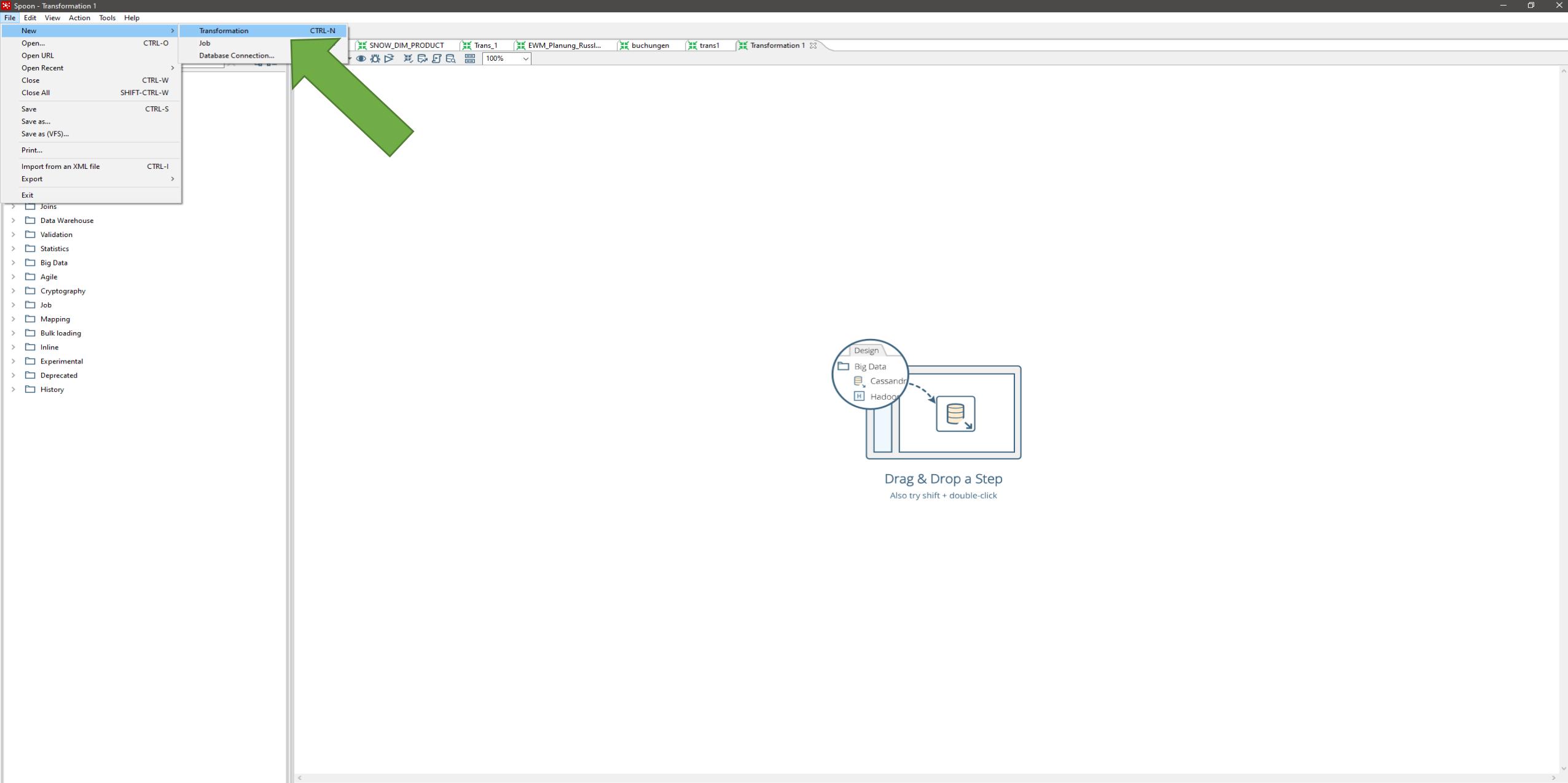
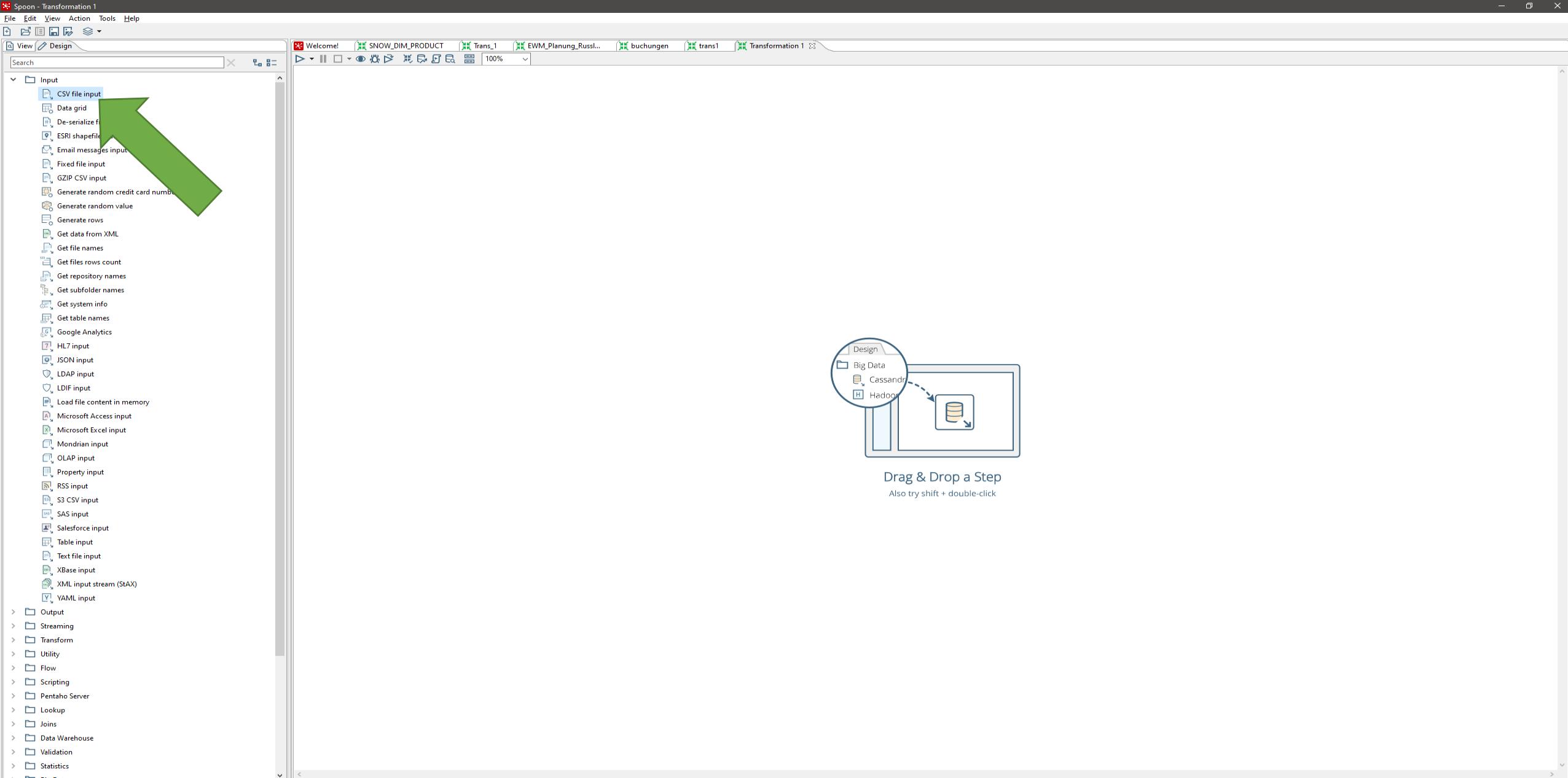


