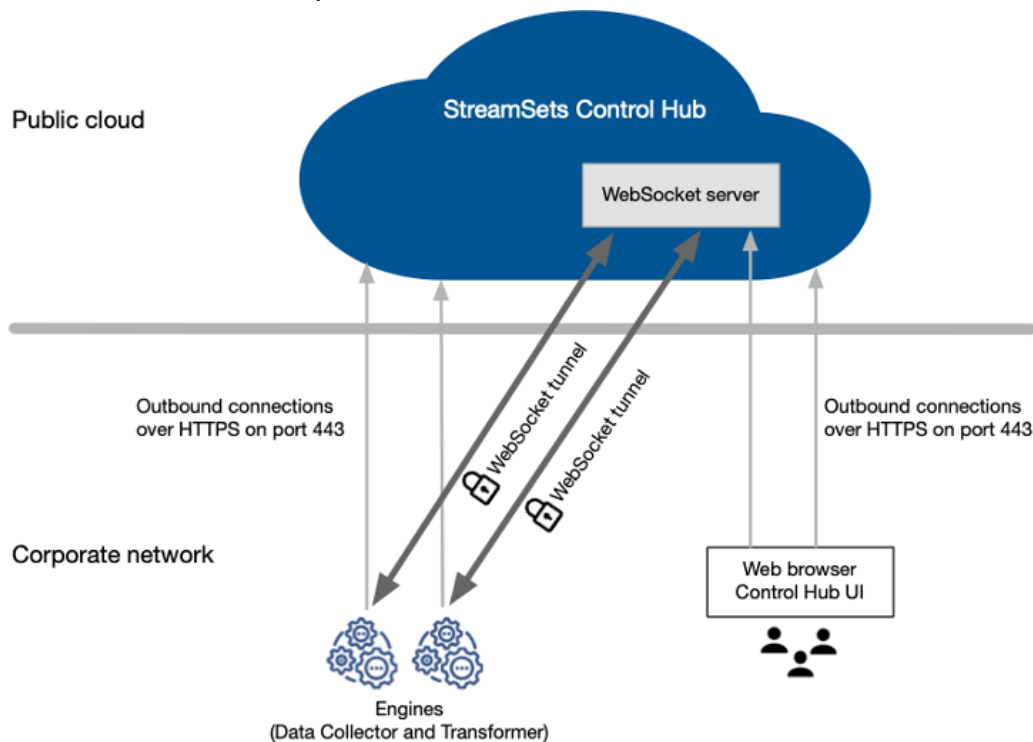


StreamSets Jumpstart

An in-depth guide into starting StreamSets, deploying engines, and running pipelines.

Overall architecture recap

The control plane of Streamsets Platform, Control Hub, hosted in the cloud. The data plane encompasses Streamsets Data Collectors (SDCs), which should be installed in a customer's network, be it in the cloud or on-premise.



Consider these options for SDC installation (not all of them are recommended for production scenario):

	Tarball (binary download)	Docker container
Windows PC	not supported	OK – the fastest to implement
Linux/Mac	OK with configured prerequisites	

Depending on your choice, find below instructions for major steps of tarball and Docker options.

Provide and prepare your machine

Use your own laptop or create a VM in the cloud of your choice. VM size should be determined by the minimum requirements for SDC installation: 2 vCores, 1GB RAM, 6GB disk space.

If on Linux, run a system update command, like *yum update* or *apt-get update*.

If your machine requires an HTTP proxy to access the Internet, make note of the proxy details, you will need them.

Docker option

Instal Docker Desktop on your Mac or Windows – <https://docs.docker.com/desktop/>

Or install on Install Docker Engine on a Linux VM – <https://docs.docker.com/engine/install/>

Tarball option

1. If not already installed, install Java 8, 11 or 17. Installation details differ between Linux versions, not mentioning MacOS, and it would be impractical to copy them here.

2. Configure the Open File Limit.

In **Linux**, it means adding a line like the following to the file */etc/security/limits.conf* (you need local admin rights for that):

```
*               -   nofile          32768
```

(you will need to log off and log on for the change to have effect)

In **MacOS**, the process has few more steps and is best covered by this link:

https://docs.streamsets.com/portal/platform-datacollector/latest/datacollector/UserGuide/Installation/Requirements.html#concept_al3_qz5_jz

Configure a Streamsets Platform Tenant

Sign-up for Streamsets Free Trial

On any computer (it doesn't have to be the machine for SDC installation), proceed to <https://cloud.login.streamsets.com/signup> and complete sign-up.

