



Talend User Component tFileInputTextFlat

Purpose

This component extracts field from flat text files with two different methods:

1. Separate fields between delimiters
2. Separate fields at absolute position

The advantages are:

- Configure field position according to the names in the header line (also by regex)
- Read only the fields needed (unlike tFileInputDelimited)
- Tolerant against enclosure chars within the content
- All error messages contains the line number and the field name or index
- After extracting field content you can use regex to do some post processing
- Can load the value extraction configuration from an external file.

Talend-Integration

This component can be found in the palette under File->Input

This component provides an output and a reject flow and several return values.

Parameters

Property	Content
File name	Full path of the file to import.
Encoding	File encoding
Lines to Skip at Start	Number of lines to skip (In case you specify later you have an header line, don't count it here)
Skip empty lines	Detect empty lines and skip them
Schema	The schema. For date and timestamps you should specify here the pattern.
Has Column Header	Check it if your file provides a column header. This enables following options. The header will be skipped in the main output flow.
Use Header line to find position	The position of a delimited field will be set according to the position of the field name in the header line.
Find column position in header by regex	This option enables you to use a regularly expression in the field configuration "Name in Header".
Load configuration from file	Allows loading the field configuration from a external file.
Configuration file (.importconfig)	The file containing the configuration. See the paragraph below.
Ignore Not Null Constraints	Avoid throwing Exceptions if a value is null but the schema defines it as not null. This works for all columns.
Field Extraction	See explaining below
Field Separator	Char to separate the fields (only for field extracting by method Delimited Fields)
Text Enclosure	Char to enclosure the field content. It is helpful to read fields with line break as content. Don't escape double quotas here! E.g. for a double quoted field set "" here. It works even if not all fields are putted in quotas.
Split Row before Field	If it is not allowed to have line breaks in field content, check this. This option helps to check the correct file structure. The performance of the line separation can be increased.
Allow Enclosure within Content	If the enclosure char can occurs in the content (not as enclosure) check this to avoid parsing problems. If switched off the file structure can be checked strict.

Locale For Number Format	To parse the number different then the local pattern, specify here the locale. The default is the English format. E.g. in German we have the following pattern: 999.999,99 In this case set the locale to “de”. Quotation required!
Default Date Pattern	For all date schema column without a defined pattern, this pattern here will be applied. If none of the provided date formats are matching, the component performs a self-test typical for English, French and German formats and applying them. Only if this last attempt fails, the component throws and exception.
Die On Error	If true all errors stops the processing. If false all malformed data sets are send to the reject output flow (if added).

Field Extraction

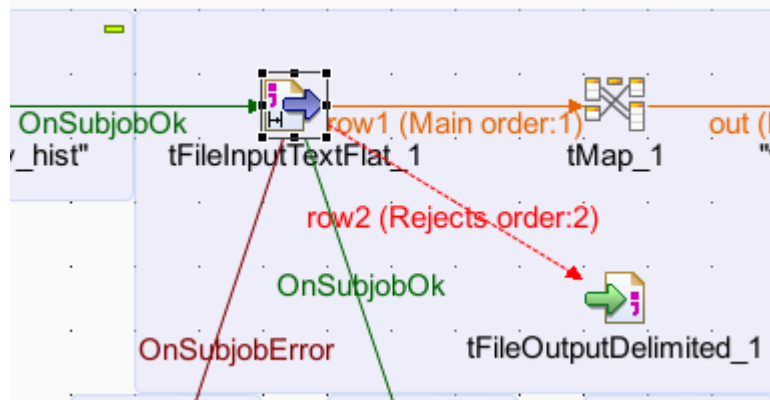
Column	Meaning
Column	Schema column name
Positioning	Field positioning method (extraction method)
Position	Position within delimiters or absolute position (depends on the chosen Positioning method) If blank for Delimited: Field position comes from field index in schema If blank for Absolute Positioning: Field position is the next after the last (it is a relative position to the previous field)
Length	Length of the field. It is required for absolute positioning. If a length is provided, the content will be trimmed to this length.
Regex	Regularly expression to post processing the field content. (No quotation needed)
Name in Header	Specify a different name for the column in the header. Field names often are not given according to the Java identifier naming rules. Therefore, you can specify the real name here (without quotation). If the option “Find column position in header by regex” is switch on you can write here regex (case insensitive). The regex expression must describe the whole possible column name in the header. It is not enough to declare only a pattern like contains. The regex expression will be surrounded from the component with ^ and \$. Example: column name in the header line can vary from: Payment_Product to P_Product. To match a schema column you can write as Name in Header this regex: <code>P[a-z]*_Product</code> to match both (and more) possibilities
Ignore If Missing	Sometimes we get fields which does not have all columns expected. You can specify here this column can be missed without problems. If missing the values remains blank.