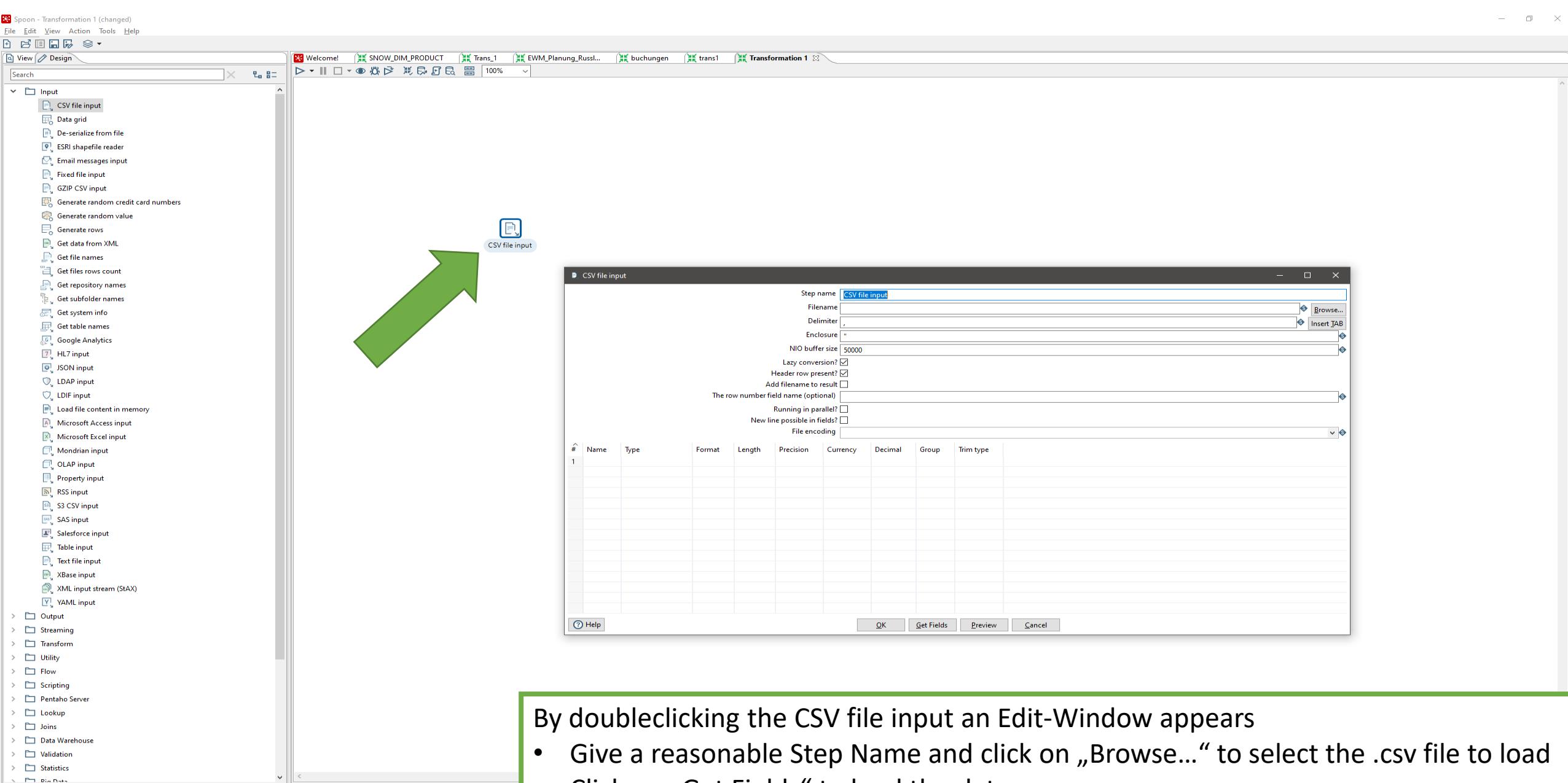
# Tutorial ETL - Processing



In order to start a new transformation, click on File > New > Transformation



- To import the operational data, click on Input and drag the CSV file input into the main window



By doubleclicking the CSV file input an Edit-Window appears

- Give a reasonable Step Name and click on „Browse...“ to select the .csv file to load
- Click on „Get Fields“ to load the data
- Press OK

View Design

Search

Input

- CSV file input
- Data grid
- De-serialize from file
- ESRI shapefile reader
- Email messages input
- Fixed file input
- GZIP CSV input
- Generate random credit card numbers
- Generate random value
- Generate rows
- Get data from XML
- Get file names
- Get files rows count
- Get repository names
- Get subfolder names
- Get system info
- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input
- SAS input
- Salesforce input
- Table input
- Text file input
- XBase input
- XML input stream (StAX)
- YAML input