Confirm your registration via the link sent to you by email, then log on to the Platform. On your first logon you will need to specify:

- the region of the Platform instance: there are instances in North America, Europe and AsiaPacific, but only the former two available for Free Trial;
- Your Company name (you can change it later if needed);
- Agree to Terms and Conditions.

Create a new Deployment

Menu **Set Up** -> **Deployments** -> Click on “**Create a deployment**” link.

Enter the new Deployment’s name.

New Deployment

1 Define Deployment

Define deployment details. Once saved, you cannot change the deployment type, the engine version, or the environment that the deployment belongs to. [Learn more](#)

Deployment Name:	<input type="text" value="Roman's First Deployment"/>	?
Deployment Type:	<input type="text" value="Self-Managed"/>	?
Environment:	<input type="text" value="Default Self-Managed Environment (SELF)"/>	?
Engine Type:	<div><input checked="" type="radio"/> Data Collector - Runs data ingestion pipelines that perform record-based data transformations in streaming, CDC, or batch modes</div> <div><input type="radio"/> Transformer - Runs data processing pipelines on Spark that perform set-based transformations such as joins, aggregates, and sorts on the entire data set</div>	
Engine Version:	<input type="text" value="5.6.1"/>	?
Deployment Tags:	<input type="text" value="Add New..."/>	?
<div><input type="button" value="Cancel"/> <input type="button" value="Save & Next"/></div>		

Leave everything else default, for example, engine version may be higher than in the screenshot. Click **Save & Next**.

The **Configure Deployment** screen requires two actions: add required libraries, and configure proxy settings (this bit can be skipped if no proxy is used).

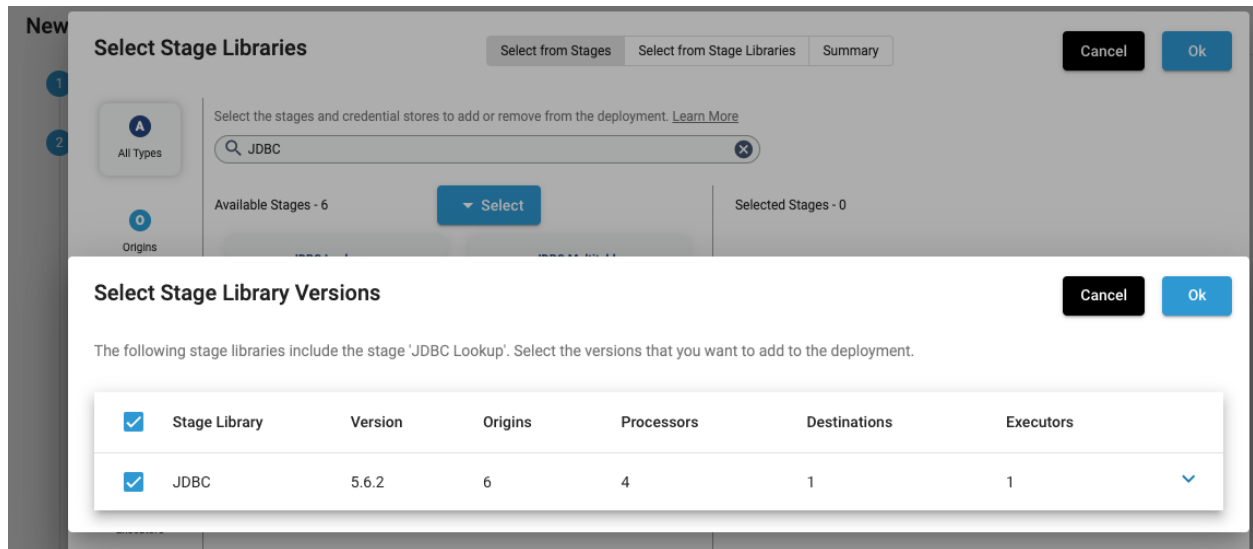
2 Configure Deployment

2a Configure Engine

Define the configuration of the engine to deploy. You can use the defaults to get started. [Learn more](#)

Stage Libraries:	<u>3 stage libraries selected</u>	?
Advanced Configuration:	<u>Click here to configure</u>	?

To add libraries, click on the link next to **Stage Libraries**. Then search for JDBC (it will be required when building a pipeline) and click the big “+” on any stage in the resulting list – all of them use the same library:



Click **Ok** to add this library, then click **Ok** to exit to the deployment wizard. If your SDC is not behind a proxy, just leave everything else default and click **Save & Next** again.

