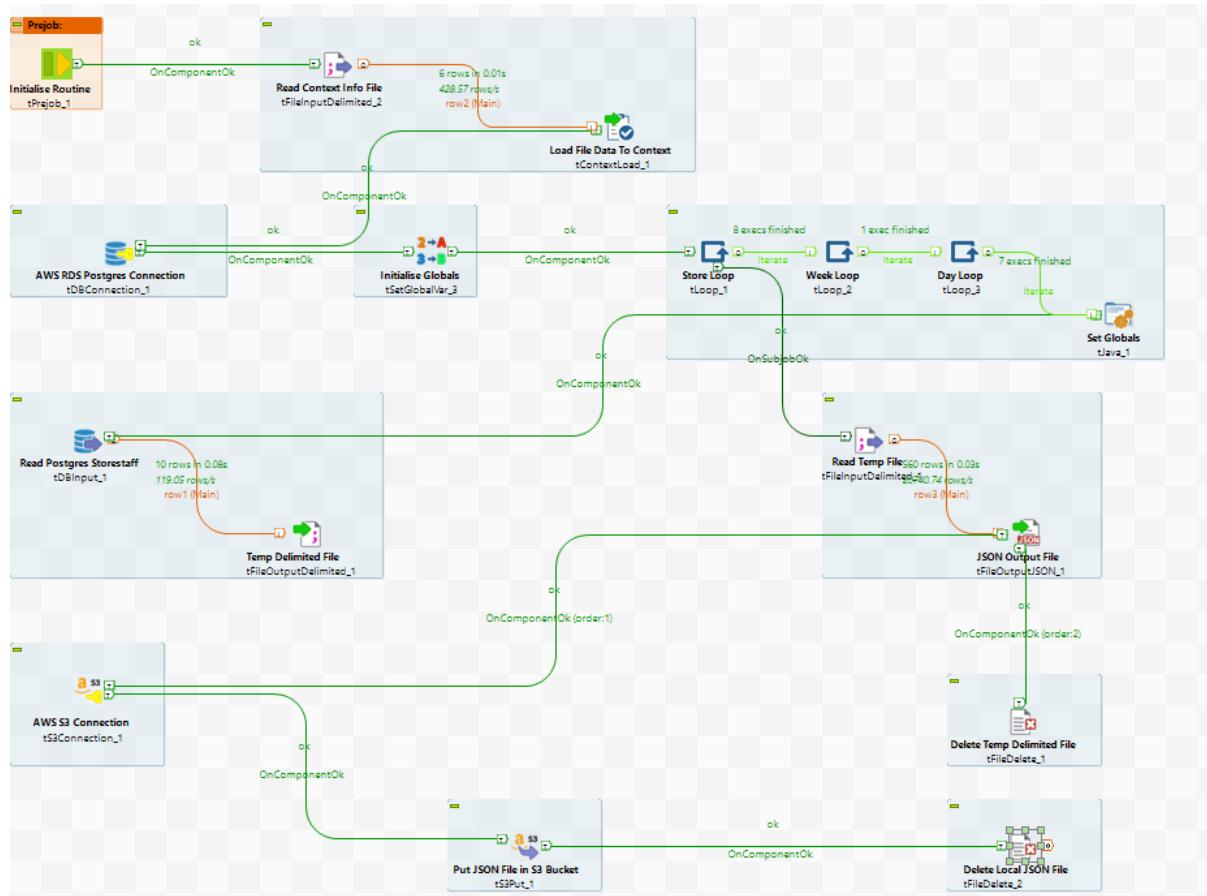


Enter the location of the JSON file and label the component “Delete Local JSON File”

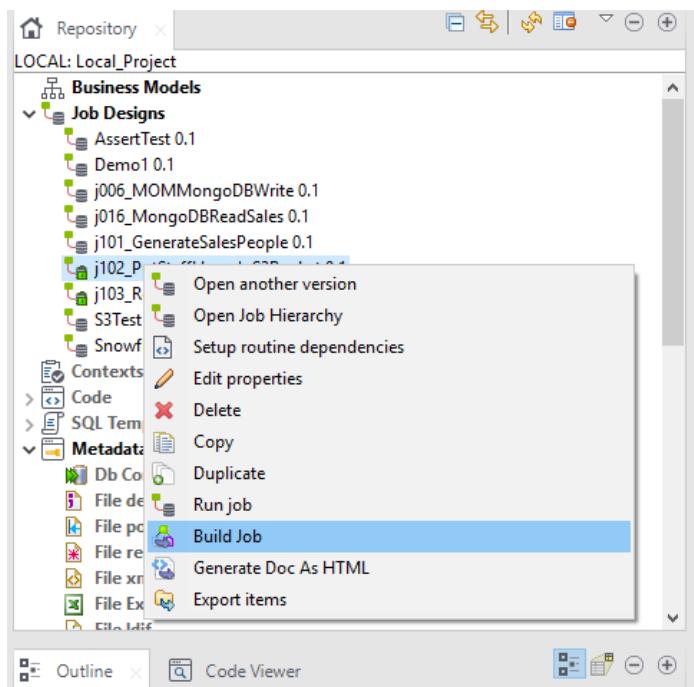
To complete the job, add a note to the design area similar to the following:



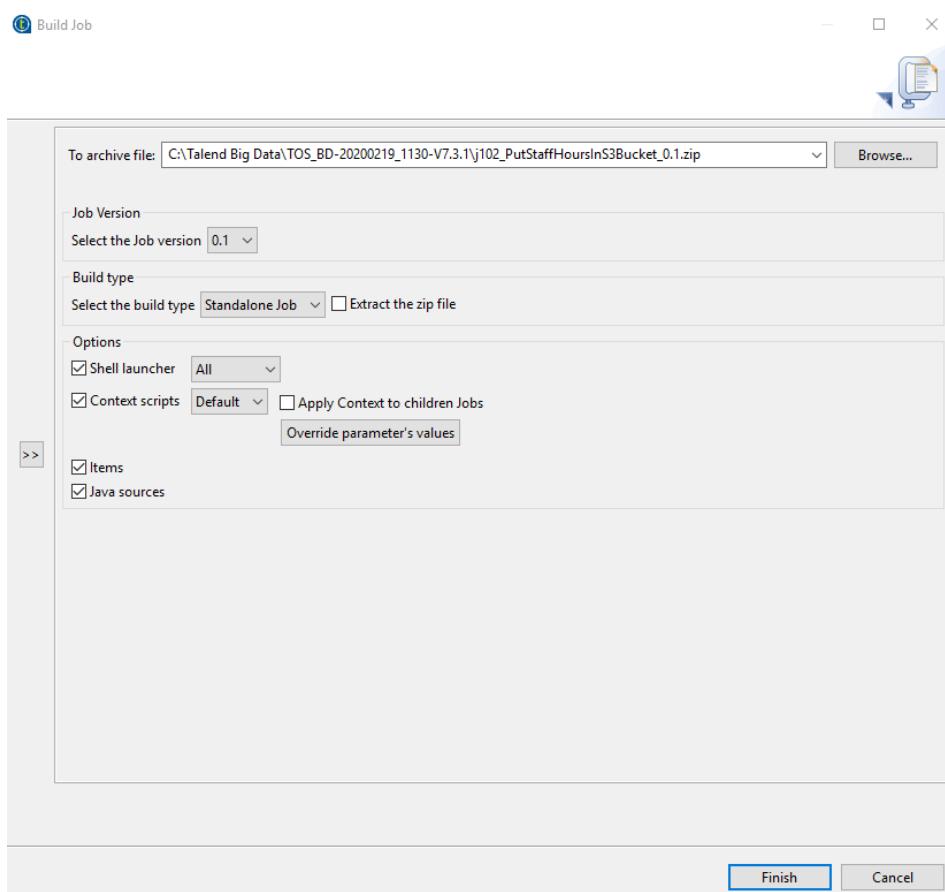
The finished procedure will look similar to the illustration below.



Unlike the previous job, which was just ran in Open Studio, this task will be deployed onto the AWS EC2 instance and ran as a stand-alone job. To begin the build, save the job then **right click** on its name under **Job Designs** in the repository pane and select **Build Job**.



The Build Job options will be displayed and in this case all default values are acceptable.



Click on **OK** to begin the build and once complete a zip file will be created in the location specified. In this case that is the root folder of Talend Studio although that can be changed if required.