Output

Streaming

Transform

Utility

Flow

Scripting

Pentaho Server

Lookup

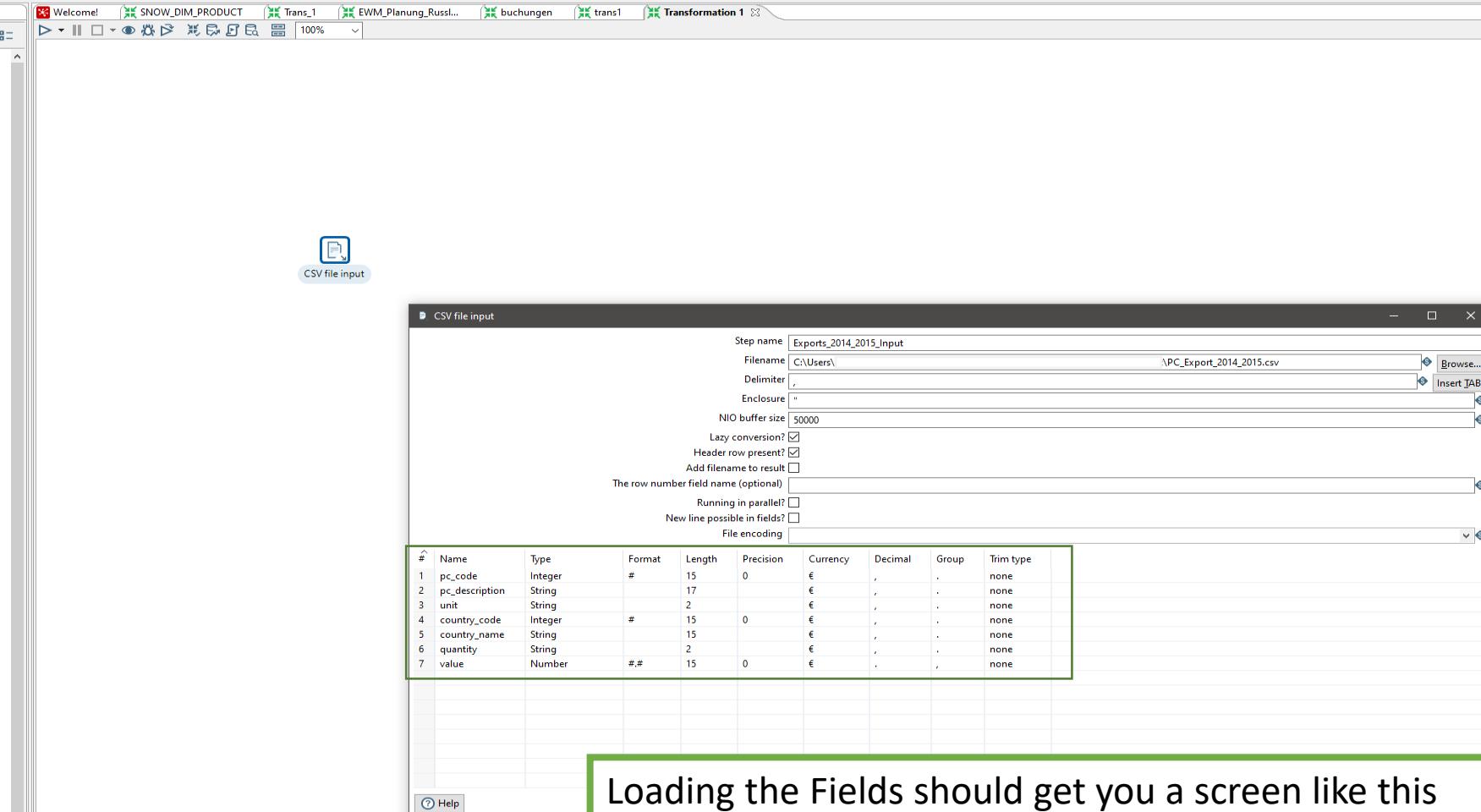
Joins

Data Warehouse

Validation

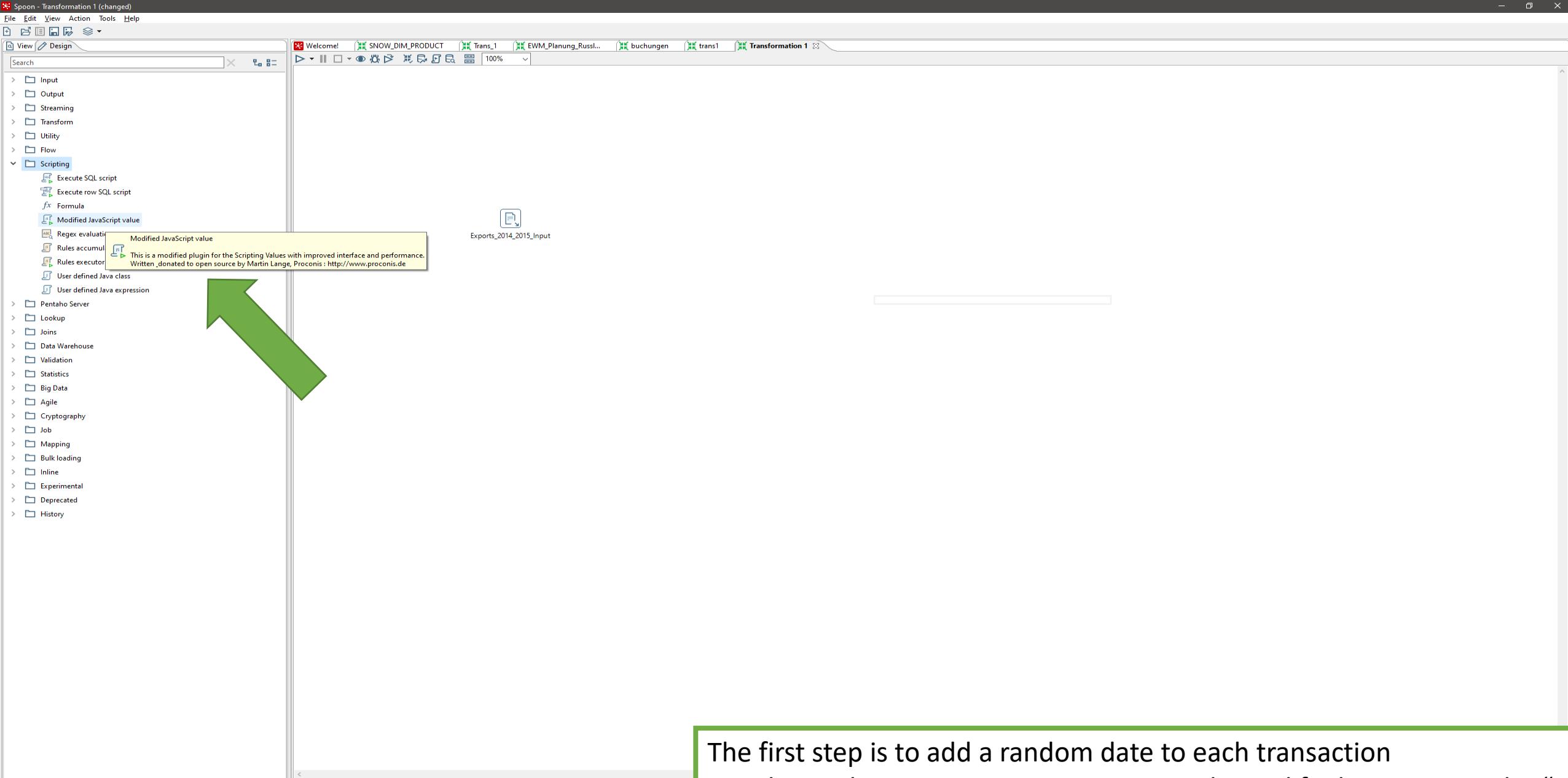
Statistics

Dim Data

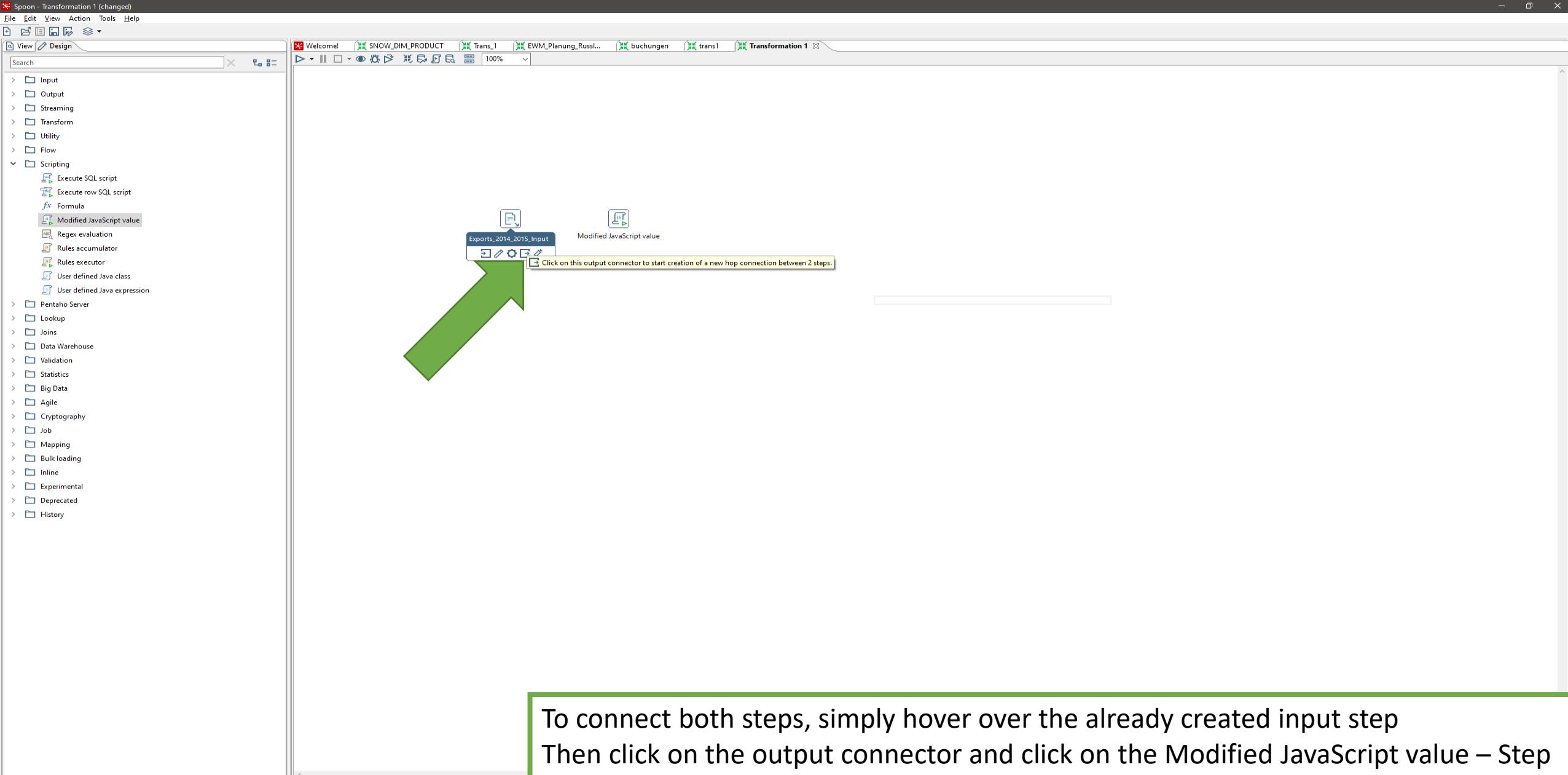


Loading the Fields should get you a screen like this  
Edit the datatype of the corresponding column to match the following:

pc_code	Product ID	A1	String
pc_description	Product Description	Tea	String
unit	Unit of quantity	Kgs	String
country_code	Country ID	1213	Integer
country_name	Country Description	Kenya	String
quantity	Count of Goods	2277547	Integer
value	Value in Mio. USD	4,276366	Number



The first step is to add a random date to each transaction  
To achieve this, we integrate a step named „Modified JavaScript value“  
This step can be found under Scripting > Modified JavaScript value  
Drag it into the main screen



To connect both steps, simply hover over the already created input step  
