



tFileExcelSheetInput

Purpose

This component reads Excel sheets.

This component needs the components tFileExcelWorkbookOpen (open a file or creates a workbook)

Advantages of this component:

- The columns to read can be set also with gaps
- The column position can automatically configured by a header line and additional by the use of regularly expressions
- Can read reliable all possible data types and tries also to convert it into the schema target column type
- Can read comments
- Can fill the output stream for empty cells with the last not empty value
- Can ignore cell read errors e.g. in case of a type conversion is impossible
- Uses always the latest Apache POI API

Talend-Integration

This component can be found in the palette under File/Spreadsheet

This component provides several return values.

Parameters for tFileExcelSheetInput

Property	Content
Workbook	Choose the tFileExcelWorkbookOpen component holding the Apache POI Excel workbook
Schema	The schema of the output
Sheet name	The name or the index of the sheet. Please take care of a valid sheet name or simply type the index of the sheet your want to write. If the sheet does not exist, it will be created automatically. You can take the sheet name from the return value of the tFileExcelSheetList.
Row start index	The component starts reading data with this row (1-based)
Stop at empty row	If no values from the row was received (only all configured columns) the component stops reading.
Skip empty rows	If the all needed values are empty, the row will be skipped.
Column start	Is visible only if no individual column configuration takes place
Use individual column configuration	If chosen it shows the configuration of columns and its position can be configured individual
User header to configure position of columns	If true the column position will be configured according to its position in the header line. The name will be found not case sensitive.
Header line	In the column configuration it is possible to configure the column position by the header line. Here set the index of the header line (1-based)
Column configuration	You can specify the columns in the Column Configuration in the column Sheet Column Name. Here you can use the Excel letter reference ("A" for the first column) or an index (0 for the first column). It is possible to have gaps between the different columns (unlike the build-in Talend components). Column: Name of the schema column Sheet Column Name: Column position as letter (starts with "A" or 0) Name in Header: if position should be found in header, set here the name in the header (also with the use of regex) Read cell comment: If checked the comment will taken as value Use last value for empty: If checked an empty cell will be filled with the last known value Ignore Errors: If something goes wrong this option avoids aborts

Advanced setting parameters for tFileExcelSheetInput

Property	Content
Language / Country for number format	Number formats are different for different languages/countries. In case of the number is stored in a text typed cell and the schema expects a number this local will be used to find the correct format pattern for the text-to-number conversion.
Return Hyperlink URL	The normal value for a cell is always the visible cell value. If the cell has an underlying hyperlink, this option must be switched on to get the hyperlink instead of the value.
Concatenate Label URL	The label of the hyperlink is the visible cell value and if both (label and hyperlink) are needed the component read both values and concatenates them with a pipe symbol.
Trim columns	If checked all textual content will be trimmed (leading and trailing spaces, tabulators or line breaks will be removed)
Use cached value if formula evaluation fails	It could happen especially if a formula references to external files, the formula fails. Excel usually keeps the value of the last evaluation and this value will be delivered.

Return values of the component:

Value	Content
NB_LINE	Number of lines read
ERROR_MESSAGE	Error message if something went wrong
LAST_ROW_INDEX	Index of the last read row in this sheet.
MAX_ROW_INDEX	The maximum row index available in the sheet (not necessarily read)
CURRENT_ROW_INDEX	The current absolute excel row index (available within the flow)

Scenario 1: Read from cells referenced by the Excel column names

Column	Sheet column name	Read cell comment	Use last value for empty	Ignore Errors
date	"A"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
value	"B"	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
id	"D"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
result	"E"	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

The cell can be addressed with the well-known Excel column name (starting with "A") or the cell index (starting with 0).

Scenario 2: Read from cells referenced by the header line

Designer | Code

Job(test_excel_input 0.1) Contexts(Job test_excel_input 0.1) Component Problems Modules Run (Job test_excel_input)

tFileExcelWorkbookOpen_1 OnSubjobOk tFileExcelSheetInput_1 row1 (Main) tLogRow_1

tFileExcelSheetInput_1

Basic settings Workbook: tFileExcelWorkbookOpen_1 | Schema: Built-In | Edit schema

Advanced settings Sheet name or index: 0

Dynamic settings Row start index (starts with 1): 2 Limit rows to:

View ☒ Stop at empty row

Documentation ☒ Use individual column configuration ☒ Use Header to config positions of columns ☒ Find column position in header by regex