### 13. Create job to write to Snowflake

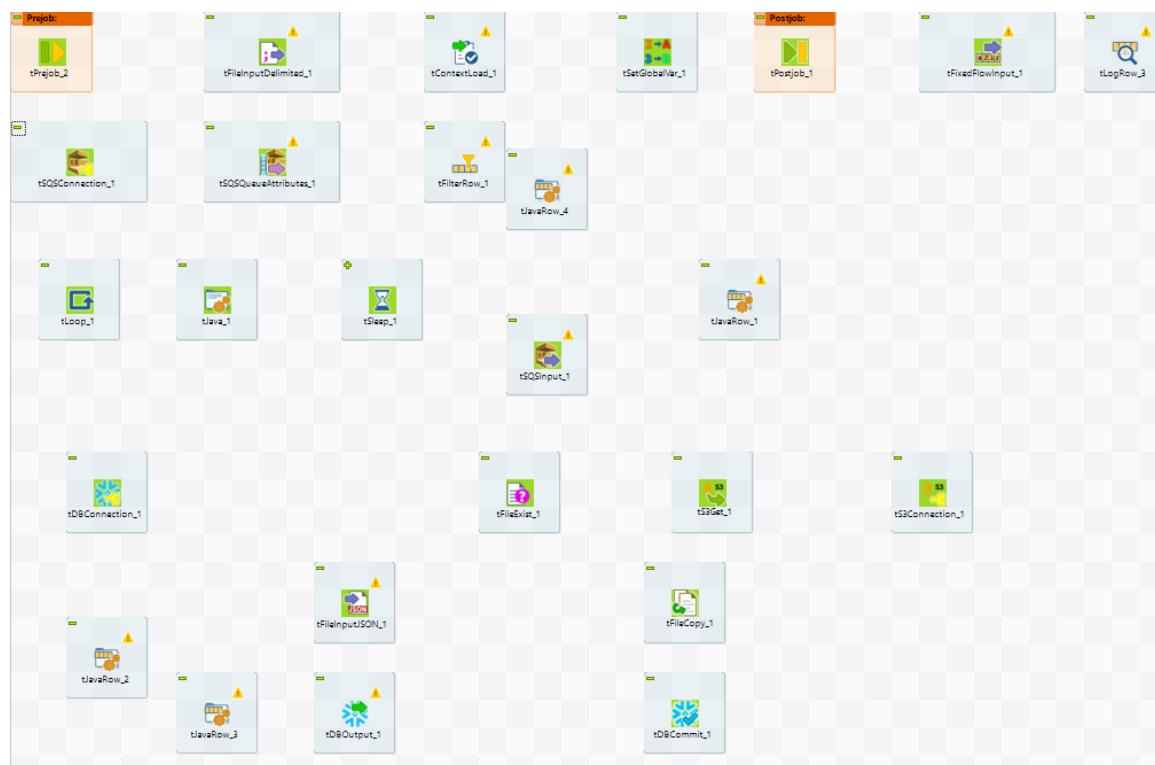
The third Talend job responds to messages arriving on the SQS queue. Messages are created by the AWS Lambda function that responds to files being placed in the S3 bucket by the previous job. After initiation this process will retrieve the file key from the message body and use it to download the JSON file from the S3 bucket. On receiving the file, the data, simulated staff timesheets, will be loaded into Snowflake for analysis later. The local JSON file will then be archived.

Create a new job in the repository called j103\_ReadSQSandWritetoSnowflake.

Drop the following components onto the design area:

tPrejob	tLoop	tS3Get
tFileInputDelimited	tJava	tFileExist
tContextLoad	tSQSQueueAttributes	tDBConnection
tSetGlobalVar	tSleep	tFileInputJSON
tPostjob	tFilterRow	tDBOutput
tFixedFlowInput	tJavaRow X 4	tFileCopy
tLogRow	tSQSInput	tDBCommit
tSQSConnection	tS3Connection	

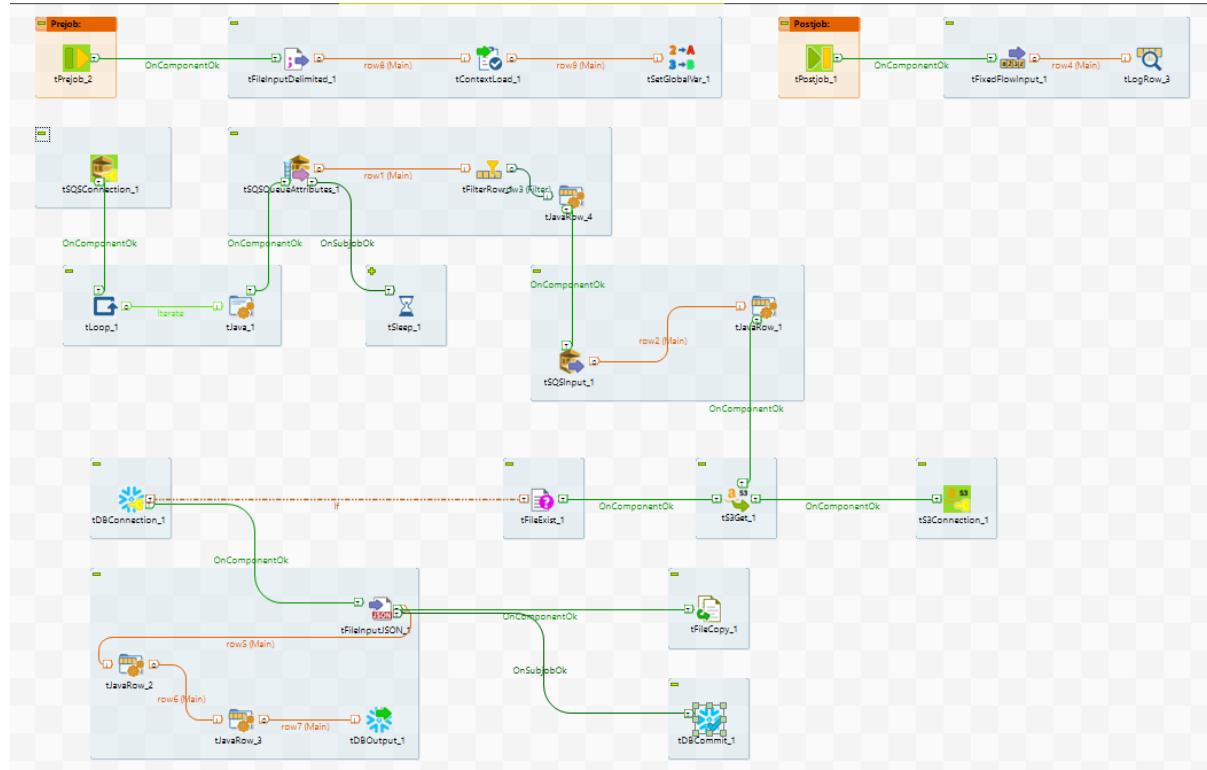
Arrange the components in a similar way to the following diagram.



Now create the joins between the components according to the following table.

From Component	To Component	Join Type
<b>tPrejob</b>	<b>tFileInputDelimited</b>	OnComponentOk trigger
<b>tFileInputDelimited</b>	<b>tContextLoad</b>	Row Main
<b>tContextLoad</b>	<b>tSetGlobalVar</b>	Row Main
<b>tPostjob</b>	<b>tFixedFlowInput</b>	OnComponentOk trigger
<b>tFixedFlowInput</b>	<b>tLogRow</b>	Row Main
<b>tSQSConnection</b>	<b>tLoop</b>	OnComponentOk trigger
<b>tLoop</b>	<b>tJava</b>	Iteration
<b>tJava</b>	<b>tSQSQueueAttributes</b>	OnComponentOk trigger
<b>tSQSQueueAttributes</b>	<b>tSleep</b>	OnSubjobOk trigger
<b>tSQSQueueAttributes</b>	<b>tFilterRow</b>	Row Main
<b>tFilterRow</b>	<b>1<sup>st</sup> tJavaRow</b>	Row Filter
<b>1<sup>st</sup> tJavaRow</b>	<b>tSQSInput</b>	OnComponentOk trigger
<b>tSQSInput</b>	<b>2<sup>nd</sup> tJavaRow</b>	Row Main
<b>2<sup>nd</sup> tJavaRow</b>	<b>tS3Get</b>	OnComponentOk trigger
<b>tS3Connection</b>	<b>tS3Get</b>	OnComponentOk trigger
<b>tS3Get</b>	<b>tFileExist</b>	OnComponentOk trigger
<b>tFileExist</b>	<b>tDBConnection</b>	Run if trigger
<b>tDBConnection</b>	<b>tFileInputJSON</b>	OnComponentOk trigger
<b>tFileInputJSON</b>	<b>tFileCopy</b>	OnComponentOk trigger
<b>tFileInputJSON</b>	<b>3<sup>rd</sup> tJavaRow</b>	Row Main
<b>3<sup>rd</sup> tJavaRow</b>	<b>4<sup>th</sup> tJavaRow</b>	Row Main
<b>4<sup>th</sup> tJavaRow</b>	<b>tDBOutput</b>	Row Main
<b>tFileInputJSON</b>	<b>tDBCommit</b>	OnSubjobOk trigger

Once complete the job in the design area should resemble the following illustration.



The components can now be configured using the following steps:

**tPrejob** requires no configuration and just initiates a task each time the job starts and is guaranteed to run prior to the main job.

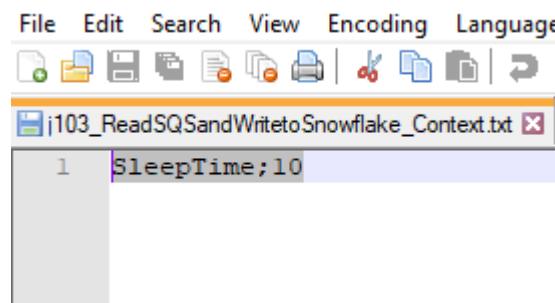
1<sup>st</sup> **tFileInputDelimited** reads in data used to populate the context variables

The screenshot shows the Talend interface with the 'Component' tab selected. A table lists context variables:

	Name	Type	Comment	Default
				Value
1	SleepTime	int   Integer		<input type="text"/>
2	filename	String	C:\talend_files\Contexts\j103_ReadSQSandWritetoSn	<input type="text"/>

Below the table are standard file operations icons (New, Open, Save, etc.) and a 'Default context environment' dropdown set to 'Default'.