If your SDC is behind a proxy, click the link next to **Advanced Configuration** and go to the **Proxy** tab like in the screenshots below.

Engine Configuration

Data Collector Configuration Credential Stores Log4j2 **Proxy** Security Policy Java Configuration

```
1 # Copyright (c) 2022 StreamSets Inc.
2 #
3 # THESE SETTINGS ARE USED ONLY BY DATAOPS PLATFORM.
4 #
5 # CHANGING THESE PROPERTIES MAY REQUIRE THAT YOU RESTART ALL ENGINES FOR THIS DEPLOYMENT.
6 #
7 # If the Data Collector engine will be behind a proxy, the proxy settings
8 # must be included here for the engine to be able to communicate
9 # with Control Hub.
10
11 # Set the proxy host to the proxy machine IP address, e.g. http.proxyHost=172.31.7.179
12 # Set the proxy host for HTTPS, e.g. https.proxyHost=172.31.7.179
13 # Set the port to use to communicate with the proxy host, e.g. http.proxyPort=3128
14 # Set the port for HTTPS communication with the proxy host
15 # Set the list of addresses for which the proxy should not be used, separated by the | (pipe) character.
16 # In addition the wildcard character '*' can be used for pattern matching.
17 # For example http.nonProxyHosts=.foo.com|192.168.12.*|localhost will indicate that every host in the
18 # foo.com domain, every IP address matching the pattern 192.168.12.* and the localhost should all be
19 # accessed directly even if a proxy server is specified.
20 # Don't use quotes around usernames or passwords unless you intend those to be part of the string properties themselves.
21
22 http.proxyHost=
23 http.proxyPort=
24 http.proxyUser=
25 http.proxyPassword=
26 https.proxyHost=
27 https.proxyPort=
28 https.proxyUser=
29 https.proxyPassword=
30 http.nonProxyHosts=
```

Fill in the required details for accessing your proxy.

In the next screen, select your **Install Type** for your data collector engine. Click **Start & Next** for the next screen, then click **Start & Generate Install Script** on the final one. Copy the install script into your clipboard.

Install Data Collector

Docker option

The generated script is just a command to run a Docker container:

```
docker run -d -e http_proxy= -e https_proxy= -e
STREAMSETS_DEPLOYMENT_SCH_URL=https://eu01.hub.streamsets.com
-e STREAMSETS_DEPLOYMENT_ID=<some random string> -e
STREAMSETS_DEPLOYMENT_TOKEN=<really long random string> -e
ENGINE_SHUTDOWN_TIMEOUT=10 streamsets/datacollector:5.6.1
```

You may make some changes to this command, for example:

- add “sudo” in the beginning, unless the Docker engine is configured to work with non-root accounts;
- container name, e.g. “--name sdc561”;

– [volume mapping](#), if local file system is going to be used.

Run this command in the terminal/command window on your VM or laptop.

Tarball option

The generated script will look like this (but if you had configured a proxy the script will include all those details):

```
http_proxy= https_proxy= bash -c 'set -eo pipefail; curl -fsS
https://eu01.hub.streamsets.com/streamsets-engine-install.sh |
bash -s -- --deployment-id="<some random string>"
--deployment-token="<really long random string>"
--sch-url="https://eu01.hub.streamsets.com"
--engine-shutdown-timeout="10" --foreground '
```

Run this command in the terminal window for your VM or Mac. The easiest/fastest (and not production-recommended) option is to run this script as your own regular user and accept all defaults that the script suggests.

Verify Data Collector availability

At the end of either Docker or Tarball install, you should have your SDC application running. 1-2 minutes after the SDC starts up, you should also be able to see it in the list of Engines in the Platform UI:

The screenshot shows the StreamSets Platform UI. On the left is a sidebar with navigation links: 'Set Up', 'Environments', 'Deployments', 'Engines' (highlighted in blue), 'Connections', 'Build', 'Run', and 'Monitor'. The main content area is titled 'What are Data Collectors?' and includes a description: 'A Data Collector engine runs data ingestion pipelines that perform record-base'. Below this is a 'Learn More' link. Further down, under the heading 'Data Collector Engines (1)', there is a list of engines. The first engine is 'Engine' with a checkbox. The second engine is 'Roman's First Deployment (Self-Managed) e5f287fcbf93:18630' with a checkbox.

Develop your first pipelines

The use case covered below is scraping semi-structured data from the Web to a relational storage.

To follow the “datalake zones” approach, we will scrape semi-structured data into the raw “zone”, and then process it to extract select information.