Return values

Return value	Content
ERROR_MESSAGE	Last error message
NB_LINE	Number of delivered lines
NB_REJECTED	Number of rejected lines

Scenario 1:

Reading a delimited file and write malformed line into a file.



Configuration with examples of different header names.

It is an example of a file, which has as first line a header line.

File Name: *

Encoding:

Lines to Skip at Start: ☒ Skip Empty Rows

Schema: Edit schema ...

☒ Has Column Header ☒ Use Column Header to Find Position ☒ Ignore Not Null Constraints

Field Extraction

Column	Positioning	Position	Length	Regex	Name in Header	<input type="checkbox"/> Ignore if missing
Bestelldatum	Delimiter separated					<input type="checkbox"/>
Einkaufspreis	Delimiter separated			"[0-9.]+"		<input type="checkbox"/>
In_Zahlung_gegeben	Delimiter separated		4		In Zahlung gegeben	<input type="checkbox"/>
Gesamtpreis	Delimiter separated					<input type="checkbox"/>
Rabatt	Delimiter separated					<input type="checkbox"/>
Stornierungsdatum	Delimiter separated					<input type="checkbox"/>
Gew__Ankauf	Delimiter separated				Gew. Ankauf	<input checked="" type="checkbox"/>

Field Separator: *

Text Enclosure: ☐ Split Row Before Field ☒ Allow Enclosure within Content

Locale for Number Format: * Default Date Format: *

☐ Die on Error

It is not necessary to specify the position if the position is identical to the schema column index.

Scenario 2:

Using regularly expression to find the correct field position by the header line of the file.

In case the input files are provided by a system or organization, which cannot be motivated to a more fixed interface design. This should be avoided but sometimes it happens.

The screenshot shows the configuration window for the **tFileInputTextFlat_5** component. The **Basic settings** tab is active, showing the following configuration:

- File Name:** `*/private/var/testdata/text/different_schemas/regex/Input_File2.csv`
- Encoding:** `UTF-8`
- Lines to Skip at Start:** `0`
- Schema:** `Built-In`
- Has Column Header:** ☒
- Use Column Header to Find Position:** ☒
- Find column position in header by regex:** ☒
- Ignore Not Null Constraints:** ☐
- Field Extraction Table:**

Column	Positioning	Position	Length	Regex	Name in Header	Ignore if missing
ContactNumber	Delimiter separated				<code>*contact.*number*</code>	<input type="checkbox"/>
ExternalReference	Delimiter separated				<code>*(external).*(reference)*</code>	<input type="checkbox"/>
PaymentPlanProduct	Delimiter separated				<code>pp_product product* purchased</code>	<input type="checkbox"/>
PaymentPlanStartDate	Delimiter separated				<code>pp.start.date date*DDI*signed</code>	<input checked="" type="checkbox"/>
PaymentPlanNumber	Delimiter separated				<code>payment.plan.number order*number</code>	<input type="checkbox"/>

Below the table, the following settings are visible:

- Field Separator:** `*,`
- Text Enclosure:** `""`
- Split Row Before Field:** ☐
- Allow Enclosure within Content:** ☐
- Locale for Number Format:** `"en"`
- Default Date Format:** `"dd-MM-yyyy"`
- Die on Error:** ☒

The current example as a header line like this:

ContactNumber, ExternalReference, Source, DateDDISigned, Title, ProductPurchased, InterestedInBeingACampaigner, PaymentFrequency, ChangeOfPaymentFrequency, StartDate, NextPaymentDue, AutoPaymentMethodStartDate, MailingDate, CallOutcome, WelcomeCallDate, PaymentMethod, AccountName, PaymentPlanDetailBalance, OrderNumber

Only 4 columns are needed and the names can be vary from file to file.

This matching can only be done by regex – see the screenshot above.

Reject flow:

The reject flow will only be filled if the option “Die On Error” is switched off.

If you add the reject flow, at first it contains all columns from the schema and these additional columns. Only the green columns will be filled, all other columns can be erased here.

Column	Key	Type	z Nullab	Date Pat	Leng	Preci	Def	Comi
lineNumber	<input type="checkbox"/>	Integ	<input checked="" type="checkbox"/>		10	0		
line	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		4000	0		
errorMessage	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		255	0		

Column	Key	Type	z Nullab	Date Pat	Leng	Preci	Def	Comi
lineNumber	<input type="checkbox"/>	Integ	<input checked="" type="checkbox"/>		10	0		
line	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		4000	0		
errorMessage	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		255	0		

The line number is the number of the extracted lines and can differ from the number in the file if fields contain line breaks.