Using a file in the following format:

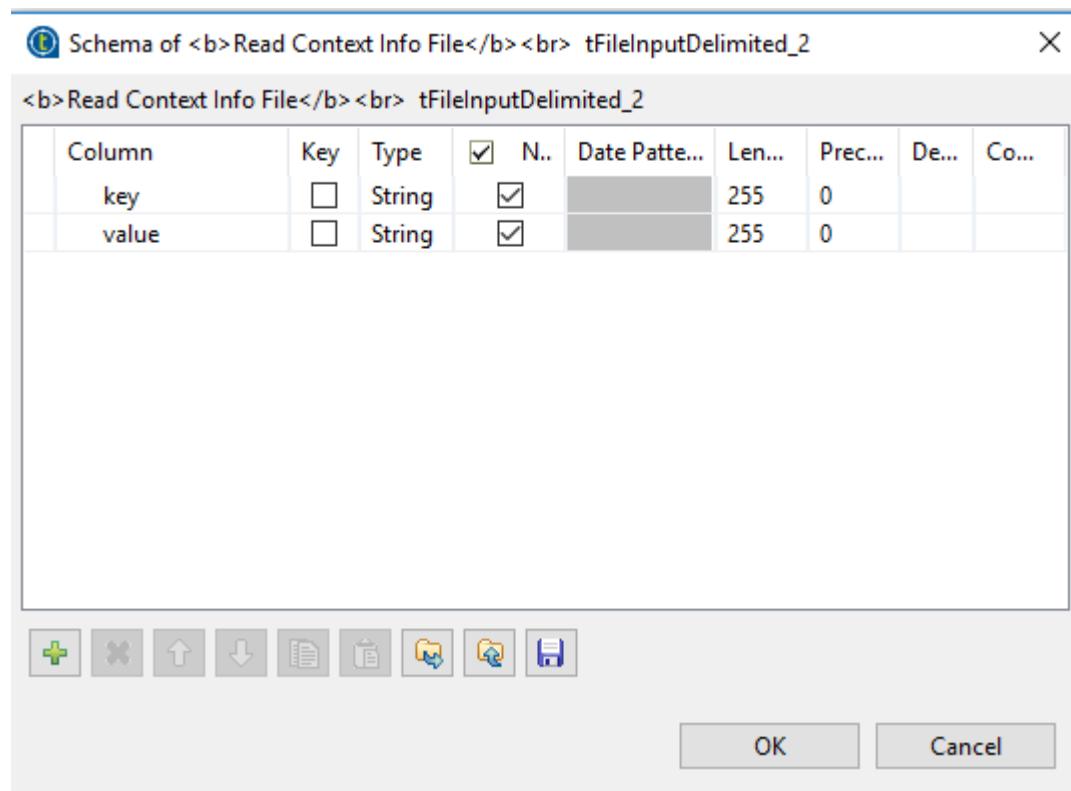


Select the component in the design area and click the **Component** tab to display the editor.

The screenshot shows the Talend interface with the 'Component' tab selected. The configuration for the 'tFileInputDelimited\_2' component is displayed:

- Basic settings** tab is active:
  - Property Type: Built-In
  - Schema: Built-In
  - File name/Stream: context.filename
  - Row Separator: "\n"
  - Field Separator: ";"
  - Header: 0
  - Footer: 0
  - Limit: (empty)
  - Skip empty rows
  - Uncompress as zip file
  - Die on error
- Advanced settings**, **Dynamic settings**, **View**, and **Documentation** tabs are visible but not active.

Use the context value of filename to identify the file and use the default values “\n” and “;” for row and field separators. There are zero header and footer lines and skip empty rows. Click on the **Edit Schema** ellipsis to check the definition.



The schema should consist of a key value pair both of string type. Press **OK** to close the editor.

Select the **View** option in the menu to display the label format.

Basic settings	Label format	<b>Read Context Info File</b>  _UNIQUE_NAME_
Advanced settings	Hint format	<b>_UNIQUE_NAME_</b>  _COMMENT_
Dynamic settings	Connection format	row
<b>View</b>		
Documentation		

Modify the label format from `_UNIQUE_NAME_` to `<b>Read Context Info File</b><br> _UNIQUE_NAME_`. This will give a user-friendly label to the component but also retain the component name to assist with monitoring. Labelling all subsequent components should follow a similar pattern.

**tContextLoad** takes the file data imported in the previous component and uses it to populate the context variables named in the file. No configuration is necessary on this component as the default values are sufficient.

Label the component “**Load File Data To Context**” using the method previously shown.

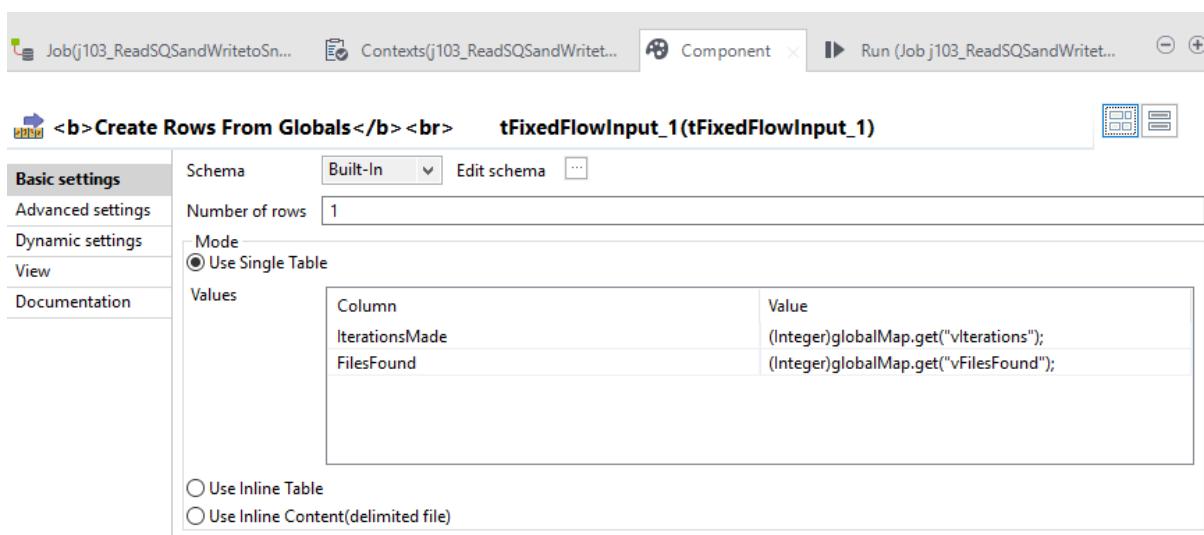
**tSetGlobalVar** initialises the variable for holding the number of iterations, file found status and S3 bucket key in the global cache. Click the component in the design area and select the **Component** tab in the bottom pane to show the editor.



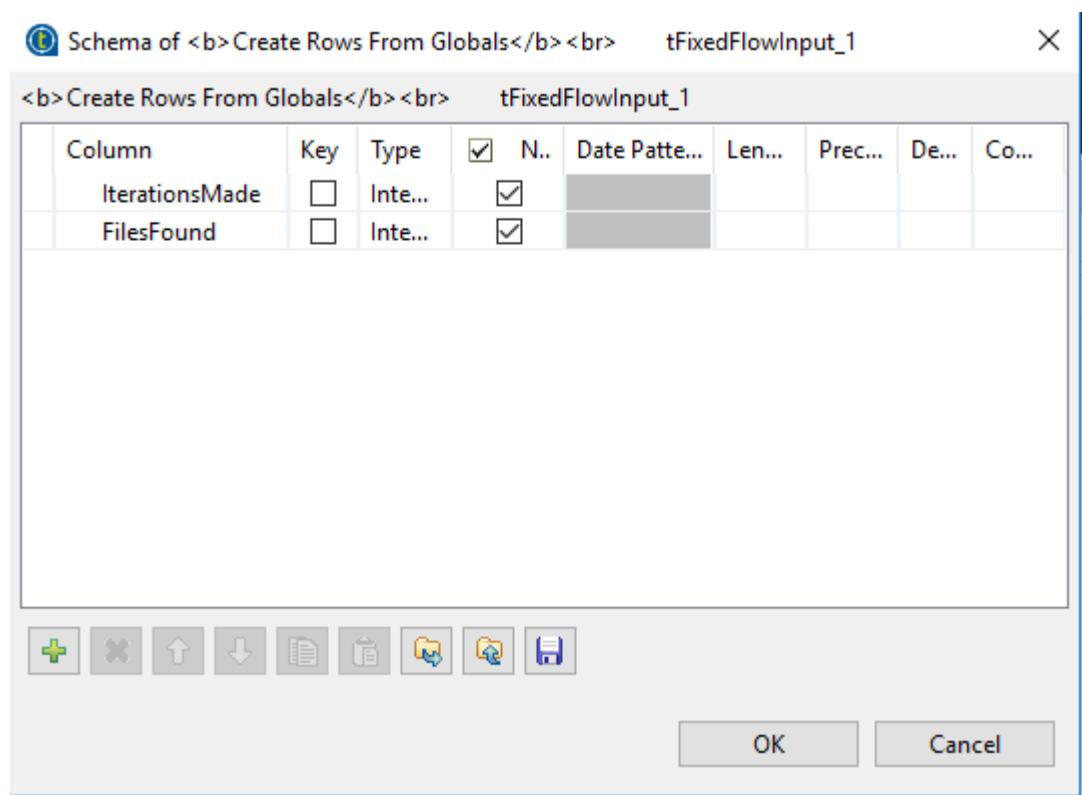
From the Basic Settings option, press the **green + icon** to create a row and name it “**vIterations**”. Assign a default value of 0. Repeat the process for “**vFilesFound**” and “**vBucketKey**” which has a default value of “” and add a user-friendly label “**Initialise Globals**” to the component.

