**ubiquis** CONSULTING

**HOME**   **OUR SERVICES**   **WHO WE ARE**   **SHOWCASE**   **BLOG**   **CONTACT US**

Home      Blog      Loading CSV files with PDI Metadata Injection
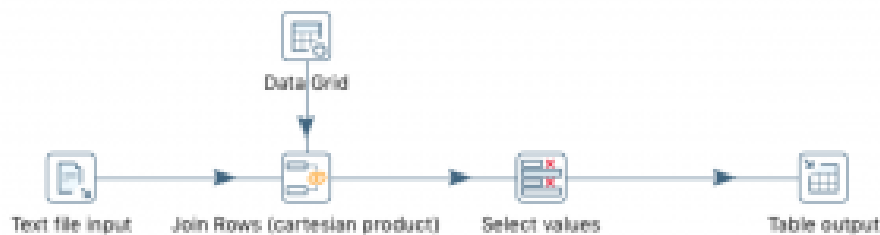
# Loading CSV files with PDI Metadata Injection



## Introduction

We often need to develop PDI jobs and transformations to load CSV files. In an ideal world, any data source will provide us with a single format for ingestion. All files look the same and a single PDI transformation can load them all into a database table.

But as is often the case, some files will have slight differences. Or not so slight ones. We've faced on several occasions the need to ingest files where the number of fields varies; or the order in which they appear; or their names.

If we only have a small number of formats to account for we can get by with two or three different transformations, one to load each format. But it may not be that simple.

If the number of formats is too large, or if there's no way to determine which format to load based on some criteria (filename, source, etc.), we need something smarter. And this is where metadata injection comes in.

## PDI Metadata Injection

PDI has supported Metadata Injection for a few years now. The basic idea consists of having a transformation parse metadata information in runtime, which may vary depending on the circumstances, and then injecting this metadata into the configuration of steps sitting in a separate "template" transformation. Depending on what the data looks like, the template transformation is injected with different values, changing parts of its behaviour. Here's a video to show you how it works, voiced by none other than Matt Casters himself.

Your mileage may vary, but PDI Metadata Injection is a very powerful ally when it comes to creating more generic algorithms.

missing fields (we'll need to add nulls for those)

extra fields (some fields are not required and should be ignored)

fields in a different order

must account for future change

must allow for easy maintenance

Here's a sample of the files we need to load:

```
1  id;name;age
2  1;john;10
3  2;jack;20
```

Some files have only 1 name
field

```
1  id;first;last;age
2  3;john;doe;30
3  4;robin;hood;50
4
```

Some have 2

```
1  id;last;first;age
2  3;federer;roger;36
3  4;nadal;rafael;31
4
```

Sometimes they're in a different order

Our goal is to parse all three files using a single PDI transformation with Metadata Injection, allow all configurations to come from an external data dictionary and support future changes to file formats as much as possible.

## The template transformation

The template transformation is quite simple:

a text file input step to read the CSV file;

a data grid step to add the fields that are missing from the input file

a select values step to re-order the fields to a fixed output order

**Remark:** not all PDI steps support Metadata Injection. In particular, the Add Constants and Generate Rows steps don't have support for it, so we can't use them to add the missing fields. The Data Grid step, however, supports Metadata Inection and can therefore be used, followed by a Join Rows, to add extra fields, which can vary from one file to the next. Also, the Select Values step isn't really necessary, as we can inject the list of output fields directly in the Table Output step. But this way we can also attach a Write to Log step right before the Table Output, if we need to debug our transformation.

## How does it work?

The Text file input step will list all fields that exist in the file, with the desired destination names; the Data grid step will include whatever fields are missing from the input file (e.g., there's no last name field in the first input file) and add a single row of data with value null; once joined to the main data flow, the select values step re-orders the fields in the order we want them and finally a table output step will write to the target database table.