Then click on the output connector and click on the Modified JavaScript value – Step  
Select „Main Output of Step“  
This should create a link between both steps  
Finally click on Modified JavaScript value

Spoon - Transformation 1 (changed)

File Edit View Action Tools Help

View Design

Search

Input Output Streaming Transform Utility Flow Scripting

Execute SQL script Execute row SQL script Formula Modified JavaScript value Regex evaluation Rules accumulator Rules executor User defined Java class User defined Java expression

Pentaho Server Lookup Joins Data Warehouse Validation Statistics Big Data Agile Cryptography Job Mapping Bulk loading Inline Experimental Deprecated History

Welcome! SNOW\_DIM\_PRODUCT Trans\_1 EWM\_Planung\_Russl... buchungen

100%

Exports\_2014\_2015\_Input Modified JavaScript value

Java script functions :

- Transform Scripts
- Transform Constants
- Transform Functions
- Input fields
  - pc\_code
  - pc\_description
  - unit
  - country\_code
  - country\_name
  - quantity
  - value
- Output fields

Step name : Modified JavaScript value

Java script :

```
//Script here
```

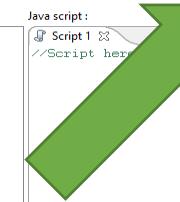
Please use the 'Replace value 'Fieldname' or 'Rename to'

Linien: 0 Compatibility mode? Optimization level: 9

Fields

#	Fieldname	Rename to	Type	Length	Precision	Replace value 'Fieldname' or 'Rename to'
1						

Help OK Cancel Get variables Test script



Again give the step a reasonable name  
We will now proceed with scripting the random date assignment

Spoon - Tuttrans (changed)  
File Edit View Action Tools Help

Modified JavaScript value

Step name Random\_Date\_Assignment\_Export\_2014\_2015

Java script functions :

Input fields

- pc\_code
- pc\_description
- unit
- country\_code
- country\_name
- quantity
- value

Output fields

Please use the 'Replace value 'Fieldname' or 'Rename To' field.