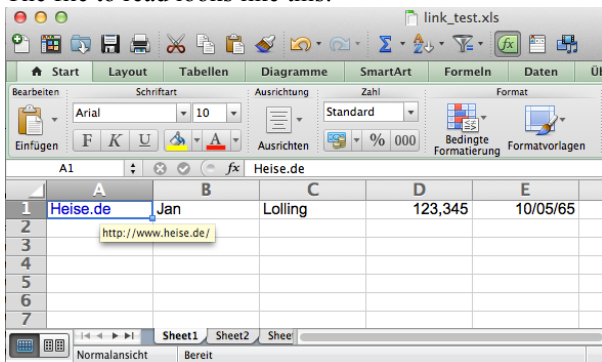
Header line: 1

Column	Name in Header	<input type="checkbox"/> Ignore if missing	<input type="checkbox"/> Read cell comment	<input type="checkbox"/> Use last value for empty	<input type="checkbox"/> Ignore Errors
date	"Create\date"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
value	"Cost per piece"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
id	"Customer\{A-Z}*"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
result	"Calculated Costs"	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

In this scenario the column the header line will configure positions. The component tries to find the column by its name (case insensitive) or by regularly expressions (also case insensitive).

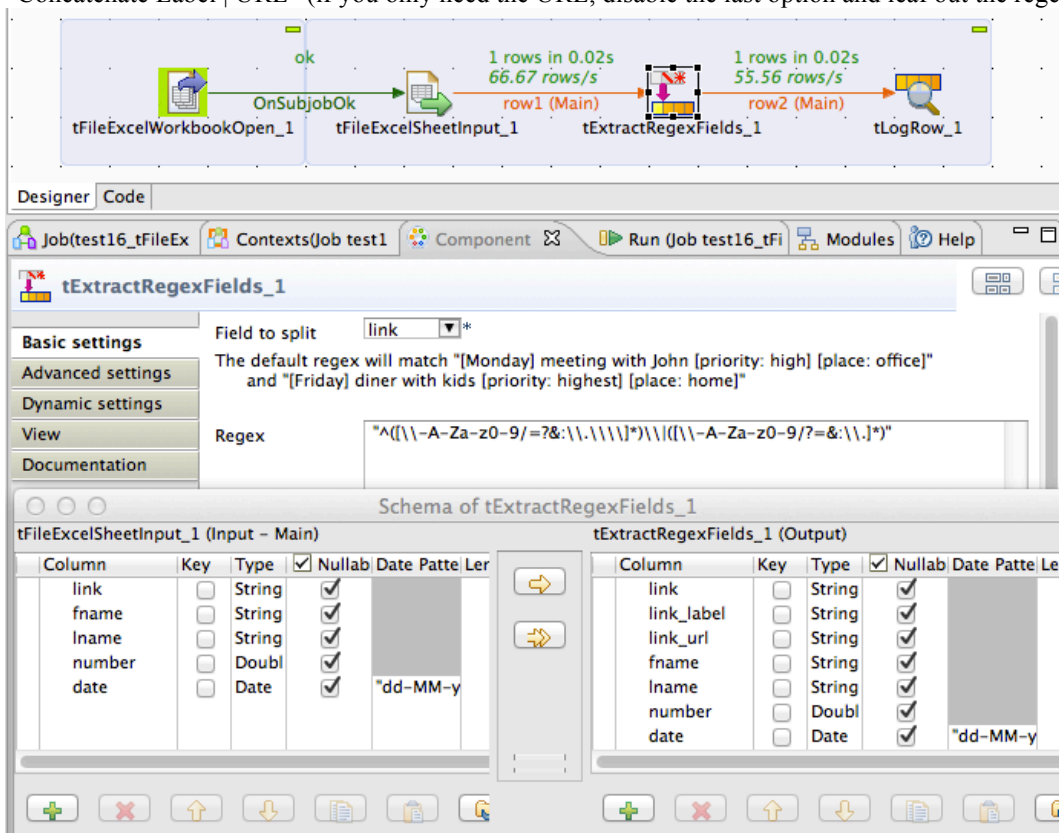
Scenario 3: Read hyperlinks from a cell

The file to read looks like this:



Here a simple job reading the hyperlink and separate label and URL with a regex expression.

In the component tFileExcelSheetInput activate in the advanced settings the options “Return Hyperlink URL” and “Concatenate Label | URL” (if you only need the URL, disable the last option and leaf out the regex component).



The tExtractRegexFields expects after the parsed field (link) additional fields as much as you want to extract content by regex groups. It is highly recommended to check the regex expression with external tools and take care you get only one regex sequence with (in this case) to groups. Please keep in mind every regex sequence causes an output record (e.g. an additional output record).

The regex expression here is: `"^([\-\A-Za-z0-9/?=&:\.\|\|]*)([\-\A-Za-z0-9/?=&:\.\|]*)"`

The output of the job:

tLogRow_1						
link	link_label	link_url	fname	lname	number	date
Heise.de	http://www.heise.de/	Heise.de	http://www.heise.de/	Jan	Lolling	123.345 10-05-1965