In the Platform UI, open menu Build -> Pipelines.

Click on **Create Pipeline** link.

New Pipeline

1

Define Pipeline

Define the pipeline name, the type of engine for the pipeline, and whether to start with a blank canvas or with a sample pipeline.
[Learn more](#)

Name:

?

Description:

?

Engine Type:

☒ **Data Collector** - Runs data ingestion pipelines that perform record-based data transformations in streaming, CDC, or batch modes

☐ **Transformer** - Runs data processing pipelines on Spark that perform set-based transformations such as joins, aggregates, and sorts on the entire data set

☐ **Transformer for Snowflake** - Runs data processing pipelines on Snowflake

Start with:

☒ **Blank Pipeline**

☐ **Sample Pipeline**

Cancel

Next

Maintain the pipeline’s Name, then click **Next**.

2 Configure Pipeline

If starting with a sample pipeline, select the sample to use. Select the authoring engine to use for pipeline design. The engine with the most recent reported time is selected by default. [Learn more](#)

Authoring Engine: Roman's First Deployment (Self-Managed) - e5f287fcbf93:18630



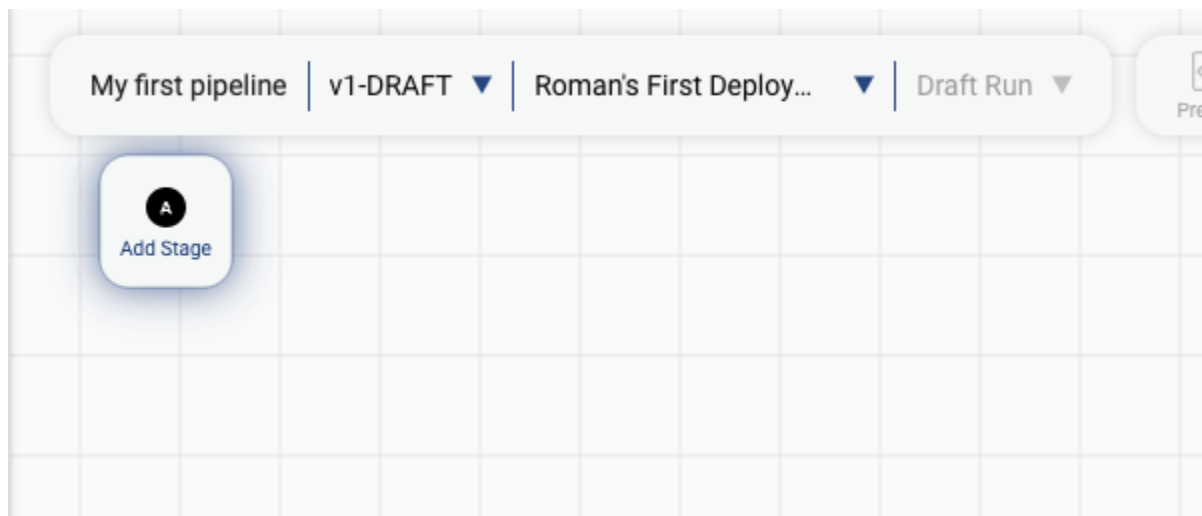
[Click here to select](#)

Back

Save & Next

Save & Open in Canvas

Select the SDC you've just installed as your Authoring Engine, and click **Save & Open in Canvas**.



In the canvas, click on **Add Stage** and find **HTTP Client** origin (you can use search). The origin's configuration will open automatically.

General	HTTP	Pagination	Credentials	OAuth 2	TLS	Data Format
Resource URL	<code>https://api.open-meteo.com/v1/forecast?latitude=-33.82&longitude=150.82&current_weather=true&hourly=temperature_2m,relativehumidity_2m,windspeed_10m</code>					
Mode	Batch					
HTTP Method	GET					
Authentication Type	None					
Use OAuth 2	<input type="checkbox"/>					
Headers	+ ADD BULK EDIT MODE					
Additional Security Headers	+ ADD					
Show Advanced Options						

Maintain configuration:

- **Resource URL** – this example is using free weather data from Open Meteo (the coordinates are of Sydney, Australia, but feel free to use your location):

https://api.open-meteo.com/v1/forecast?latitude=-33.82&longitude=150.82¤t_weather=true&hourly=temperature_2m,relativehumidity_2m,windspeed_10m

- **Mode** (Batch).