**tPostjob** requires no configuration initiating a task each time the job completes and is guaranteed to run after the main job. Since the job runs on an infinite loop this component and its associated sub job will only be called in the event of termination, but it is useful for debugging.

**tFixedFlowInput** generates a data flow using the values of the global variables.



Click the **Edit schema** ellipsis to create the schema.



Click the green + icon to add 2 rows for "IterationsMade" and "FilesFound", both of integer type.  
Press **OK** to exit the editor.

Select single table mode and assign the following values to the columns:

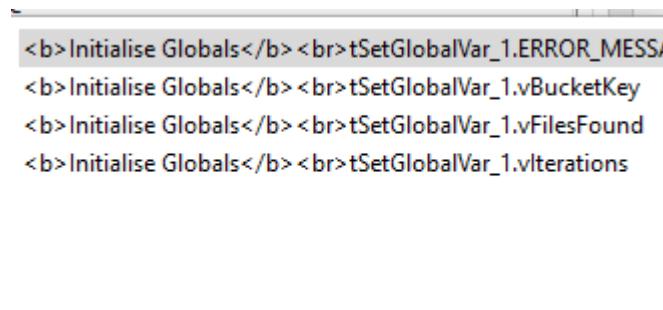
Column	Value
IterationsMade	((Integer)globalMap.get("vIterations"))
FilesFound	((Integer)globalMap.get("vFilesFound"))



A more reliable way to enter the global variable values is to allow Talend to look them up, avoiding transcription errors. Place the cursor in the value field and press CTRL and Enter simultaneously to bring up a list of system variables.

Column	Description: Error Message Global variable, property of component tSQSQueueAttributes [tSQS Queue Attributes tSQSQueueAttributes_1] Type: String Availability: After  Variable Name: ERROR_MESSAGE	Value
		<b> SQS Queue Attributes</b> tSQSQueueAttributes_1 (<b> Always True Loop</b>  tLoop_1.CURRENT_I <b> Always True Loop</b>  tLoop_1.ERROR_ME: <b> Archive JSON File</b>  tFileCopy_1.DESTINAT <b> Archive JSON File</b>  tFileCopy_1.DESTINAT <b> Archive JSON File</b>  tFileCopy_1.DESTINAT <b> Archive JSON File</b>  tFileCopy_1.ERROR_M <b> Archive JSON File</b>  tFileCopy_1.SOURCE_I <b> AWS S3 Connection</b>  tS3Connection_1.ERR <b> AWS SQS Connection</b>  tSQSConnection_1.I

Scrolling through the list will eventually find the correct variable but a filter can be applied to make the process easier. Find the component name where the variable is initialised which in this case is **tSetGlobalVar**. When the list is displayed type “**tse**” (case insensitive) and the list will reduce to the entries applicable to that component only including the global variables.



Move the cursor to the required entry and click to insert the value. Notice that by default it always casts the variable to a String. To change the value to an integer simply overtype the cast from String to Integer. Repeat the process for both global variables and label the component “**Create Rows From Globals**”.

A **tLogRow** component is used to display the data flow values on the standard output.

The screenshot shows the configuration dialog for the **tLogRow** component. The left sidebar has tabs for **Basic settings**, **Advanced settings**, **Dynamic settings**, **View**, and **Documentation**. The main area has the following settings:

- Schema:** Built-In
- Mode:**  Vertical (each row is a key/value list)
- Title printing mode:**  Print unique name
- Print content with log4j

Select the **Vertical** mode option for clarity and label the component “**Write To Log**”.

Connection to the AWS SQS queue is configured via the **tSQSConnection** component.

The screenshot shows the configuration dialog for the **tSQSConnection** component. The left sidebar has tabs for **Basic settings**, **Advanced settings**, **Dynamic settings**, **View**, and **Documentation**. The main area has the following settings:

Access Key	"AKIAYYXALVCPRF4W3MXO"
Secret Key	*****
Inherit credentials from AWS role	<input type="checkbox"/>
Assume Role	<input type="checkbox"/>
Region	EU (Ireland) <input type="button" value="▼"/>

The Access key identifier can be retrieved by viewing the “**My Security Credentials**” option for your account in the AWS console.

The screenshot shows the AWS Identity and Access Management (IAM) service in the AWS console. The left sidebar is titled 'Identity and Access Management (IAM)' and includes sections for 'Dashboard', 'Access management', 'Access reports', and 'Documentation'. The main content area is titled 'Your Security Credentials' and contains a list of credential types: 'Password', 'Multi-factor authentication (MFA)', and 'Access keys (access key ID and secret access key)'. A note below states: 'Use access keys to make programmatic calls to AWS from the AWS CLI, Tools for PowerShell, AWS SDKs, or direct AWS API calls. You can have a maximum of two access keys (active or inactive) at a time.' Below this is another note: 'For your protection, you should never share your secret keys with anyone. As a best practice, we recommend frequent key rotation. If you lose or forget your secret key, you cannot retrieve it. Instead, create a new access key and make the old key inactive. Learn more'.

Created	Access Key ID	Last Used	Last Used Region	Last Used Service	Status	Actions
Apr 14th 2021	AKIAYYYXALVCPRF4W3MXO	2021-05-05 19:39 UTC+0100	eu-west-2	s3	Active	<a href="#">Make Inactive</a>   <a href="#">Delete</a>

A 'Create New Access Key' button is located at the bottom of the access keys section. A note at the bottom right of the page says: 'Root user access keys provide unrestricted access to your entire AWS account. If you need long-term access keys, we recommend creating a new IAM user with limited permissions and generating access keys for that user instead. [Learn more](#)'.

The secret key will be in the file produced when the key pair were generated and should have been stored in a safe place. Select the region that matches the SQS queue location and label the component “**AWS SQS Connection**”.

Unlike the JMS and MOM components for queue monitoring, the SQS Input version does not have an option to keep monitoring the queue and will terminate after use. To provide monitoring functionality a **tLoop** is used with a perpetual condition to constantly run the job unless manually terminated.

The screenshot shows the configuration interface for a Talend component. On the left is a vertical toolbar with buttons for 'Basic settings', 'Advanced settings', 'Dynamic settings', 'View', and 'Documentation'. The main panel is titled '**Always True Loop**' and contains a sub-section '**tLoop\_1(tLoop\_1)**'. The configuration fields are as follows:

- Loop Type**: A radio button group where 'While' is selected.
- Declaration**: A code editor containing the Java declaration: `int i=11`.
- Condition**: A code editor containing the Java condition: `i>10`.
- Iteration**: A code editor containing the Java iteration: `i++`.

A **While** loop is selected and arbitrary conditions applied such that it was always equate to **true**, causing the loop to continue indefinitely. A label “**Always True Loop**” should be assigned to the component.

The **tJava** component is called on each iteration of the loop to increment the counter in the global map,

```

tJava_1(tJava_1)
Code


```

Integer vIts = (Integer)globalMap.get("vIterations");
vIts++;
globalMap.put("vIterations", vIts);

```


```

The java snippet retrieves the value of the global variable **vIterations** as an integer, increments it and sends the new value back to the global map. Label the component “**Increment Iterations Count**”.

**tSQSQueueAttributes** allows a peek at the SQS queue status allowing the program to decide whether to apply message retrieval logic for the current iteration of the loop or not. A lot of information about the queue is returned but the item of interest is the **number of messages**.

```

tSQSQueueAttributes_1(tSQSQueueAttributes_1)
Connection
 Use an existing connection
Component List tSQSConnection_1 - AWS SQS Connection
Queue (Name or URL) https://sns.eu-west-2.amazonaws.com/602839689375/RetailPOC.fifo
Schema Built-In Edit schema ...
 Die on error

```

Check the “**Use an existing connection**” option and select the connection component from the list. Add the queue name which can be retrieved from the AWS SQS console.

Details	
Name	RetailPOC.fifo
Type	FIFO
Encryption	-
URL	https://sns.eu-west-2.amazonaws.com/602839689375/RetailPOC.fifo
More	

Make sure to place the value in double quotes as it is a string, otherwise Talend will treat it as a variable and cause an error. Clicking the **Edit schema** ellipsis will show the information returned from SQS.

The first item “**ApproximateNumberOfMessages**” is the metric that is of interest to this procedure. Label the component “**SQS Queue Attributes**”.

The **tSleep** component will introduce a pause into the loop iteration controlled by the value of the context variable defined in the context file.

Attached by an OnComponentOK it will be fired every iteration of the loop. Label the component “**Sleep Between Iterations**”

A **tFilterRow** is used to control the program flow based on the SQS status data returned.