The line is the input for the field extraction. You can find here the probably malformed content.

The error message contains the description of the problem occurred while parsing the line.

Possible error messages:

If a column is missing in the header line (assuming we digger for it) an exception will be thrown which says the missing column.

If a column matches to the an already used delimiter position the read of the first line fails with exceptions like this:
`java.lang.Exception: nextDataRow failed in line 0:Extract field
pp.start.date|date*DDI*signed failed:Current field index 8 is lower then last
field index:9`

This is a failed check of the parser to avoid reading the content with a wrong configuration.

The component sorts all extractions internal by the position. It is a design decision to avoid reading content unwanted twice.

Loading the column configuration from a file

The component has the capability to load the column configuration from a file with the extension `.importconfig`. Typical use case is there are many different files which contains the same information but in different columns and positions. Instead of writing different jobs you can write one and describe every file with an import configuration.

The configuration file contains key value pairs. Do not quote (") the values!
It contains keys to describe the columns and keys to define delimiters and enclosures.

The column index must match to the column index of the schema.

The basic type describes the principle type of the column and can be left out as long as the data class in the key `CLASS` is set.

Key	Type	Description
<code>COLUMN_x_BASICTYPE</code>	Integer	See table below
<code>COLUMN_x_CLASS</code>	String	See table below
<code>COLUMN_x_DEFAULT</code>	String	This is the textual replacement for an empty value. Take as a replacement text in the file.
<code>COLUMN_x_DELIMITERCOUNT</code>	Integer	The position of the field in the row. It starts with 0 and allows gaps.
<code>COLUMN_x_ENABLED</code>	Boolean	Set it to false if you do not want to read this field
<code>COLUMN_x_IGNORE_DATASET_IF_INVALID</code>	Boolean	Allow continuing with the next rows if this field has an invalid value.
<code>COLUMN_x_IGNORE_MISSING_COLUMN</code>	Boolean	If the column position will be found in the header line this flag allows this column to be missing.
<code>COLUMN_x_LOCALE</code>	String	To interpret numbers correctly the component needs to know the country or language. Set here the ISO 2-letter or 5-letter country code. Examples: en or en_UK or de_DE
<code>COLUMN_x_NAME</code>	String	The name of the column in the header. Also if no header search is intended the name is mandatory. As name is also possible a regex expression if the component has to search the column in the header by the help of regular expressions.
<code>COLUMN_x_NULL_ENABLED</code>	Boolean	Null value allowed or not
<code>COLUMN_x_POSITIONTYPE</code>	Integer	Absolute position = 0 Delimited = 2 Delimited with max. length = 3
<code>COLUMN_x_TRIM</code>	Boolean	Trim the value if true. Regardless of this setting empty values (number white spaces == 0) will always be returned as null!
<code>CHARSET</code>	String	The charset of the file. Typical values are: UTF-8, UTF-16, Cp-1252, ISO-8859-15, ASCII
<code>ALLOW_ENCLASURE_IN_TEXT</code>	Boolean	If enclosures are used this flag allows enclosures also in the text itself.
<code>DELIMITER</code>	char	The delimiter character
<code>ENCLOSURE</code>	char	The character used to enclose the field content. Especially useful if the delimiter can also be part of the content or if the content contains line breaks.
<code>IGNORE_ENCLOSED_LINE_BREAK</code>	Boolean	If the content can contain line breaks and is enclosed, set this true
<code>IGNORE_NOT_NULL_CONSTRAINTS</code>	Boolean	true = switch off the check of nullable or not.
<code>SKIP_EMPTY_LINES</code>	Boolean	Empty lines will be skipped
<code>SKIP_ROWS</code>	Integer	Skip number of lines in the file before starting parsing it. The header-line will be taken after skipping lines.

These are the possible classes:

Data class	Corresponding basic type
String	0
Date	1
BigDecimal	2
Long	2
Integer	2
Double (default)	2
Float	2
Short	2
Boolean	8

Scenario for loading the configuration from an external file

At first to create the configuration file (as an template for your own changes) you can use the component with a job specific field configuration and if it works you can save it in the advanced settings.

The configuration file can be created in the advanced settings (only if the Load-option is switched off).

