

# tFileInputKeyValue



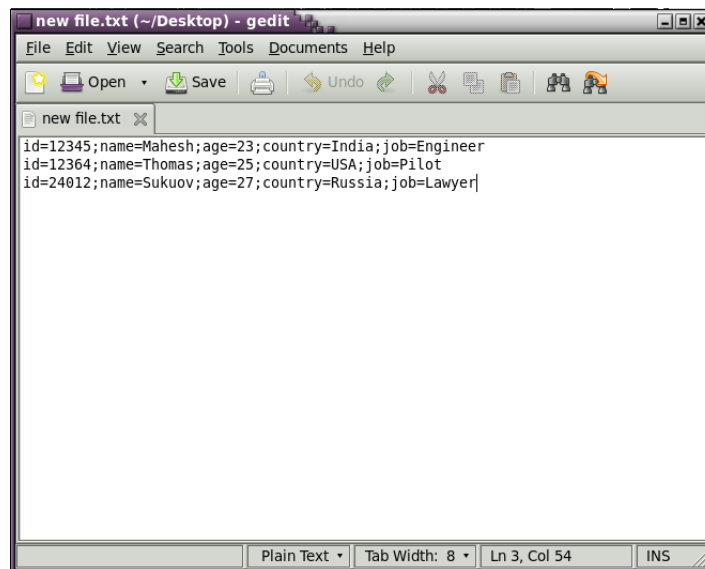
## tFileInputKeyValue properties

<b>Version</b>	2.0	
<b>Component family</b>	File/Input	
<b>Function</b>	<b>tFileInputKeyValue</b> reads a given file containing key-value pairs, row by row with separated fields.	
<b>Purpose</b>	Opens a file containing “key-value” pairs in each row, and extracts the “values” from the “keys”. The “keys” are taken as fields and “values” as its contents. (This component allows empty key-value pairs in the input). It then sends the extracted values as defined in the Schema to the next Job component, via a Row link.	
<b>Basic settings</b>	<i>File Name</i>	<b>File name:</b> Name and path of the file to be processed.  In order to avoid the inconvenience of hand writing, you could select the variable of interest from the auto-completion list (Ctrl +Space) to fill the current field on condition that this variable has been properly defined.
	<i>Field Separator(Regex)</i>	It is used to separate key-value pairs in each row. It has to be a Regex Expression.
	<i>Row Separator(Regex)</i>	It is used to distinguish rows. It has to be a Regex Expression.
	<i>Key-Value Separator(Regex)</i>	It separates a Value from a key. It has to be a Regex Expression.
	<i>Key Name</i>	<b>Column:</b> This field is automatically populated with the columns defined in the schema that you propagated.
		<b>Key (Regex):</b> This field should contain the “key” names as present in the input file. It has to be a Regex Expression.
		<b>Trim Value:</b> Select this check box to remove leading and trailing whitespaces from defined columns.
	<i>Header</i>	Number of rows to be skipped in the beginning of file
	<i>Footer</i>	Number of rows to be skipped at the end of the file.
	<i>Limit</i>	Maximum number of rows to be processed. If Limit = 0, no row is read or processed.
	<i>Schema and Edit Schema</i>	A schema is a row description, i.e., it defines the number of fields that will be processed and passed on to the next component. The schema is either built-in or remote in the Repository.
		<b>Built-in:</b> The schema will be created and stored locally for this component only.
		<b>Repository:</b> The schema already exists and is stored in the

		Repository, hence can be reused in various projects and Job flowcharts.
	<i>Die on error</i>	Select this check box to stop the execution of the Job when an error occurs. Clear the check box to skip the row on error and complete the process for error-free rows.
<b>Advanced settings</b>	<i>tStatCatcher Statistics</i>	Select this check box to gather the processing metadata at the Job level as well as at each component level.
<b>Usage</b>	<p>Use this component to read a file containing “key-value” pairs in each row. It extracts the “values” from the “keys”.</p> <p>The “keys” are taken as fields and “values” as its contents. (This component even allows empty key-value pairs in the input).</p> <p>It then sends the extracted values as defined in the Schema to the next Job component, via a Row link.</p> <p>For further information, please see <a href="#">“Scenario 2”</a> and <a href="#">“Scenario 3”</a>.</p>	
<b>New Features</b>	<p>Allows the key to keep on changing for each row.</p> <p>Allows the key-value separator to duplicate in the value.</p> <p>Allows the key-value separator to duplicate in the key.</p>	
<b>Author</b>	Mahesh M. Pillai	