Click the green + icon to add a filter row and select the “**ApproximateNumberOfMessages**” option from the InputColumn drop-down list. Leave the Function empty and select “**Greater than**” as the Operator. Add a Value of 0 to allow the filter when message(s) are present and add the label “**SQS Messages Found**”.

A **tJavaRow** component is used to increment the message found count in the global map.

**tJavaRow\_4(tJavaRow\_4)**

**Code**

```
Integer vFf = (Integer)globalMap.get("vFilesFound");
vFf++;
globalMap.put("vFilesFound", vFf);
```

Very similar to the earlier **tJava** component but called via a filer Row rather than a trigger. The java snippet retrieves the value of the global variable **vFilesFound** as an integer, increments it and sends the new value back to the global map. Label the component “**Increment Found Count**”.

SQS Messages are consumed by the **tSQSInput** component. It is configured using the same information as the **tSQSQueueAttributes** component.

**tSQSInput\_1(tSQSInput\_1)**

**Basic settings**

Connection  
 Use an existing connection  
Component List tSQSConnection\_1 - AWS SQS Connection

Queue (Name or URL) https://sqs.eu-west-2.amazonaws.com/602839689375/RetailPOC.fifo

Schema Built-In

Read standard attributes (ApproximateFirstReceiveTimestamp, ApproximateReceiveCount, SenderId, SentTimestamp)  
 Read custom user attributes  
 Custom visibility timeout  
 Custom wait time  
 Delete the messages while streaming  
 Read all messages from the queue

Max number of message to return per request (from 1 to 10) 1

Die on error

Check the “**Use an existing connection**” option and select the connection component from the list. Add the queue name which can be retrieved from the AWS SQS console as previously shown. Label the component “**Consume SQS Queue**”.

The object key for the file in S3 storage is contained in the body of the message, having been placed there by the AWS Lambda function. A second **tJavaRow** component is used to assign the key value to the **vBucketKey** variable in the global map.

The simple java snippet overwrites the current value with the contents of the message body. Label the component “**Write Bucket Key To Global**”.

The simple java snippet overwrites the current value with the contents of the message body. Label the component “**Write Bucket Key To Global**”.

Connection to the S3 bucket is initiated by the **tS3Connection** component.

Use the keys generated when creating the bucket, as described earlier in this document and select the AWS Region which should match the bucket location. Label the component “**AWS S3 Connection**”

The **tS3Get** component is used to retrieve the file from the S3 bucket using the global bucket key variable assigned previously.

Check the “**Use an existing connection**” option and select the connection component from the list. Add the bucket name “**emeraldmill.sales**” and retrieve the variable from the global map for the key, **((String)globalMap.get("vBucketKey"))**. In the File editor add the full path name of the local file that it will be saved to, “**C:/talend\_files/POC/Files/In/staffhours.json**”. Label the component “**Get JSON File From S3 Bucket**”.

A **tFileExist** component is used to check the success of the previous operation by validating the existence of the local file copy after download from S3.

**<b>Check Staff Hours File Exists</b> <br> tFileExist\_1(tFileExist\_1)**

**Basic settings**

File name/Stream: "C:/talend\_files/POC/Files/In/staffhours.json"

Advanced settings

Dynamic settings

View

Documentation

Add the name of the local file into the **File name/stream** editor and label the component “**Check Staff Hours File Exists**”

A **tDBConnection** component is configured to access the database set up earlier in this document using the **SnowflakeDB** account.

 This proof of concept is set up to use a trial Snowflake account, but it could easily be converted to use a similar service such as AWS Redshift or Google cloud platform (GCP) Big Query. Alternatively, any ANSI compliant RDBMS, cloud based or local could be used.

**<b>Snowflake Connection</b> <br> tDBConnection\_1(tDBConnection\_1)(Snowflake)**

**Basic settings**

Database: Snowflake

Property Type: Built-In

Account: "ZY40898"

Authentication Type: Basic

User Id: "JohnTucker1961"

Password: \*\*\*\*\*

Warehouse: "COMPUTE\_WH"

Schema: "PUBLIC"

Database: "RETAIL\_POCT"

Advanced settings

Dynamic settings

View

Documentation

Select **Snowflake** as the database and enter account number together with **User Id & Password**. The **Warehouse**, **Schema** and **Database** fields refer to the Snowflake setup earlier, “**COMPUTE\_WH**”, “**PUBLIC**” and “**RETAIL\_POCT**” in this case. Label the component “**Snowflake Connection**”.

The local JSON file downloaded from S3 is accessed by the **tFileInputJSON** component.

**<b>Read JSON Staff Hours File</b> <br> tFileInputJSON\_1(tFileInputJSON\_1)**

**Basic settings**

Property Type: Built-In

Schema: Repository

GENERIC:staffhours - metadata

Read By: JsonPath

API version: 2.1.0

Use Url

Filename: "C:/talend\_files/POC/Files/In/staffhours.json"

Loop Json query: "\$.data[\*]"

Mapping

Column	Json query
week_no	"week_no"
day_no	"day_no"
store_code	"store_code"
hours_worked	"hours_worked"
employee_no	"employee_no"
first_name	"first_name"
last_name	"last_name"

Die on error

Advanced settings

Dynamic settings

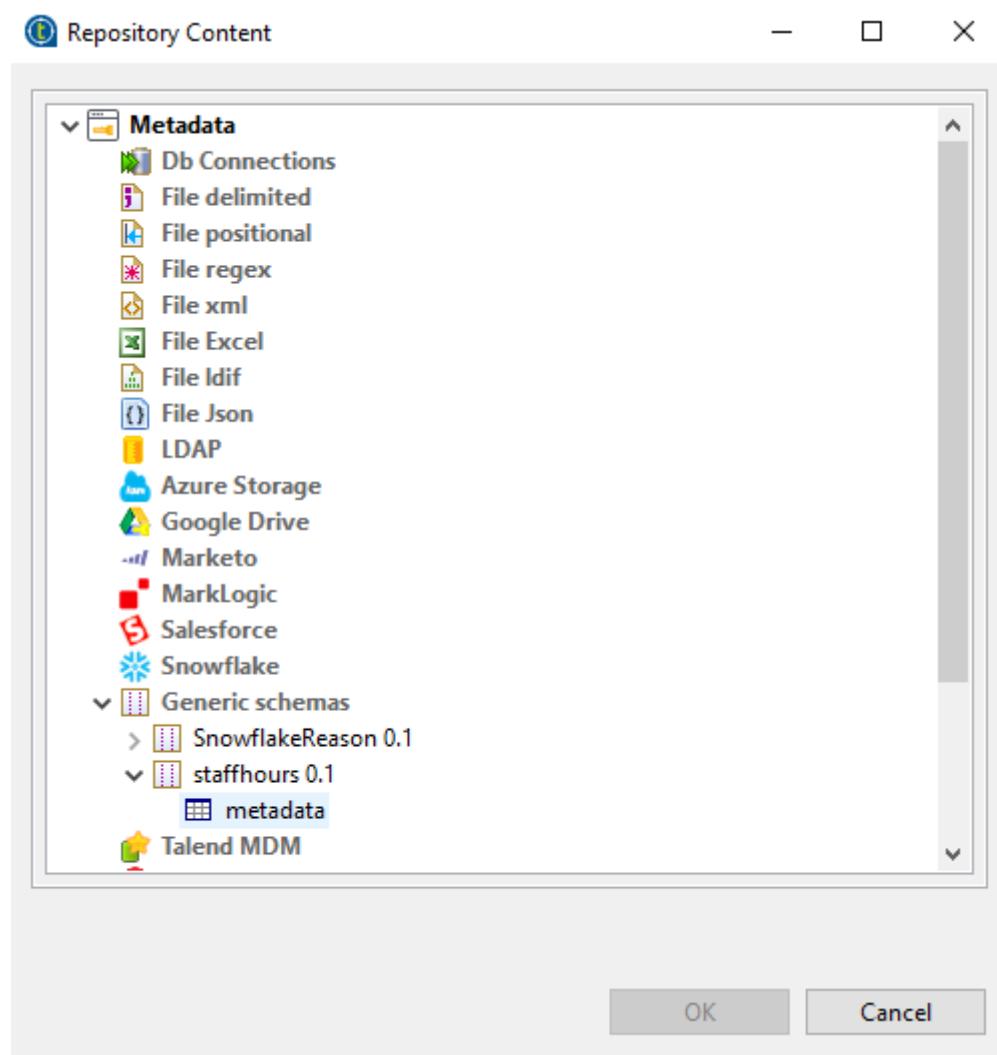
View

Documentation

A schema created for the previous job and saved in the repository as “**staffhours**” can be reused as the JSON file contains the same field layout. Select the **Repository** option from the Schema drop-down list which will add an edit control for the repository schema name.



Click the **ellipsis** next to the control to call the repository viewer.



Expand the **Generic schemas** section. Select **staffhours** and click **OK** to complete. You can check the schema by clicking **Edit schema** and selecting the **View Schema** option.