Set the expected incoming **Data Format** to JSON format, and **Max Object Length** (max size of a single record) to 10000 characters:

The screenshot shows the 'Data Format' configuration tab. It includes a 'Data Format' dropdown set to 'JSON', a 'JSON Content' dropdown set to 'Multiple JSON objects', a 'Compression Format' dropdown set to 'None', a 'Max Object Length (chars)' input field set to '10000', a 'Charset' dropdown set to 'UTF-8', and an 'Ignore Control Characters' checkbox which is unchecked. There are also links for 'Hide Advanced Options' and an information icon.

Click on **Add Stage** and select **Local FS** destination. Set the target location to `/tmp/raw_weather`:

The screenshot shows the 'Output Files' configuration tab. It includes a 'Directory Template' input field set to `/tmp/raw_weather`, a 'File Type' dropdown set to 'Text files', and a 'Files Prefix' input field set to `sdc-${sdc:id() }`. There are also links for 'Hide Advanced Options' and an information icon.

... and specify the output Data Format (choose JSON).

Execute the pipeline, it should finish within a minute, processing 1 record and generating a file in `/tmp/raw_weather` directory (if your SDC runs in a Docker container, the file will be created in that container, too).

Create a second pipeline, and use Directory origin, use `/tmp/raw_weather` in **Files Directory** configuration setting, `*` as files pattern, and Last Modified Timestamp for Read Order.

The screenshot shows a pipeline configuration interface. At the top, a visual diagram shows a 'Directory 1' origin connected to an 'Add Stage' action. Below this, a status bar indicates 'Advanced Options', 'Edit Mode', and 'All Changes Saved'. The 'Files' tab is selected, showing configuration fields: 'Files Directory' (set to `/tmp/raw_weather`), 'File Name Pattern' (set to `*`), 'First File to Process' (set to `First File to Process`), 'Number of Threads' (set to `1`), 'File Name Pattern Mode' (set to `Glob`), and 'Read Order' (set to `Last Modified Timestamp`). A 'Hide Advanced Options' link is also visible.

Set **Data Format** for JSON again, and set max record size to 10000.

Let's have a look at the data. Click on **Preview** button in the toolbar, and select **Run Preview** action:

The screenshot shows the same pipeline configuration interface, but with a toolbar at the top. The toolbar includes a dropdown menu with 'Second pipeline', 'v1-DRAFT', and 'Roman's First Deploy...'. To the right of the dropdown are two buttons: 'Preview' (with an eye icon) and 'Validate' (with a checkmark icon). Below the toolbar, the visual diagram of the pipeline is shown again.

After a few seconds, you will see data preview in the lower half of the screen. The data is nested:

► OUTPUT

▼ Record1 : {MAP}

latitude : {DOUBLE} -33.875  
longitude : {DOUBLE} 150.875  
generationtime_ms : {DOUBLE} 0.08404254913330078  
utc_offset_seconds : {INTEGER} 0  
timezone : {STRING} "GMT"  
timezone_abbreviation : {STRING} "GMT"  
 elevation : {DOUBLE} 60.0  
► current_weather : {MAP} 
► hourly_units : {MAP} 
▼ hourly : {MAP} 
► time : {LIST[168]} 
► temperature_2m : {LIST[168]} 
► relativehumidity_2m : {LIST[168]} 
► windspeed_10m : {LIST[168]} 

▼ Event Record1 (new-file) : {MAP}

filepath : {STRING} "/tmp/raw_weather/sdc-b467102b-42c8-4988-a814-2596fc429a56_622c164d-9164-4020-ae57-c81838918a83" 

► Event Record2 (finished-file) : {MAP}

Expand the *hourly* node. It contains several lists – columns of values, where the first values from each list would form the first record, etc. If only there was a way to zip those values together..

Add a new stage **Field Zip**, specify */hourly/time* as the first field to zip, and */hourly/temperature_2m* as the second field, set the zipped field to */zipped*:

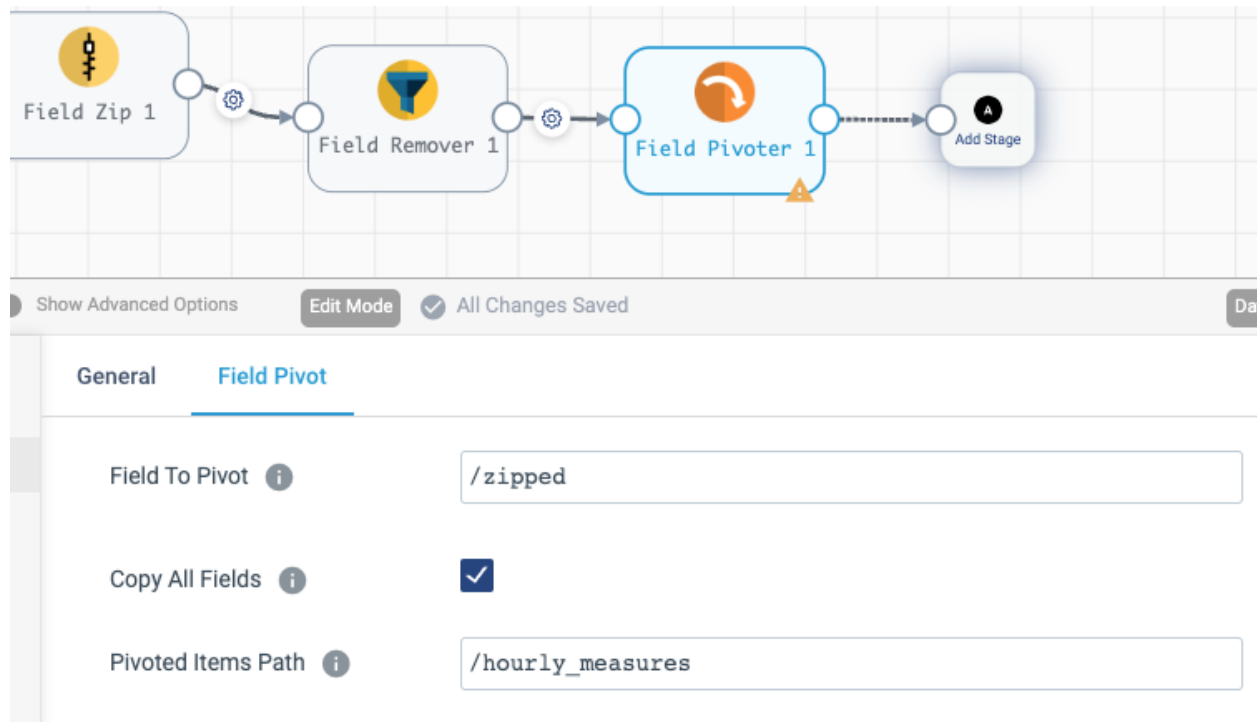
The screenshot shows the 'Field Zip' stage configuration in the Data Collector Pipeline. The pipeline consists of three stages: 'Directory 1', 'Field Zip 1', and 'Add Stage'. The 'Field Zip 1' stage is highlighted with a blue border and a warning icon. Below the pipeline view, the 'General' tab is selected, and the 'Zip' sub-tab is active. The configuration includes three input fields: 'First Field' with the value '/hourly/time', 'Second Field' with the value '/hourly/temperature_2m', and 'Path for Zipped Field' with the value '/zipped'. There are buttons for '+ ADD ANOTHER' and 'BULK EDIT MODE'. Below these, there is a checkbox for 'Zip Values Only' which is currently unchecked, and a dropdown menu for 'Field Does Not Exist' set to 'Send to Error'.

You may want to run **Preview** again to check how resulting data looks. To zip in more fields, you need to chain more Field Zip processors, but at the moment let's just proceed with only two.