## Scenario 1: Reading data from a file containing multiple Key-Value Pairs in a Row

The following scenario creates a two-component Job, which aims at reading each row (containing multiple key-value pair) of a file, extracts the “value” from the “keys” and displays the output in the Run log console. The contents of the Input Text File are shown below:

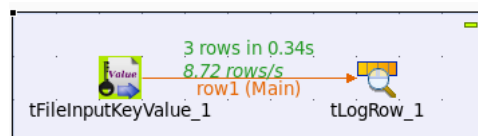


```

id=12345;name=Mahesh;age=23;country=India;job=Engineer
id=12364;name=Thomas;age=25;country=USA;job=Pilot
id=24012;name=Sukuov;age=27;country=Russia;job=Lawyer

```

- Here the keys are “id”, “name”, “age”, “country”, and “job”.



- Drop a **tFileInputKeyValue** component from the Palette to the design workspace.
- Drop a **tLogRow** component the same way.
- Right-click on the **tFileInputKeyValue** component and select **Row > Main**. Then drag it onto the **tLogRow** component and release when the plug symbol shows up.
- Select the **tFileInputKeyValue** component again, and define its **Basic** settings:

Component: tFileInputKeyValue\_1

**Basic settings**

File Name: /home/550778/Desktop/new file.txt

Schema: Built-in Edit schema

Field Separator(Regex): ;, Row Separator(Regex): \n Key-Value Separator(Regex): =

Column	Key	Trim Value
EmployeeId	"id"	<input checked="" type="checkbox"/>
EmployeeName	"name"	<input type="checkbox"/>
EmployeeAge	"age"	<input type="checkbox"/>
EmployeeCountry	"country"	<input type="checkbox"/>
EmployeeDesignation	"job"	<input type="checkbox"/>

Header: 0 Footer: 0 Limit:

☒ Die on error

- Fill in a path to the file in the **File Name** field. This field is mandatory.
- If the path of the file contains some accented characters, you will get an error message when executing your Job. For more information regarding the procedures to follow when the support of accented characters is missing, see Talend Open Studio for Big Data Installation Guide.
- Set the **Schema** as either a local (Built-in) or a remotely managed (Repository) to define the data to pass on to the **tLogRow** component.
  - You can load and/or edit the schema via the **Edit Schema** function.
- In the present scenario we have in total a maximum of 5 key-value pairs in a row.
- You may define them as shown below.

Schema oftFileInputKeyValue\_1

tFileInputKeyValue\_1

Column	Key	Type	Nullable	Date Pattern	Length	Precision	Default	Comment
EmployeeId	<input type="checkbox"/>	int	<input type="checkbox"/>					
EmployeeName	<input type="checkbox"/>	String	<input type="checkbox"/>					
EmployeeAge	<input type="checkbox"/>	short	<input type="checkbox"/>					
EmployeeCountry	<input type="checkbox"/>	String	<input type="checkbox"/>					
EmployeeDesignation	<input type="checkbox"/>	String	<input type="checkbox"/>					

OK Cancel

- The column name can be any valid name.
- Define the **Field separator** used to delimit key-value pairs in a row (Here it is “;”). Then define the **Row Separator**

allowing to identify the end of a row (Here it is “\n”). Also define the **Key-Value Separator** (Here it is “=”).

These fields should be in Regex Patterns.

- Now the **Key Name** Field should be filled.