The file will be read using **JsonPath** which is the JSON equivalent of XPath, utilising the latest version of the API 2.1.0. Filename will be the local JSON file path and the loop expression should be **"\$.data[\*]"** which will retrieve data for each occurrence of the “**data**” node.

```

1  [
2    {
3      "data": [
4        {
5          "store_code": 1,
6          "day_no": 1,
7          "hours_worked": 6,
8          "employee_no": 6,
9          "last_name": "Nixon",
10         "week_no": 1,
11         "first_name": "Dwight"
12       },
13       {
14         "store_code": 1,
15         "day_no": 1,
16         "hours_worked": 6,
17         "employee_no": 9,
18         "last_name": "Roosevelt",
19         "week_no": 1,
20         "first_name": "Dwight"
21     }
22   ]
23 ]

```

In the mapping section the values in the Json query column should match the previous column but be enclosed in double quotes.

**tFileInputJSON\_1(tFileInputJSON\_1)**

**Advanced settings**

- Advanced separator (for numbers)
- Use the loop node as root

Encoding: UTF-8

Ensure that the “**Use the loop node as root**” is checked in the **Advanced settings** section then label the component “**Read JSON Staff Hours File**”.

A third **tJavaRow** component is used to clear the value of the bucket key global variable.

**tJavaRow\_2(tJavaRow\_2)**

**Basic settings**

Schema: Built-In

Code:

```

//Code generated according to input schema and output schema
output_row.week_no = input_row.week_no;
output_row.day_no = input_row.day_no;
output_row.store_code = input_row.store_code;
output_row.hours_worked = input_row.hours_worked;
output_row.employee_no = input_row.employee_no;
output_row.first_name = input_row.first_name;
output_row.last_name = input_row.last_name;

globalMap.put("vBucketKey", "");

```

To ensure the data flow passes through the component click the **Sync columns** button followed by **Generate code**. This will generate code assigning each input value to a corresponding output. Failure to complete these steps would mean the input data was lost to the chain downstream from this component.



*Always do the previous steps first prior to any alterations otherwise Generate code will delete any work you have done in the editor and replace with the input output assignment.*

Add this line below the generated code to clear the variable “`globalMap.put("vBucketKey", "");`” and label the component “**Clear S3 Bucket Key**”.

A fourth **tJavaRow** component is used to re-order the output into the format requires by **Snowflake**.

Click on the **Edit schema** ellipsis to call the editor.

Column	Key	Type	N.	Date Pattern...	Length	Precisi...	Def...	Com...
week_no	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
day_no	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
store_code	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
hours_worked	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
employee_no	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
first_name	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
last_name	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					

Column	Key	Type	N.	Date Pattern...	Length	Precisi...	Def...	Com...
STORE_CODE	<input type="checkbox"/>	Inte...	<input checked="" type="checkbox"/>					
DAY_NO	<input type="checkbox"/>	Inte...	<input checked="" type="checkbox"/>					
HOURS_WORK...	<input type="checkbox"/>	Inte...	<input checked="" type="checkbox"/>					
EMPLOYEE_NO	<input type="checkbox"/>	Inte...	<input checked="" type="checkbox"/>					
LAST_NAME	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
WEEK_NO	<input type="checkbox"/>	Inte...	<input checked="" type="checkbox"/>					
FIRST_NAME	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					

Modify the output schema by changing the column name to **upper case** and changing the order as shown in the diagram above. Click **OK** to close the editor then **Generate code**. The correct output should be produced. Label the component “**Reorder Output**”.

Writing the data to the SnowflakeDB table is handled by a **tDBOutput** component.

<b>Basic settings</b>	Database <input type="button" value="Snowflake"/> Apply
Advanced settings	Connection Component <input type="button" value="Snowflake Connection"/> tDBConnection_1
Dynamic settings	Table "TBSTAFFHOURS"
View	Schema Built-In <input type="button" value="Edit schema"/> ... <input type="button" value="Sync columns"/>
Documentation	Table Action NONE <input type="button" value="NONE"/> Output Action INSERT <input type="button" value="INSERT"/>

Click **Sync columns** to incorporate any changes made to the previous component. Select **Snowflake** as the database and the Snowflake connection component from the list. Set the table name to “**TBSTAFFHOURS**” with the Output Action of “**INSERT**”. Label the component “**Write to Snowflake**”.

A **tFileCopy** component is used to archive the local JSON file after use.

Enter the full path of JSON file "**C:/talend\_files/POC/Files/In/staffhours.json**" and the destination directory "**C:/talend\_files/POC/Files/In/done**". Check the **Rename** option and add the following expression for the Destination filename "**staffhours\_** +   
**TalendDate.formatDate("ddMMyyyy\_HHmmss",TalendDate.getCurrentDate()) + ".json"**.

Check the **Remove source file**, **Replace existing file** and **Create the directory if it doesn't exist** options and label the component "**Archive JSON File**".

The final component is **tDBCommit** to commit Snowflake changes to the database.

Select **Snowflake** as the database and the DB connection from the list. Check the **Close Connection** box.



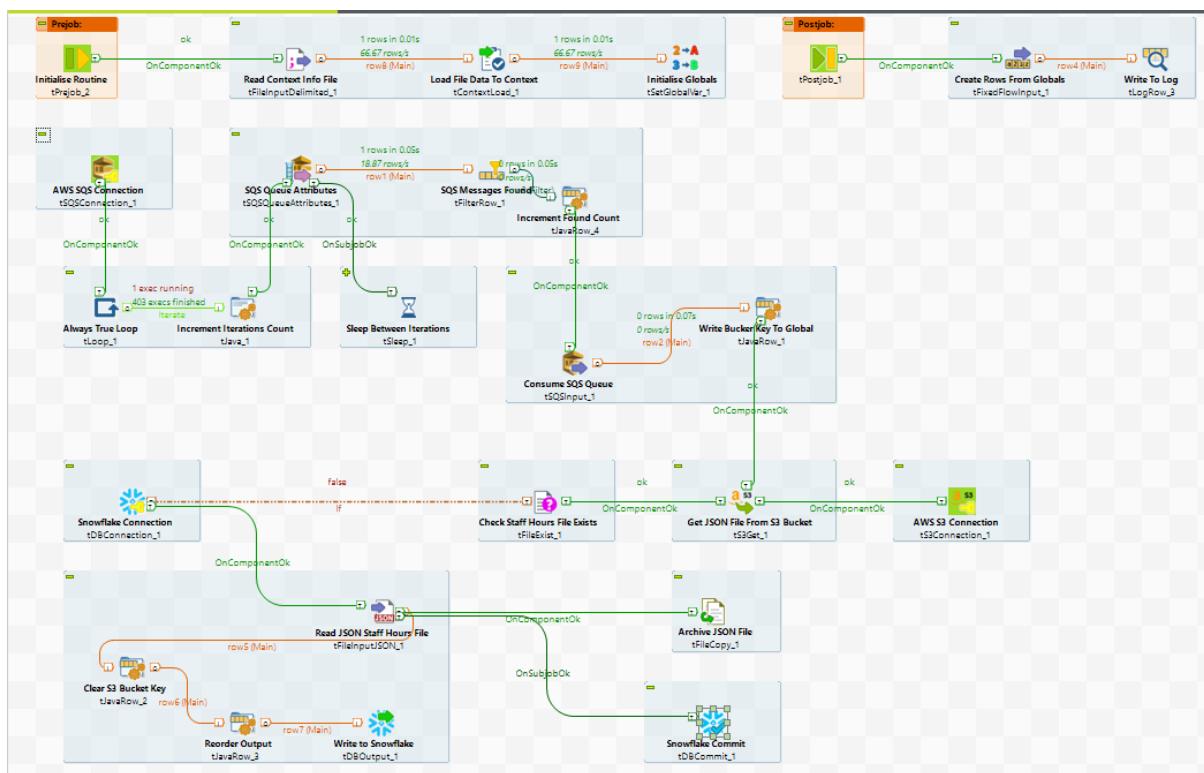
*As the commit operation is initiated by an OnSubjobOk trigger from the routine that includes the Snowflake output component, completion of the database write operation is guaranteed prior to the commit. It is therefore safe to close the connection as part of the commit.*

Label the component as "**Snowflake Commit**" and this completes the configuration of the components. For documentation purposes add the following **note** to the design area.