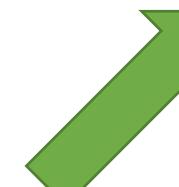
```
// Script for random assignation of month and year  
// d1 creates present date, sets its Day to the 1st, then randomizes its month  
// if month is less than April, assign 2015; if month is more April and higher assign 2014  
// reason is CSV-entries are aggregated from range April to March  
// d3 finally creates a string representing the adequate year + randomized month: "YYYYMM" <- Month_ID for DIM_Zeit  
  
var d1 = new Date();  
d1.setDate(1);  
d1.setMonth((Math.random()*12));  
var d2 = date2str(d1,"MM");  
if (str2num(d2) >= 4) {  
    d1.setFullYear(2014);  
} else {  
    d1.setFullYear(2015);  
}  
var d3 = date2str(d1,"yyyy") + d2;
```

Position: 17, 32  
Compatibility mode? Optimization level 9

Fields

#	Fieldname	Rename to	Type	Length	Precision	Replace value 'Fieldname' or 'Rename to'
1	d1		Date			N
2	d2		String			N
3	d3		String			N

OK Cancel Get variables Test script



The procedure of the script is described in the comment section above the code  
After typing the script a click on „Get variables“ should give you an output like this  
Change the datatype of d3 to Integer and click on „Test script“ and press OK

Spoon - Tuttrans (changed)  
 File Edit View Action Tools Help

View Design

Search

Input  
 Output  
 Streaming  
**Transform**  
 Add XML  
 Add a checksum  
 Add constants  
 Add sequence  
 Add value fields changing sequence  
 Calculator  
 Closure generator  
 Concat fields  
 Get ID from slave server  
 Number range  
 Replace in string  
 Row denormaliser  
 Row flattener  
 Row normaliser  
 Select values  
 Set field value  
 Set field value to a constant  
 Sort rows  
 Split field to rows  
 Split fields  
 String operations  
 Strings cut  
 Unique rows  
 Unique rows (HashSet)  
 Value mapper  
 XSL transformation  
 Utility  
 Flow  
 Scripting  
 Pentaho Server  
 Lookup  
 Joins  
 Data Warehouse  
 Validation  
 Statistics  
 Big Data  
 Agile  
 Cryptography  
 Job  
 Mapping  
 Bulk loading  
 Inline  
 Experimental  
 Deprecated  
 History

Modified JavaScript value

Step name Random\_Date\_Assignment\_Export\_2014\_2015

Java script functions:

- > Transform Scripts
- > Transform Constants
- > Transform Functions
- > Input fields
  - pc\_code
  - pc\_description
  - unit
  - country\_code
  - country\_name
  - quantity
  - value
- > Output fields
 

Please use the 'Replace value 'Fieldname' or 'Rename To' field.

Java script:

```

// Script for random assignation of month and year
// d1 creates present date, sets its Day to the 1st, then randomizes its month
// if month is less than April, assign 2015; if month is more April and higher assign 2014
// reason is CSV-entries are aggregated from range April to March
// d3 finally creates a string representing the adequate year + randomized month: "YYYYMM" <- Month_ID for DIM_Zeit

var d1 = new Date();
d1.setDate(1);
d1.setMonth((Math.random()*12));
var d2 = date2str(d1,"MM");
if (str2num(d2) >= 4) {
    d1.setFullYear(2014);
} else {
    d1.setFullYear(2015);
}
var d3 = date2str(d1,"yyyy") + d2;

```

Examine preview data

Rows of step: Random\_Date\_Assignment\_Export\_2014\_2015 (10 rows)

#	pc_code	pc_description	unit	country_code	country_name	quantity	value	d1	d2	d3
1	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/08/01 22:31:06.132	08	201408
2	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/09/01 22:31:06.133	09	201409
3	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2015/03/01 22:31:06.133	03	201503
4	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/05/01 22:31:06.133	05	201405
5	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2015/01/01 22:31:06.133	01	201501
6	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/04/01 22:31:06.133	05	201405
7	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/04/01 22:31:06.133	04	201404
8	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2015/01/01 22:31:06.133	01	201501
9	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2014/11/01 22:31:06.133	11	201411
10	test value test value	test value test value	test value test value	0	test value test value	test value test value	0	2015/01/01 22:31:06.133	01	201501

**Close** **Show Log**



The Script test should prompt something like this  
 Variable d3 is everything we are interested in for now  
 Press on the Close button

Spoon - Tuttrans (changed)  
File Edit View Action Tools Help

Modified JavaScript value

Step name Random\_Date\_Assignment\_Export\_2014\_2015

Java script functions :

Input fields

- pc\_code
- pc\_description
- unit
- country\_code
- country\_name
- quantity
- value