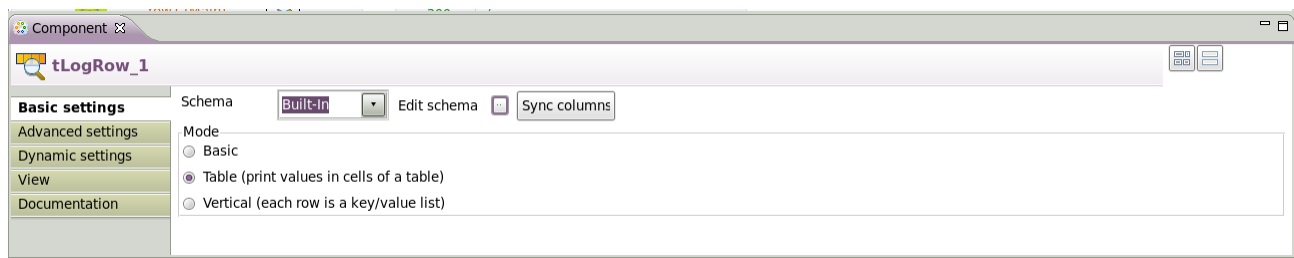
**Column** field is automatically populated with the columns defined in the schema that you propagated.

The **Key** field corresponds to the keys in the file that we have opened.

In the current scenario the keys are “id”, “name”, “age”, “country”, and “job”.

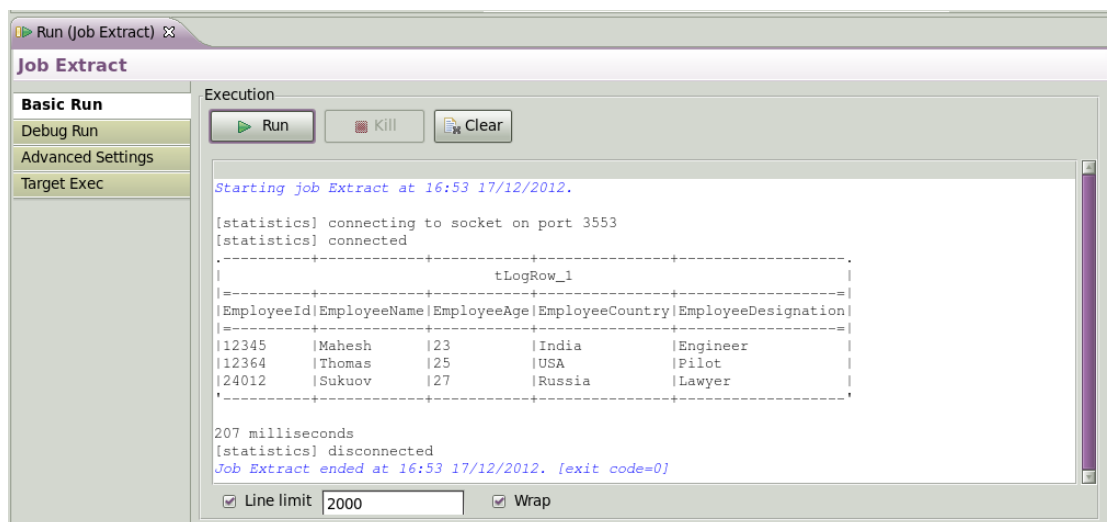
The **Trim Value** field can be used to remove leading and trailing whitespaces from the values corresponding to the fields.

- In this scenario, the **Header**, **Footer**, **Limit** numbers are not set.
- Die on Error can be set as per need. (Select this check box to stop the execution of the Job when an error occurs. Clear the check box to skip the row on error and complete the process for error-free rows. )
- Select the tLogRow and define the Field separator to use for the output display.



- Go to **Run** tab, and click on Run to execute the Job.

The file is read row by row and the extracted fields are displayed on the Run log as defined in both components Basic settings.