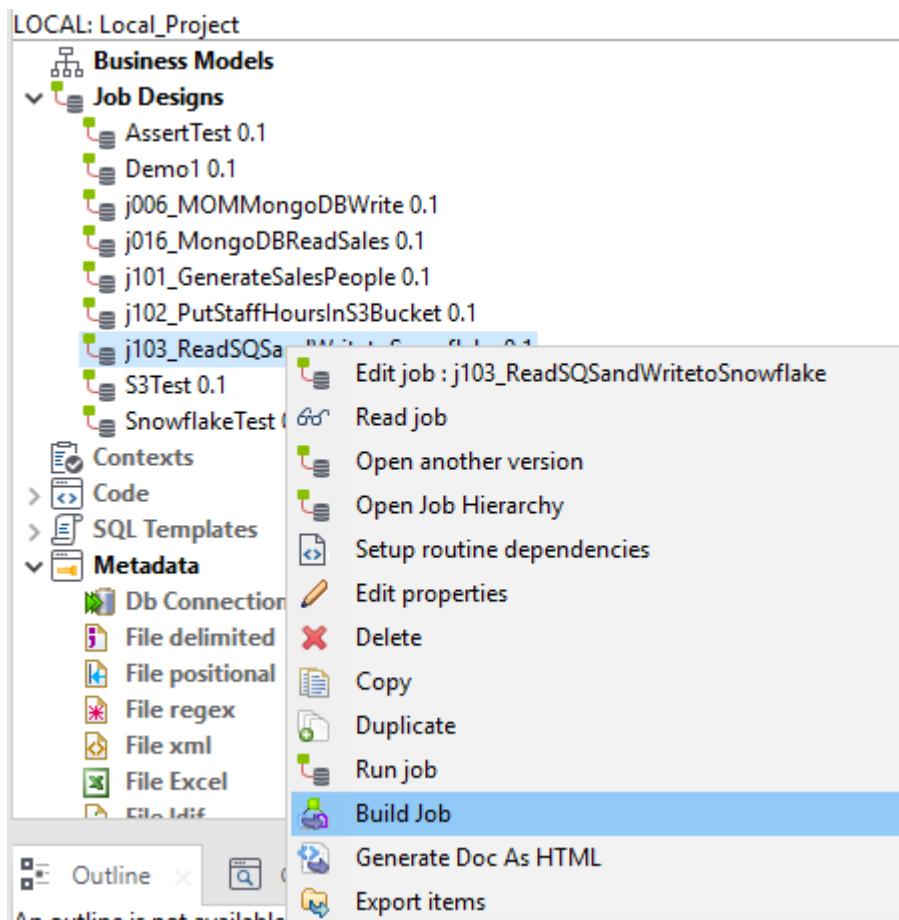
Job using infinite loop to keep monitoring AWS SQS queue.  
When message is found extracts S3bucket key from body.  
Uses key to extract JSON file from bucket then clears key  
 Parses JSON file and reorders rows into Snowflake format  
 Write data from JSON into snowflake table and archives local JSON file

The completed job should resemble the illustration below.

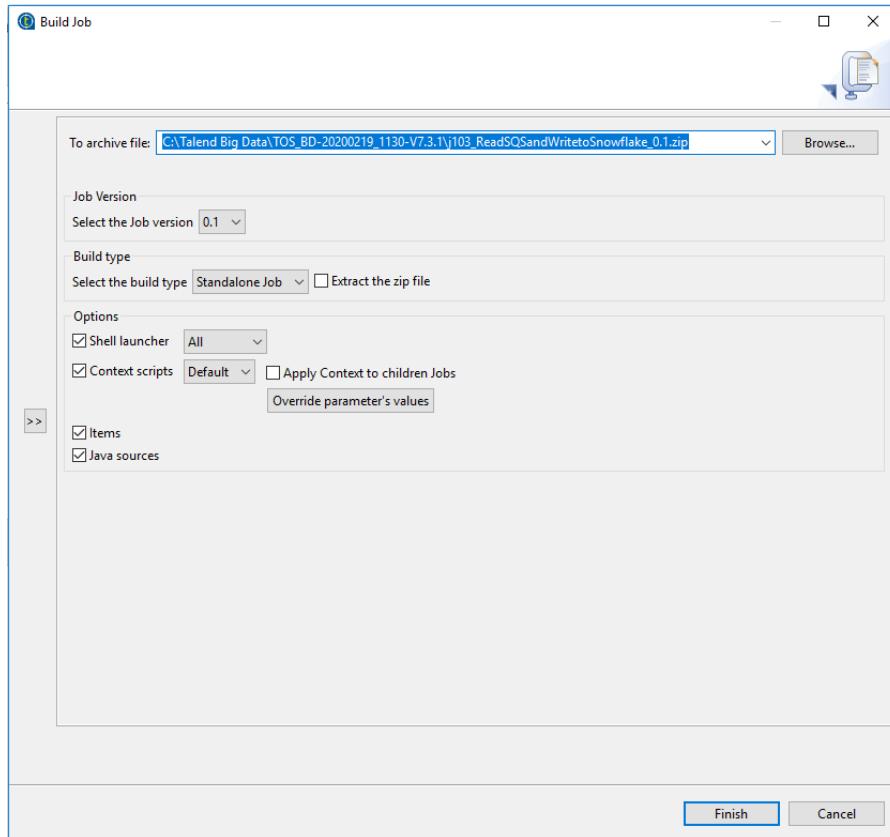
## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL



Build the job as **stand-alone** to run on the EC2 instance. Begin the build by saving the job then right click on its name under **Job Designs** in the repository pane and select **Build Job**.



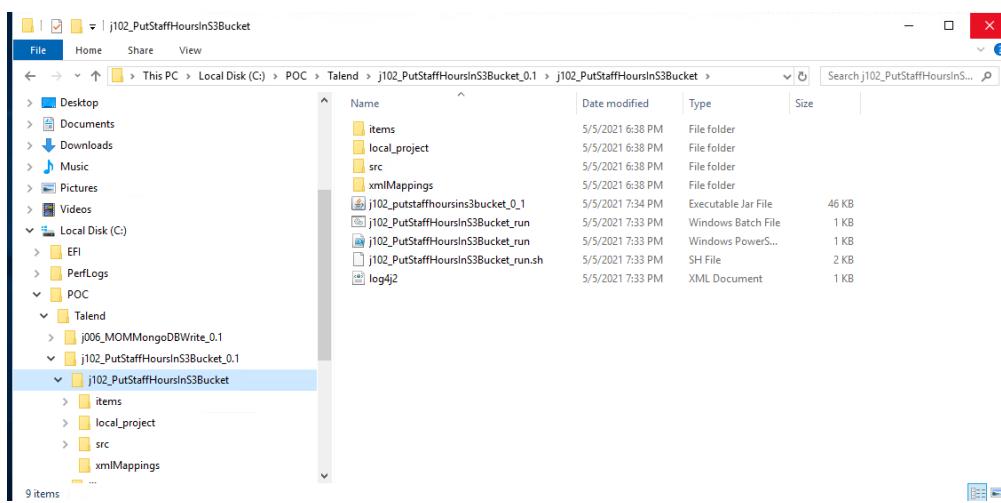
Accept all default values in the build screen and click **Finish** to complete.



Once complete a zip file will be created in the location specified. In this case that is the root folder of Talend Studio although that can be changed if required.

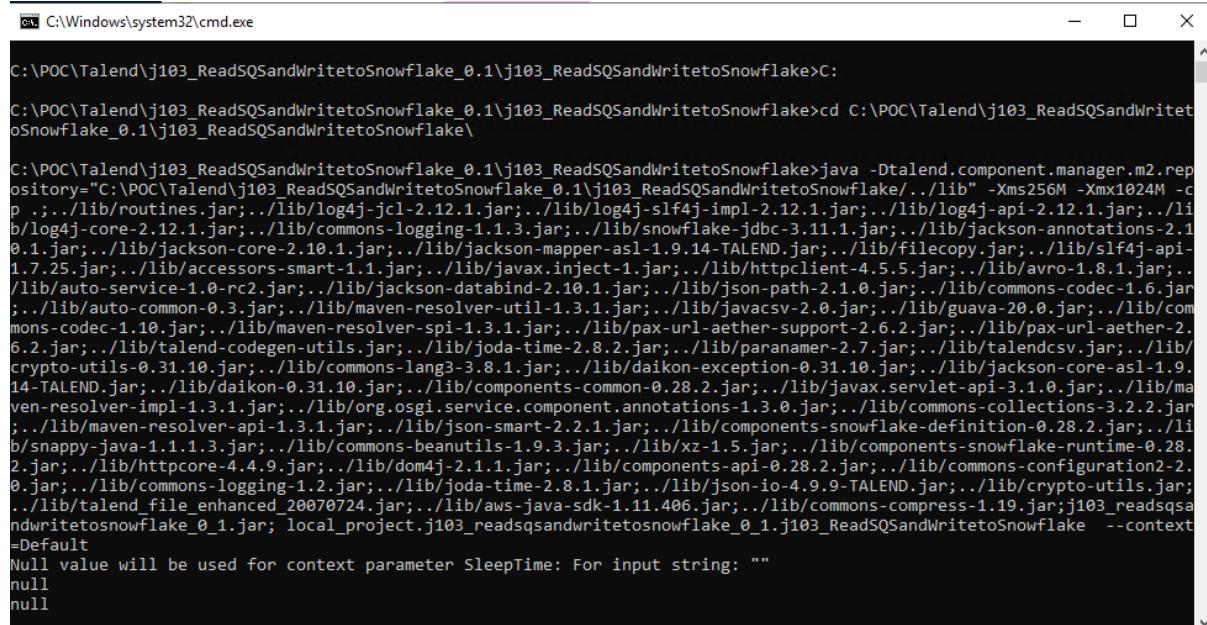
## 14. Deploy jobs to EC2

The deployment process is identical for both Talend jobs. Copy the entire zip files produced by the build process to the **EC2** server and extract the contents to the desired location. Once extracted look in the folder structures for the windows batch files **j102\_PutStaffHoursInS3Bucket\_run.bat** and **j103\_ReadSQSandWritetoSnowflake\_run.bat**.



Clicking on either batch file will run the respective job although they could also be running via **Task Scheduler** or if the subscription product was used, be submitted by **Talend Job Server**.