If we take our data dictionary to include 4 fields, **id**, **first_name**, **last_name**, and **age**, the Text File Input step must have 3 fields configured, **id**, **first_name**, **age**; the data grid would have a single field, **last_name**, with a single row of data and null as value; the Select values step would then re-order the fields as **id**, **first_name**, **last_name**, **age**.
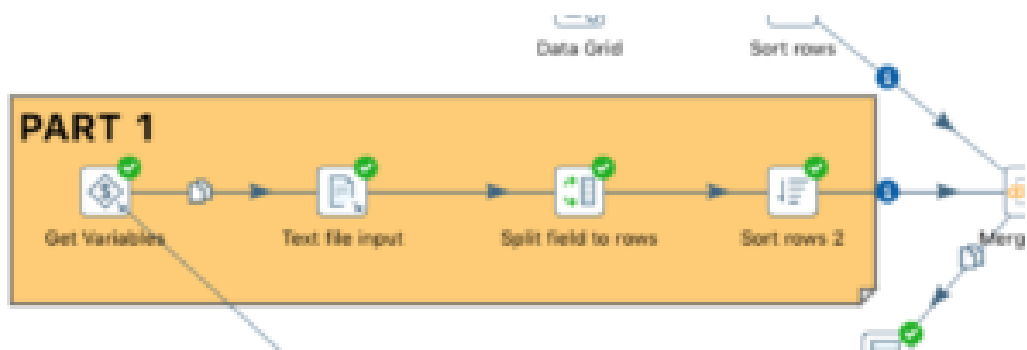
## The injector transformation

This transformation is the brains of the process. It must determine the file structure, fill in the gaps, and then inject the missing bits in the template transformation so it can load the files one by one.

Here's the structure of the transformation. We'll go over the various steps, one by one, and see what's happening at each stage of the process.
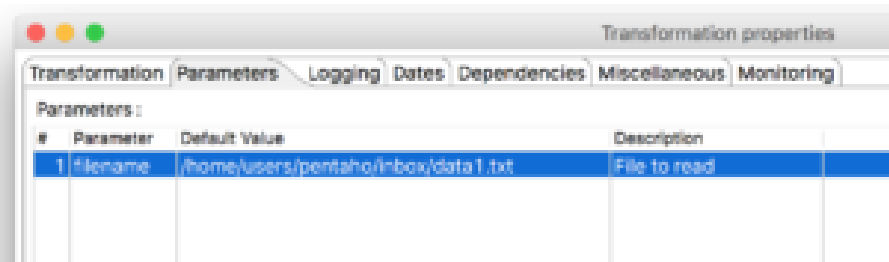
As often happens, injector transformations tend to have multiple streams of data, each carrying part of the set up to the metadata injection step.

## Parsing the list of fields



## Passing the filename

The filename is first set as a parameter, so that this transformation can then be called by a parent job, once for each of the files waiting to be ingested. The get variables step then retrieves this value.



This website uses cookies to improve your experience. We'll assume you're ok with this, but you can opt-out if you wish. Find out more    ACCEPT

## The list of fields

We then read the header row using a Text File Input step.

| | |
|---|---|
| Accept filenames from previous step | ☑ |
| Pass through fields from previous step | ☐ |
| Step to read filenames from | Get Variables |
| Field in the input to use as filename | filename |

Reading the filename from the previous step

| | |
|---|---|
| Filetype | CSV |
| Separator | $[00] |
| Enclosure | " |
| Allow breaks in enclosed fields? | ☐ |
| Escape | |
| Header | ☐  Number of header lines  1 |
| Footer | ☐  Number of footer lines  1 |
| Wrapped lines? | ☐  Number of times wrapped  1 |
| Paged layout (printout)? | ☐  Number of lines per page  80 |
| | Document header lines  0 |
| Compression | None |
| No empty rows | ☑ |
| Include filename in output? | ☐  Filename fieldname |
| Rownum in output? | ☐  Rownum fieldname |
| | Rownum by file?  ☐ |
| Format | mixed |
| Encoding | |
| Limit | 1 |
| Be lenient when parsing dates? | ☑ |
| The date format Locale | en_GB |

The contents tab

| File | Content | Error Handling | Filters | Fields | Additional output fields |
|---|---|---|---|---|---|

| # | Name | Type | Format | Position | Length | Precision | Currency |
|---|---|---|---|---|---|---|---|
| 1 | line | String | | | | | |

The fields tab

The important bits here:

> We use character **$[00]** as the separator, as it's a character that does not occur naturally in the data. This character is ASCII's character 00, or the null character. It should be safe under most uses, but your mileage may vary.