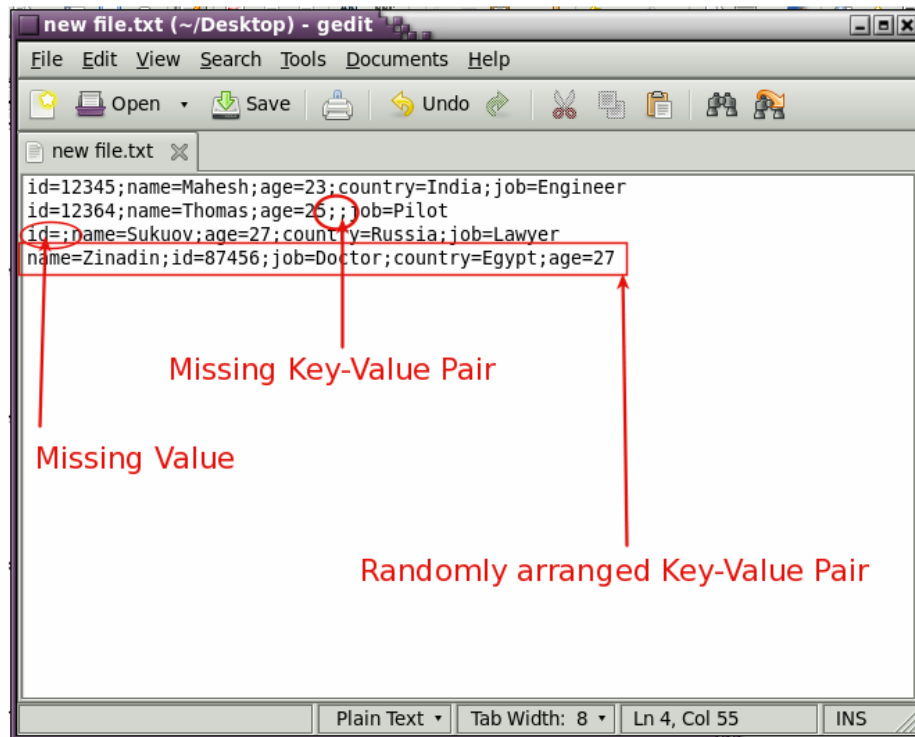
## Scenario 2: Reading data from a file that contains

- **Key-Value Pairs randomly distributed in a row**
- **Key-Value Pairs missing in a row**

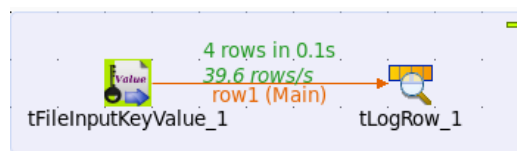
The following scenario describes the **Flexibility** this component offers:

- It allows missing key value pairs.
- It allows missing values.
- It allows random distribution of key-value pairs within the row.

The following scenario depicts a two-component Job, which aims at reading each row (containing multiple key-value pair) of a file, extracts the “value” from the “keys” and displays the output in the Run log console. The above mentioned (*in italics*) conditions are created in the input file. The Input File (with Annotations) is shown below:



- Drop a **tFileInputKeyValue** component from the Palette to the design workspace.
- Drop a **tLogRow** component the same way.
- Right-click on the **tFileInputKeyValue** component and select **Row > Main**. Then drag it onto the **tLogRow** component and release when the plug symbol shows up.



- Select the **tFileInputKeyValue** component again, and define its **Basic** settings:

**tFileInputKeyValue\_1**

**Basic settings**

File Name:

Schema: **Built-in**

Field Separator(Regex):  Row Separator(Regex):  Key-Value Separator(Regex): 

Key Name

Column	Key	Trim Value
EmployeeId	"id"	<input checked="" type="checkbox"/>
EmployeeName	"name"	<input type="checkbox"/>
EmployeeAge	"age"	<input type="checkbox"/>
EmployeeCountry	"country"	<input type="checkbox"/>
EmployeeDesignation	"job"	<input type="checkbox"/>

Header:  Footer:  Limit:

☒ Die on error

- Fill in a path to the file in the **File Name** field. This field is mandatory.  
If the path of the file contains some accented characters, you will get an error message when executing your Job. For more information regarding the procedures to follow when the support of accented characters is missing, see Talend Open Studio for Big Data Installation Guide.
- Set the **Schema** as either a local (Built-in) or a remotely managed (Repository) to define the data to pass on to the **tLogRow** component.
- You can load and/or edit the schema via the **Edit Schema** function.  
In the present scenario we have in total a maximum of 5 key-value pairs in a row.  
You may define them as shown below.