The manual configuration of the field extraction:

Save the configuration file:

Here the configuration file matching to the manual configuration above.

```
ALLOW_ENCLASURE_IN_TEXT=false
CHARSET=UTF-8
COLUMN_0_BASICTYPE=0
COLUMN_0_CLASS=String
COLUMN_0_DEFAULT=99
COLUMN_0_DELIMITERCOUNT=0
COLUMN_0_ENABLED=true
COLUMN_0_IGNORE_DATASET_IF_INVALID=false
COLUMN_0_IGNORE_MISSING_COLUMN=false
COLUMN_0_LOCALE=en_US
COLUMN_0_NAME=*contact.*number*
COLUMN_0_NULL_ENABLED=true
COLUMN_0_POSITIONTYPE=2
COLUMN_0_TRIM=false
COLUMN_1_BASICTYPE=0
COLUMN_1_CLASS=String
```

```

COLUMN_1_DELIMITERCOUNT=1
COLUMN_1_ENABLED=true
COLUMN_1_IGNORE_DATASET_IF_INVALID=false
COLUMN_1_IGNORE_MISSING_COLUMN=false
COLUMN_1_LOCALE=en_US
COLUMN_1_NAME=*(external).* (reference)*
COLUMN_1_NULL_ENABLED=true
COLUMN_1_POSITIONTYPE=2
COLUMN_1_TRIM=false
COLUMN_2_BASICTYPE=1
COLUMN_2_CLASS=Date
COLUMN_2_DELIMITERCOUNT=8
COLUMN_2_ENABLED=true
COLUMN_2_FORMAT=dd.MM.yyyy
COLUMN_2_IGNORE_DATASET_IF_INVALID=false
COLUMN_2_IGNORE_MISSING_COLUMN=false
COLUMN_2_LOCALE=en_US
COLUMN_2_NAME=pp.start.date|date.DDI.signed
COLUMN_2_NULL_ENABLED=true
COLUMN_2_POSITIONTYPE=2
COLUMN_2_TRIM=false
COLUMN_3_BASICTYPE=2
COLUMN_3_CLASS=Long
COLUMN_3_DELIMITERCOUNT=71
COLUMN_3_ENABLED=true
COLUMN_3_FORMAT=en
COLUMN_3_IGNORE_DATASET_IF_INVALID=false
COLUMN_3_IGNORE_MISSING_COLUMN=false
COLUMN_3_LOCALE=en_US
COLUMN_3_NAME=payment.plan.number|order.number
COLUMN_3_NULL_ENABLED=true
COLUMN_3_POSITIONTYPE=2
COLUMN_3_TRIM=false
COLUMN_4_BASICTYPE=0
COLUMN_4_CLASS=String
COLUMN_4_DELIMITERCOUNT=45
COLUMN_4_ENABLED=true
COLUMN_4_IGNORE_DATASET_IF_INVALID=false
COLUMN_4_IGNORE_MISSING_COLUMN=false
COLUMN_4_LOCALE=en_US
COLUMN_4_NAME=pp.product|product.purchased
COLUMN_4_NULL_ENABLED=true
COLUMN_4_POSITIONTYPE=2
COLUMN_4_TRIM=false
COLUMN_5_BASICTYPE=0
COLUMN_5_CLASS=String
COLUMN_5_DELIMITERCOUNT=4
COLUMN_5_ENABLED=true
COLUMN_5_IGNORE_DATASET_IF_INVALID=false
COLUMN_5_IGNORE_MISSING_COLUMN=false
COLUMN_5_LOCALE=en_US
COLUMN_5_NAME=Fundraiser
COLUMN_5_NULL_ENABLED=true
COLUMN_5_POSITIONTYPE=2
COLUMN_5_TRIM=false
DELIMITER=,
ENCLOSURE="
IGNORE_BOM=false
IGNORE_ENCLOSED_LINE_BREAK=true
IGNORE_NOT_NULL_CONSTRAINTS=false
SKIP_EMPTY_LINES=true
SKIP_ROWS=0

```

Now switch on the Load option (as in the screenshot above) and it should work in the same way.

The screenshot shows the Talend Job Configuration window for a job named "Job test_tFileInputTextFlat 0.1". The "Advanced settings" tab is active. The "Load" checkbox is checked. The "Config file (.importconfig)" field is set to "/Volumes/Data/Talend/testdata/text/test.importconfig". Other settings include "File Name" as "/private/var/testdata/text/different_schemas/regex/Input_File.csv", "Encoding" as "UTF-8", "Lines to Skip at Start" as 0, and "Skip Empty Rows" checked. The "Field Separator" is set to ",", "Text Enclosure" is set to "\"", "Locale for Number Format" is set to "en", and "Default Date Format" is set to "dd-MM-yyyy".