As we want to read the header row, we must disable the header option

Next we split use a Split field to Rows step, using semi-colon (the separator in our data files) to split the field:



Split field to rows

And we sort the rows of data by field name:



Sort fields alphabetically

This is our result. The **field_ordinal** field tells us the order in which the fields occur in the data file.



One row per field

The reason for that sort step will become clear in the next part.

## Preparing the data dictionary

We need a data dictionary, where the relationship between data fields coming from the files are mapped to the target fields on the destination database. Here we use a data grid, but on a production ETL this information would be stored in a database.

In this simple example we need the following information:



Mapping between source and target field names

We then sort by field name, ahead of joining it back to the field names coming from the file.

## Joining the source and target fields

The next step is joining the two streams of data, the field list coming from the data file and the field list coming from the dictionary.



The merge join step requires both input streams of data to be sorted by the join keys, which is why we needed to sort both data streams immediately before.

The need for a full outer join will be clear in a while.

This is the resulting data after the join:

We now start preparing the various data streams to be injected. The first one is the list of fields to inject in the template's Text File Input step, which will contain all the input fields present in the data file.

A Filter Rows step allows us to select only those rows that correspond to fields present in the data file.



Once sorted by the field_ordinal value (the order in which they appear in the data file, we have:



This list of fields exactly matches the order in which they occur in the data file.

## The list of target fields

The next step is to prepare the list of target fields. These will be injected in the Select Values and Table output steps.

We start by filtering only those fields that have a match in the data dictionary:



The output of this step will include all fields that found a match in the data dictionary:



Note that multiple occurrences of **first_name** are displayed in the output. This is because or our full outer join above. Though we need a full outer join to get both the fields present in the input data file as well as those missing, this then causes the duplicates. We get rid of them with a Group by step (actually, a Memory Group by that doesn't require a sort beforehand),

This website uses cookies to improve your experience. We'll assume you're ok with this, but you can opt-out if you wish. Find out

more     ACCEPT

This way we get a single row of output for each combination **target_field_name-target_field_ordinal**. The Aggregate fields are not used for this section but will be necessary for the next one, where we inject the Data Grid step to add the missing field information.
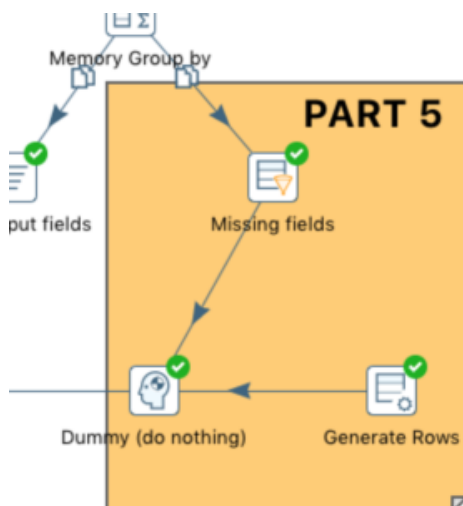
We finally sort by the **target_field_ordinal** field, which gives us



This is the list of output fields, in the order in which we want them. The order is determined by the ordinal field coming from the data dictionary, so it'll be the same order regardless of which file we're reading.

## Adding missing fields



In some cases we have files that are missing some of the fields. We can identify those rows, since they will have a null **field_name** after the Group by of the previous section (that's why we chose the "First non-null value" aggregation; **target fields** with non-null **field names** exist in the data file. **Target fields** with only null **field names** do not exist in the data file.

So, we first filter out only those rows for which **field_name** is null:

Which gives this result:



And we'll use this information (only **target_field_name** and **field_type**) to inject into the Data Grid step.