**Schema of tFileInputKeyValue\_1**

tFileInputKeyValue\_1

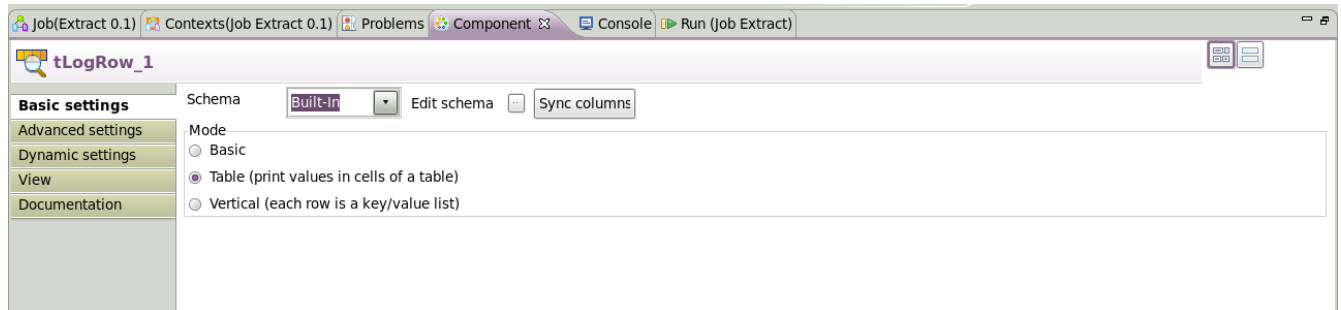
Column	Key	Type	Nullable	Date Pattern	Length	Precision	Default	Comment
EmployeeId	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
EmployeeName	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
EmployeeAge	<input type="checkbox"/>	Short	<input checked="" type="checkbox"/>					
EmployeeCountry	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
EmployeeDesignation	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					

OK Cancel

- The column name can be any valid name.
- Define the **Field separator** used to delimit key-value pairs in a row (Here it is “;”). Then define the **Row Separator** allowing to identify the end of a row (Here it is “\n”). Also define the **Key-Value Separator** (Here it is “=”).
- These fields should be in Regex Patterns.
- Now the **Key Name** Field should be filled.  
**Column** field is automatically populated with the columns defined in the schema that you propagated.  
The **Key** field corresponds to the keys in the file that we have opened.  
In the current scenario the keys are “id”, “name”, “age”, “country”, and “job”.  
The **Trim Value** field can be used to remove leading and trailing whitespaces from the values corresponding to the fields.
- In this scenario, the **Header**, **Footer**, **Limit** numbers are not set.
- Die on Error can be set as per need. (Select this check box to stop the execution of the Job when an error occurs. Clear the