Output fields

Please use the 'Replace value 'Fieldname' or 'Rename To' field.

```
// Script for random assignation of month and year  
// d1 creates present date, sets its Day to the 1st, then randomizes its month  
// if month is less than April, assign 2015; if month is more April and higher assign 2014  
// reason is CSV-entries are aggregated from range April to March  
// d3 finally creates a string representing the adequate year + randomized month: "YYYYMM" <- Month_ID for DIM_Zeit  
  
var d1 = new Date();  
d1.setDate(1);  
d1.setMonth((Math.random()*12));  
var d2 = date2str(d1,"MM");  
if (str2num(d2) >= 4) {  
    d1.setFullYear(2014);  
} else {  
    d1.setFullYear(2015);  
}  
var d3 = date2str(d1,"yyyy") + d2;
```

Position: 17, 32  
Compatibility mode? Optimization level 9

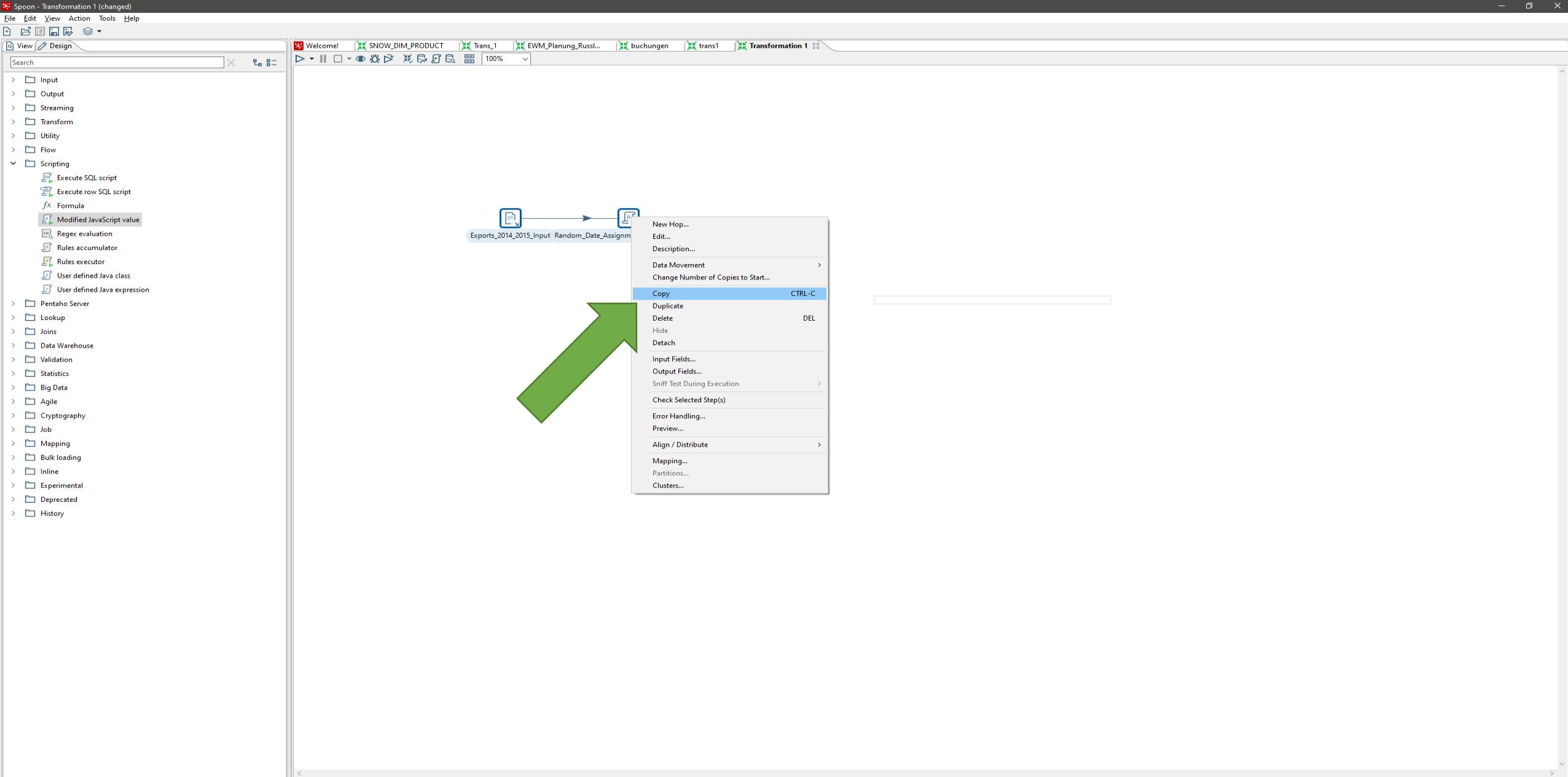
Fields

#	Fieldname	Rename to	Type	Length	Precision	Replace value 'Fieldname' or 'Rename to'
1	d1		Date			N
2	d2		String			N
3	d3		String			N