However, since PDI 7.1, if nothing is injected we get an error. If there's nothing to inject, then Data Grid returns 0 rows. When we then cross-join them in the injected transformation the result is 0 rows of data. This behaviour is not observed in a non-injected data grid: even an empty Data Grid step will always output 1 row of data (with no fields and no values). But on Metadata Injection it will output 0 rows.

Because of this we need to always inject at least one row of data in the Data Grid. That's what the Generate Rows step is doing:



With this fix we ensure there's always at least 1 row of data arriving at the Metadata Injection step from this data stream.

## Metadata Injection

At this point we have 3 incoming data streams. We add a new one, directly from the Get Variables step (this will inject the filename), and we're ready to go.

Here is the configuragion of the Metadata Injection step:

**ubiquis** CONSULTING

HOME    OUR SERVICES    WHO WE ARE    SHOWCASE    BLOG    CONTACT US



Metadata injection: Select Values

In the Select Values step we inject all the **target_field_names** coming from above, regardless of them existing in the data file or not (as missing fields are added in the Data Grid step).



Metadata injection: Text File Input

On the text file input step we inject the filename coming from the Get Variables step right at the beginning, and the field name and data type coming from the dictionary (we use the **target_field_name** directly and that way avoid renaming the field later).

And that's it (almost!). Lets test it.

## Running the transformation

We first set the filename parameter to point to our first file, **data1.txt**, and run the transformation. Once successful, we can go and see the data on the target database:

## UBIQUIS CONSULTING

HOME    OUR SERVICES    WHO WE ARE    SHOWCASE    BLOG    CONTACT US

```
1|john||10
2|jack||20
3|john|doe|30
4|robin|hood|50
sqlite>
```

And finally, after changing it again to point to file **data2a.txt**, we have:

```
sqlite> select * from users;
1|john||10
2|jack||20
3|john|doe|30
4|robin|hood|50
3|roger|federer|36
4|rafael|nadal|31
sqlite>
```

All 3 files loaded correctly, even though they had different structures of data. Our work is done.

## Future proofing

Or is it?

We should test whether our transformation can cope with new fields showing up in the files; fields that the data dictionary isn't yet aware of.

So, let's load one more file, this time with an extra field (and to make it interesting, we'll just add the new field in the beginning):

```
1  active;id;first;last;age;
2  Y;3;john;doe;30
3  N;4;robin;hood;50
4
```

If we switch the transformation to the new file, it'll run successfully and this is what the data shows:

```
sqlite> select * from users;
1|john||10
2|jack||20
3|john|doe|30
4|robin|hood|50
3|roger|federer|36
4|rafael|nadal|31
3|john|doe|30
4|robin|hood|50
sqlite>
```

Indeed, two more rows of data were inserted, and the right fields were written to the right columns. Doesn't have the new column, obviously, as the ETL isn't yet aware of it.

In order to load the new column, we need only two things:

> add a new column to the target table

```
sqlite> alter table users add column active string;
```

**HOME**    **OUR SERVICES**    **WHO WE ARE**    **SHOWCASE**    **BLOG**    **CONTACT US**



Now we can re-run that last file and see the new column populated:



## Conclusion

With this transformation we are able to ingest CSV files with variable schema, and add more columns as required. All the data dictionary information can be put in an external database, to make maintenance easier. This allows our ETL to accommodate changes in the data format as they occur.

Of course this is not the end of it. We've only scratched the surface in terms of detecting and parsing the various changes in format. Files can have different separators, encodings, and even date formats, which need to be accounted for. But it's a start, and allows us to, at the cost of spending a bit more time in the initial development, implement an ETL that can, hopefully, survive the test of time.

All this is possible because of PDI's Metadata Injection step. Unfortunately, not all PDI steps can be injected yet so our template algorithm may need some adjustments in order to only injectable steps. However it's a much better position to be in than having to develop different transformations for each possible format and having ETL errors trigger more and more development work.

## Download the files

The files used in this article can be downloaded below.

PDI and data files

data.db (SQLite db; put it in your /tmp folder)

**Remark for Windows users:** to use the transformation as is you need to have SQLite installed, and you'll need to change the database connection in the template transformation to point to a valid path for the data.db file.