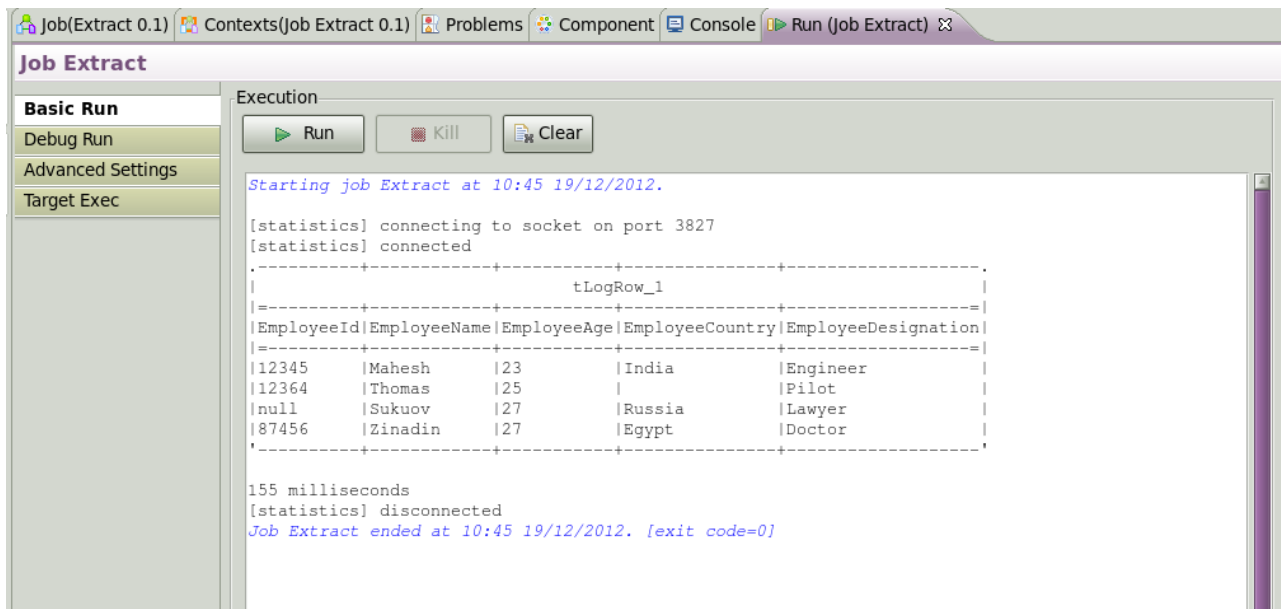
check box to skip the row on error and complete the process for error-free rows. )

- Select the tLogRow and define the Field separator to use for the output display.



- Go to **Run** tab, and click on Run to execute the Job.

The file is read row by row and the extracted fields are displayed on the Run log as defined in both components Basic settings.



In spite of the input being in an unordered format, the output is in the correct format.

## Scenario 3: Reading data from a file that contains

- **Dynamic Key Name**
- **Key-Value Separator duplicated in the “value”**
- **Key-Value Separator duplicated in the “key”**

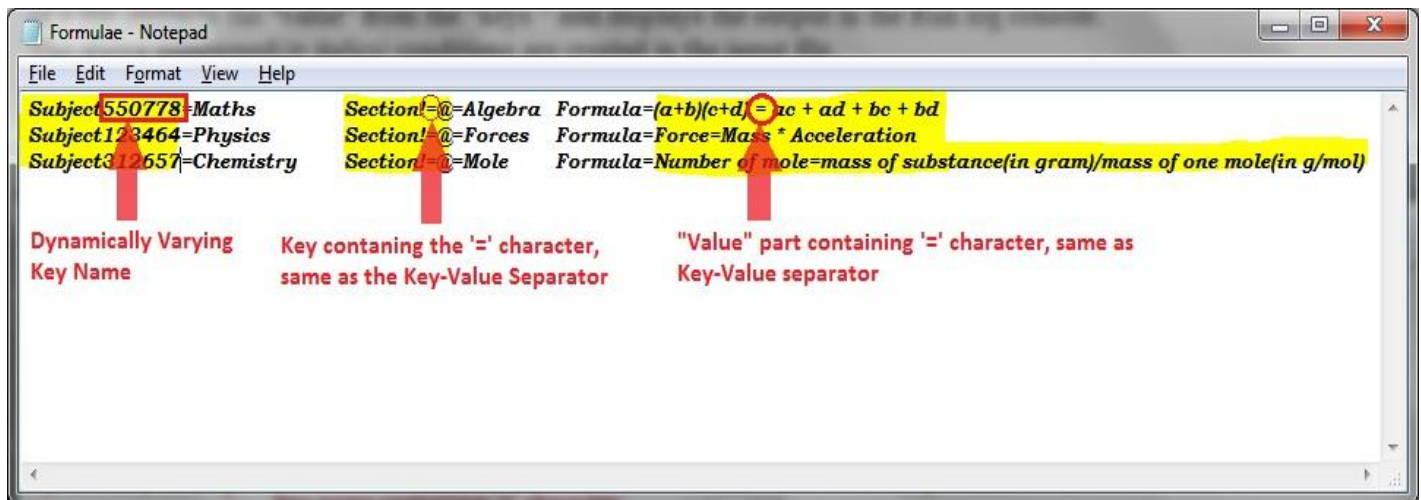
The following scenario describes the *latest Flexibility* this component offers:

- It allows the key to keep on changing for each row.
- It allows the key-value separator to be present in the value.
- It allows the key-value separator to be present in the key.