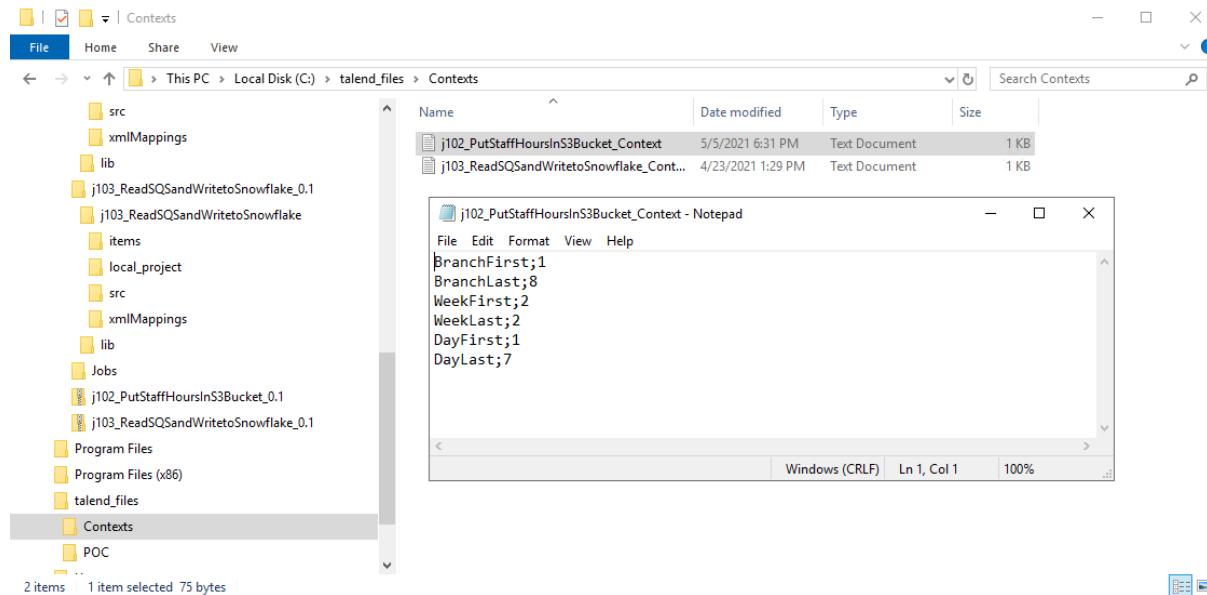
Store staff records have already been generated. To run the rest of the process, double-click the **j103\_ReadSQSandWritetoSnowflake\_run.bat** file which will start the SQS queue monitoring job in a loop. Once initiated this will remain active until manually terminated assuming no errors intervene.



```
C:\Windows\system32\cmd.exe
C:\POC\Talend\j103_ReadSQSandWritetoSnowflake_0.1\j103_ReadSQSandWritetoSnowflake>C:
C:\POC\Talend\j103_ReadSQSandWritetoSnowflake_0.1\j103_ReadSQSandWritetoSnowflake>cd C:\POC\Talend\j103_ReadSQSandWritetoSnowflake_0.1\j103_ReadSQSandWritetoSnowflake\

C:\POC\Talend\j103_ReadSQSandWritetoSnowflake_0.1\j103_ReadSQSandWritetoSnowflake>java -Dtalend.component.manager.m2.repository=C:\POC\Talend\j103_ReadSQSandWritetoSnowflake_0.1\j103_ReadSQSandWritetoSnowflake/../lib -Xms256M -Xmx1024M -cp ..;/lib/routines.jar;..;/lib/log4j-jcl-2.12.1.jar;..;/lib/log4j-slf4j-impl-2.12.1.jar;..;/lib/log4j-api-2.12.1.jar;..;/lib/log4j-core-2.12.1.jar;..;/lib/commons-logging-1.1.3.jar;..;/lib/snowflake-jdbc-3.11.1.jar;..;/lib/jackson-annotations-2.1.0.1.jar;..;/lib/jackson-core-2.10.1.jar;..;/lib/jackson-mapper-asl-1.9.14-TALEND.jar;..;/lib/filecopy.jar;..;/lib/slf4j-api-1.7.25.jar;..;/lib/accessors-smart-1.1.jar;..;/lib/javax.inject-1.jar;..;/lib/httpclient-4.5.5.jar;..;/lib/avro-1.8.1.jar;..;/lib/auto-service-1.0-rc2.jar;..;/lib/jackson-databind-2.10.1.jar;..;/lib/json-path-2.1.0.jar;..;/lib/commons-codec-1.6.jar;..;/lib/auto-common-0.3.jar;..;/lib/maven-resolver-util-1.3.1.jar;..;/lib/javacsv-2.0.jar;..;/lib/guava-20.0.jar;..;/lib/commons-codec-1.10.jar;..;/lib/maven-resolver-spi-1.3.1.jar;..;/lib/pax-url-aether-support-2.6.2.jar;..;/lib/pax-url-aether-2.6.2.jar;..;/lib/talend-codegen-utils.jar;..;/lib/joda-time-2.8.2.jar;..;/lib/paranamer-2.7.jar;..;/lib/talendcsv.jar;..;/lib/crypto-utils-0.31.10.jar;..;/lib/commons-lang3-3.8.1.jar;..;/lib/daikon-exception-0.31.10.jar;..;/lib/jackson-core-asl-1.9.14-TALEND.jar;..;/lib/daikon-0.31.10.jar;..;/lib/components-common-0.28.2.jar;..;/lib/javax.servlet-api-3.1.0.jar;..;/lib/maven-resolver-impl-1.3.1.jar;..;/lib/org.osgi.service.component.annotations-1.3.0.jar;..;/lib/commons-collections-3.2.2.jar;..;/lib/maven-resolver-api-1.3.1.jar;..;/lib/json-smart-2.2.1.jar;..;/lib/components-snowflake-definition-0.28.2.jar;..;/lib/snappy-java-1.1.3.jar;..;/lib/commons-beanutils-1.9.3.jar;..;/lib/xz-1.5.jar;..;/lib/components-snowflake-runtime-0.28.2.jar;..;/lib/httpcore-4.4.9.jar;..;/lib/dom4j-2.1.1.jar;..;/lib/components-api-0.28.2.jar;..;/lib/commons-configuration2-2.0.jar;..;/lib/commons-logging-1.2.jar;..;/lib/joda-time-2.8.1.jar;..;/lib/json-io-4.9.9-TALEND.jar;..;/lib/crypto-utils.jar;..;/lib/talend_file_enhanced_20070724.jar;..;/lib/aws-java-sdk-1.11.406.jar;..;/lib/commons-compress-1.19.jar;j103_readsqsandwritetosnowflake_0.1.jar; local_project.j103_readsqsandwritetosnowflake_0.1.j103_ReadSQSandWritetoSnowflake --context=Default
Null value will be used for context parameter SleepTime: For input string: ""
null
null
```

Check the context file for the Talend job to generate time sheet data.



Timesheets for all stores for week 2 should be generated. Click the **j102\_PutStaffHoursInS3Bucket\_run.bat** file to initiate the process which will launch a separate command window which closes on completion.

The timesheet data is generated in Talend and sent as a **JSON** format file to an **S3** bucket. Arrival in the bucket triggers an **AWS Lambda** function which writes the file name to an **SQS** message. The

## TALEND CONNECT TO AWS AND SNOWFLAKE TUTORIAL

other Talend job is monitoring the **SQS** queue and on receipt of a new message uses it to retrieve the **JSON** file from **S3**. The data is then parsed and written to **SnowflakeDB**.

The screenshot shows the Snowflake Worksheet interface. On the left, there's a sidebar with a tree view of databases (DEMO\_DB, RETAIL\_POC), schemas (INFORMATION\_SCHEMA, PUBLIC), and tables (TBPRODUCT, TBSTAFFHOURS). The main area has a toolbar with icons for Databases, Shares, Data Marketplace, Warehouses, Worksheets (which is selected), and History. Below the toolbar, a message says "Enjoy your free trial! Visit our documentation to learn more about using Snowflake or contact our support team with any questions." A user profile for "JOHTUCKER1961" (SYSADMIN) is at the top right. The central workspace shows a query editor with the following SQL code:

```
1 SELECT * FROM "RETAIL_POC"."PUBLIC"."TBSTAFFHOURS"
2 WHERE WEEK_NO = 2
```

Below the query editor is a "Results" section titled "Data Preview". It shows a table with 560 rows. The columns are Row, STORE\_CODE, DAY\_NO, HOURS\_WORKED, EMPLOYEE\_NO, LAST\_NAME, WEEK\_NO, and FIRST\_NAME. The data preview table is as follows:

Row	STORE_CODE	DAY_NO	HOURS_WORKED	EMPLOYEE_NO	LAST_NAME	WEEK_NO	FIRST_NAME
1	1	1	6	16	Johnson	2	Woodrow
2	1	1	6	7	Grant	2	Grover
3	1	1	6	13	Harding	2	Rutherford
4	1	1	6	10	Arthur	2	Lyndon
5	1	1	6	2	Washington	2	Theodore
6	1	1	6	17	Quincy	2	Thomas
7	1	1	6	15	Kennedy	2	Benjamin
8	1	1	6	19	Harrison	2	Ronald
9	1	1	6	3	Washington	2	Lyndon
10	1	1	6	4	Tyler	2	Warren
11	1	2	6	19	Harrison	2	Ronald

That completes this tutorial. In the real world the data would have come from multiple sources rather than being simulated by Talend but the purpose of this proof of concept was to demonstrate the connectivity capabilities of the product using modern cloud-based technologies.