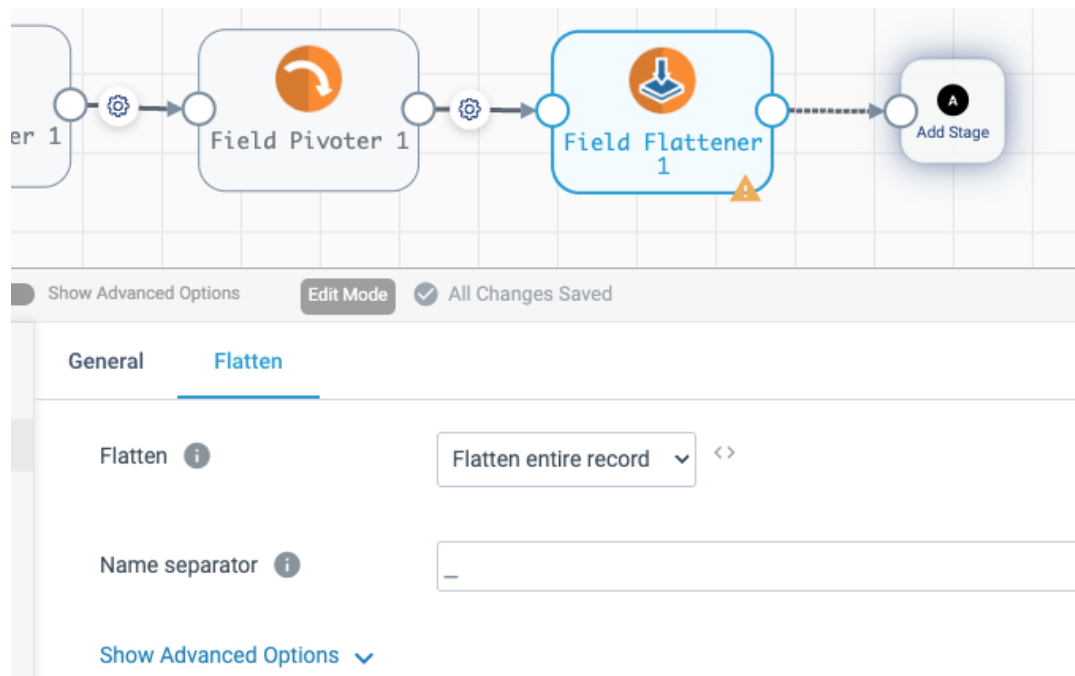
As we don't need the `/hourly` field anymore, let's drop it using **Field Remover** stage. Getting rid of unneeded fields early helps to keep the memory footprint of the pipeline low.

The screenshot shows the 'Field Remover' stage configuration in the Data Collector Pipeline. The pipeline now consists of four stages: 'Directory 1', 'Field Zip 1', 'Field Remover 1', and 'Add Stage'. The 'Field Remover 1' stage is highlighted with a blue border and a warning icon. Below the pipeline view, the 'General' tab is selected, and the 'Remove/Keep' sub-tab is active. The configuration includes a dropdown menu for 'Action' set to 'Remove Listed Fields', and a text input field for 'Fields' containing the value '/hourly' with a delete button (X) and an 'Add New...' button.

The zipped list of values still resides in a single record. In a relational database model, each measurement should be a separate record – This is achieved with the **Field Pivoter** processor. The field to pivot is */zipped*, and let's put the pivoted result to */hourly_measures*.