The following scenario depicts a two-component Job, which aims at reading each row (containing multiple key-value pair) of a file, extracts the “value” from the “keys” and displays the output in the Run log console.

The above mentioned (*in italics*) conditions are created in the input file.

The Input File (with Annotations) is shown below:



There are three tab separated fields in the above file.

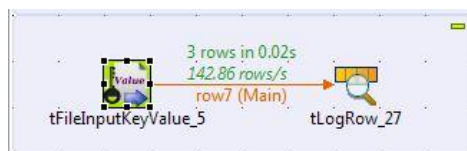
The key-value separator is ‘=’.

The first one is the *Subject* field .The key name is varying for each row. viz., Subject550778, Subject123464, etc. The second one is the *Section* field. The key name (Section!@) contains ‘=’ which is same as the key-value separator.

The third field is the *Formula*. The value itself contain the key-value separator(=). viz., Force=Mass \* Acceleration.

In these scenarios, the **tFileInputKeyValue** component can be leveraged to filter off the value from the key.

- Drop a **tFileInputKeyValue** component from the Palette to the design workspace.
- Drop a **tLogRow** component the same way.
- Right-click on the **tFileInputKeyValue** component and select **Row > Main**. Then drag it onto the **tLogRow** component and release when the plug symbol shows up.



- Select the **tFileInputKeyValue** component again, and define its **Basic** settings:



**tFileInputKeyValue\_5**

File Name: C:/Users/User/Desktop/Formulae.txt

Schema: Built-In Edit schema

Field Separator(Regex): \t Row Separator(Regex): \n Key-Value Separator(Regex): =

Column	Key(Regex)	Trim Value
Subject	Subject\\d{6}	<input checked="" type="checkbox"/>
Section	Section!=@	<input checked="" type="checkbox"/>
Formula	Formula	<input checked="" type="checkbox"/>

Header: 0 Footer: 0 Limit:

☐ Die on error

- Fill in a path to the file in the **File Name** field. This field is mandatory. If the path of the file contains some accented characters, you will get an error message when executing your Job. For more information regarding the procedures to follow when the support of accented characters is missing, see Talend Open Studio for Big Data Installation Guide.
  - Set the **Schema** as either a local (Built-in) or a remotely managed (Repository) to define the data to pass on to the **tLogRow** component.
  - You can load and/or edit the schema via the **Edit Schema** function.
- In the present scenario we have in total a maximum of 3 key-value pairs in a row. You may define them as shown below.

Schema of tFileInputKeyValue\_5

Column	Key	Type	N.	Date Patte...	Len...	Prec...	De...	Co...
Subject	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Section	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Formula	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					