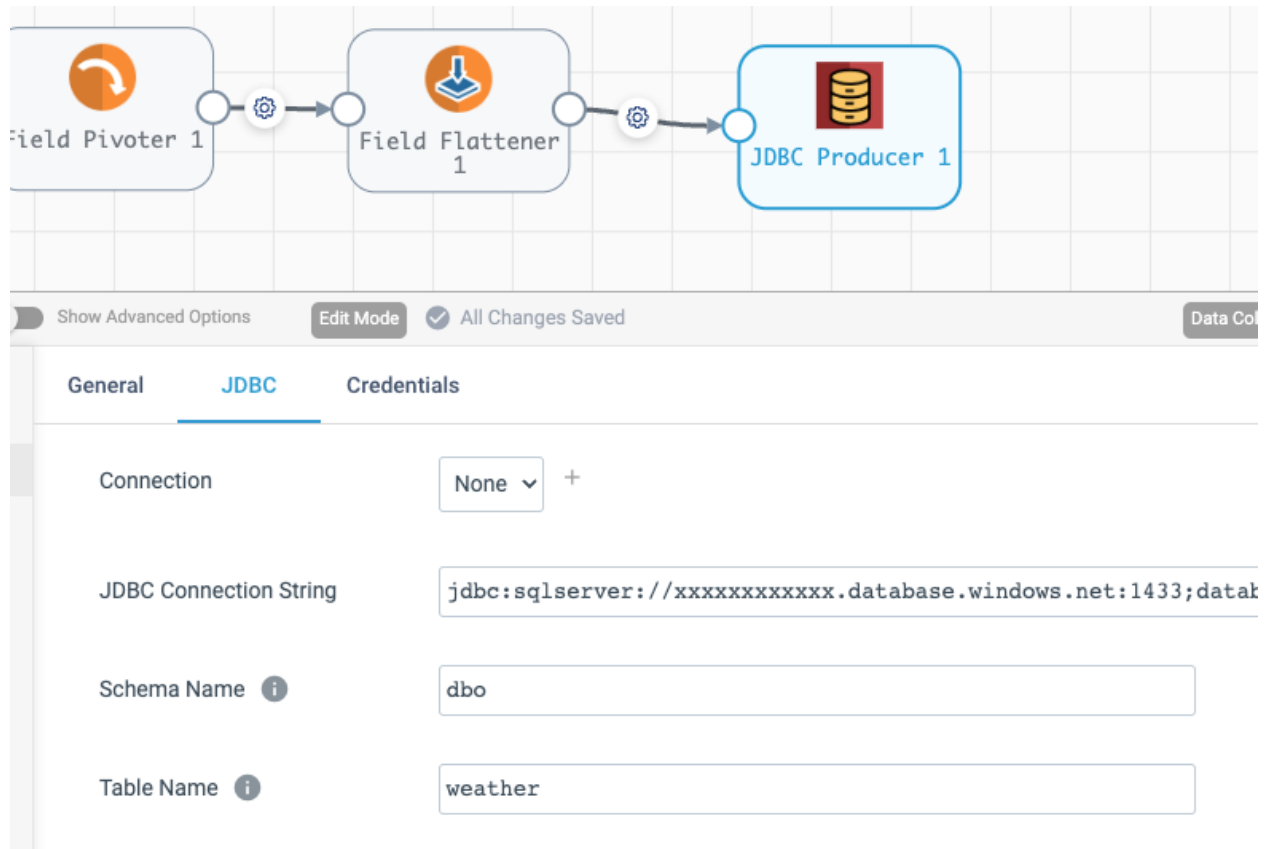


The last thing to do is convert nested fields to flat records suitable for loading into a relational database. Add **Field Flattener** processor and configure it to flatten the entire record. To simplify further processing, use underscore as the name separator, instead of the default dot.



At this stage, the data is already suitable for relational database load. One extra step that might be required is renaming fields with **Field Renamer** to match the target schema.

Finish your pipeline with a [destination](#) stage of your choice, e.g. **Local FS** to output flattened data to a directory on SDC machine – by now, you should be able to configure it without much guidance. Alternatively, to output to a SQL database, use **JDBC Producer** – the example below demonstrates how to configure it for Azure SQL Database.



That's it, your pipeline to process scraped nested data is ready. Before running it, make sure to create the target table in the database:

```
create table dbo.weather (  
  latitude real not null,  
  longitude real not null,  
  generationtime_ms real null,  
  utc_offset_seconds int null,  
  timezone nvarchar(3),  
  timezone_abbreviation nvarchar(3),  
  elevation real null,  
  current_weather_temperature real not null,  
  current_weather_windspeed real not null,  
  current_weather_winddirection real not null,  
  current_weather_weathercode int null,  
  current_weather_is_day int null,  
  current_weather_time nvarchar(16) null,  
  hourly_units_time nvarchar(10) null,  
  hourly_units_temperature_2m nvarchar(3) null,  
  hourly_units_relativehumidity_2m nvarchar(3) null,  
  hourly_units_windspeed_10m nvarchar(6) null,  
  hourly_measures_time nvarchar(16) not null,  
  hourly_measures_temperature_2m real not null  
);
```

Start the pipeline, and wait until it reports 168 processed records and pauses. Actually, it's not a pause, the pipeline is listening for new files to arrive, and you can test it by running the first (data scraping) pipeline in a new browser window and then returning to this one: from 168

processed records, the number will change to 336. You can manually stop the pipeline by clicking the **Stop** button.

Conclusion And Next Steps

Congratulations, you have created a new Streamsets Platform tenant, installed a Data Collector engine and used it to extract, transform and load weather data! If nothing unexpected happened, you should have done all this in less than an hour.

There are a couple of things here to reflect about.

This ETL process can be executed as a batch. For this, you need to create a [Job](#) for each pipeline, and then create an [orchestration pipeline](#) running transformation/load after extraction. To make the second pipeline stop after all available data is processed, a [Pipeline Finisher](#) would be required. This might be a good option if the source is updated infrequently, or subsequent analytics does not require most recent data.

Alternatively, this process could be used for more modern, continuous data delivery with decoupled architecture. In the first pipeline, use **Polling** mode of HTTP Client origin, to make it scrape data at regular intervals. Create and start Jobs for both pipelines – they are not tightly coupled, but the second pipeline would process a new file every time after the first one creates it. Meanwhile, the storage in the middle doesn't have to be file-based: these pipelines can be easily changed to use a queueing application like [Kafka](#) or [Amazon Kinesis](#)

Here some more ideas to explore, in no special order:

- Platform management features, like [version control](#) or [Python SDK](#),
- other data sources and destinations, like databases (batch and change data capture) or cloud storage/services,
- if Snowflake is the central data storage of choice, build enterprise data warehouse using pipeline created in [Transformer for Snowflake](#).