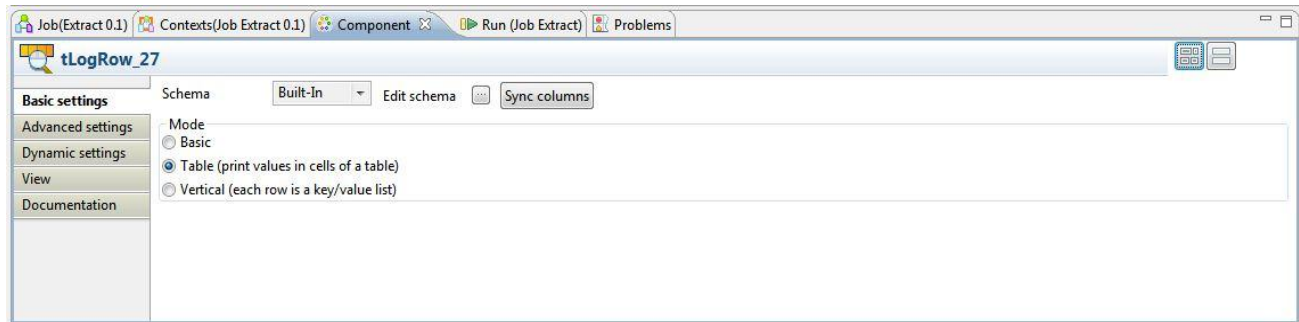
OK Cancel

- The column name can be any valid name.
  - Define the **Field separator** used to delimit key-value pairs in a row (Here it is "\t"). Then define the **Row Separator** allowing to identify the end of a row (Here it is "\n"). Also define the **Key-Value Separator** (Here it is "=").
  - These fields should be in Regex Patterns.
  - Now the **Key Name** Field should be filled.
- Column** field is automatically populated with the columns defined in the schema that you propagated. The **Key** field corresponds to the keys in the file that we have opened. The **Key** should be in Regex Patterns. This allows the user to have control over defining patterns.

In the current scenario the key for the first field is "Subject" followed by a six-digit number. So, the corresponding Regex Pattern would be "Subject\\d{6}". The second key is "Section!=@" and the third is "Formula".

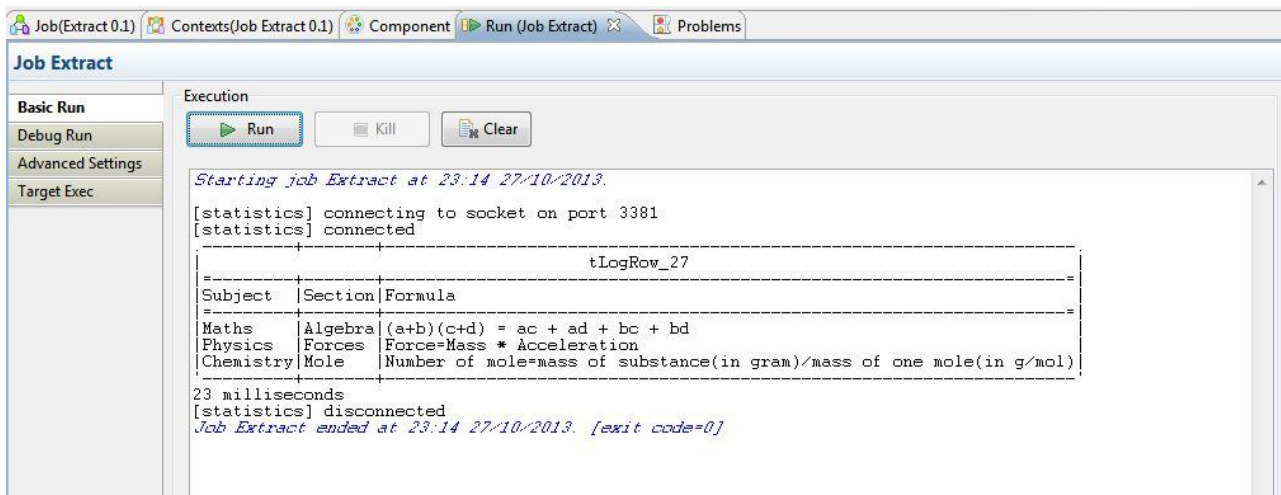
The **Trim Value** field can be used to remove leading and trailing whitespaces from the values corresponding to the fields.

- In this scenario, the **Header**, **Footer**, **Limit** numbers are not set.
- Die on Error can be set as per need. (Select this check box to stop the execution of the Job when an error occurs. Clear the check box to skip the row on error and complete the process for error-free rows. )
- Select the tLogRow and define the Field separator to use for the output display.



- Go to **Run** tab, and click on Run to execute the Job.

The file is read row by row and the extracted fields are displayed on the Run log as defined in both components Basic settings.



In spite of the key name being dynamic and the key/value containing the key-value separator (=), the output is obtained in the required format.