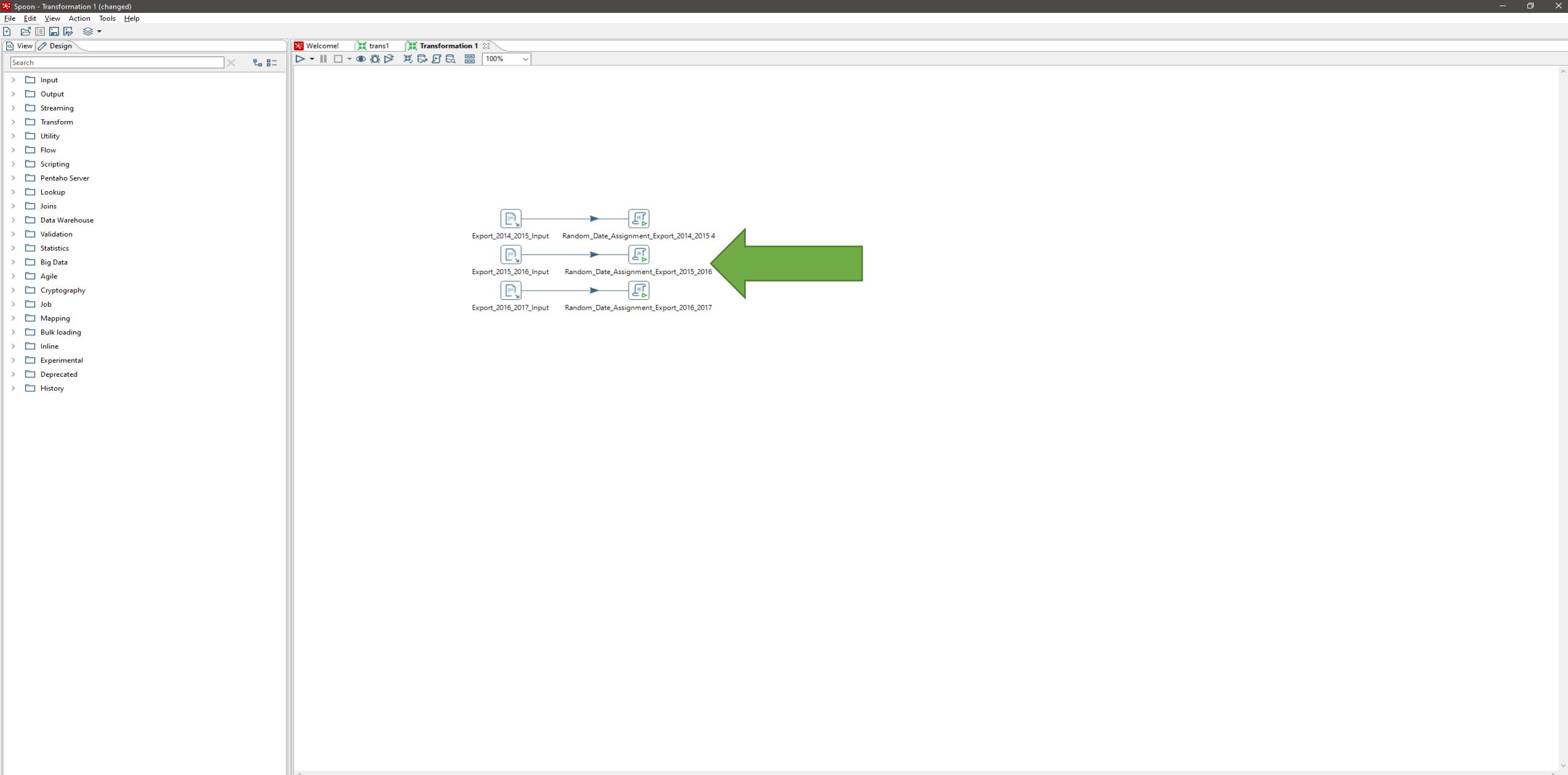
OK Cancel Get variables Test script



Delete both variables d1 and d2 since we are not interested in it by now  
Press OK



Now copy both steps and paste them to create inputs for the other two export files



Rearrange the steps and rename them according to the picture above  
Then double-click on „Random\_Date\_Assignment\_Export\_2015\_2016“

Spoon - Tuttrans (changed)  
 File Edit View Action Tools Help

View Design  
 Search

> Input  
 Output  
 Streaming  
 > Transform  
 Add XML  
 Add a checksum  
 Add constants  
 Add sequence  
 Add value fields changing sequence  
 Calculator  
 Closure generator  
 Concat fields  
 Get ID from slave server  
 Number range  
 Replace in string  
 Row denormaliser  
 Row flattener  
 Row normaliser  
 Select values  
 Set field value  
 Set field value to a constant  
 Sort rows  
 Split field to rows  
 Split fields  
 String operations  
 Strings cut  
 Unique rows  
 Unique rows (HashSet)  
 Value mapper  
 XSL transformation

> Utility  
 Flow  
 Scripting  
 Pentaho Server  
 Lookup  
 Joins  
 Data Warehouse  
 Validation  
 Statistics  
 Big Data  
 Agile  
 Cryptography  
 Job  
 Mapping  
 Bulk loading  
 Inline  
 Experimental  
 Deprecated  
 History

Welcome! trans1 Tuttrans  
 100%  
 Export\_2014\_2015\_Input Random\_Date\_Assignment\_Export\_2014  
 Export\_2015\_2016\_Input Random\_Date\_Assignment\_Export\_2015  
 Export\_2016\_2017\_Input Random\_Date\_Assignment\_Export\_2016

Modified JavaScript value  
 Step name: Random\_Date\_Assignment\_Export\_2015\_2016  
 Java script functions:  
 > Transform Scripts  
 > Transform Constants  
 > Transform Functions  
 > Input fields  
 pc\_code  
 pc\_description  
 unit  
 country\_code  
 country\_name  
 quantity  
 value  
 > Output fields  
 Please use the 'Replace value 'Fieldname''

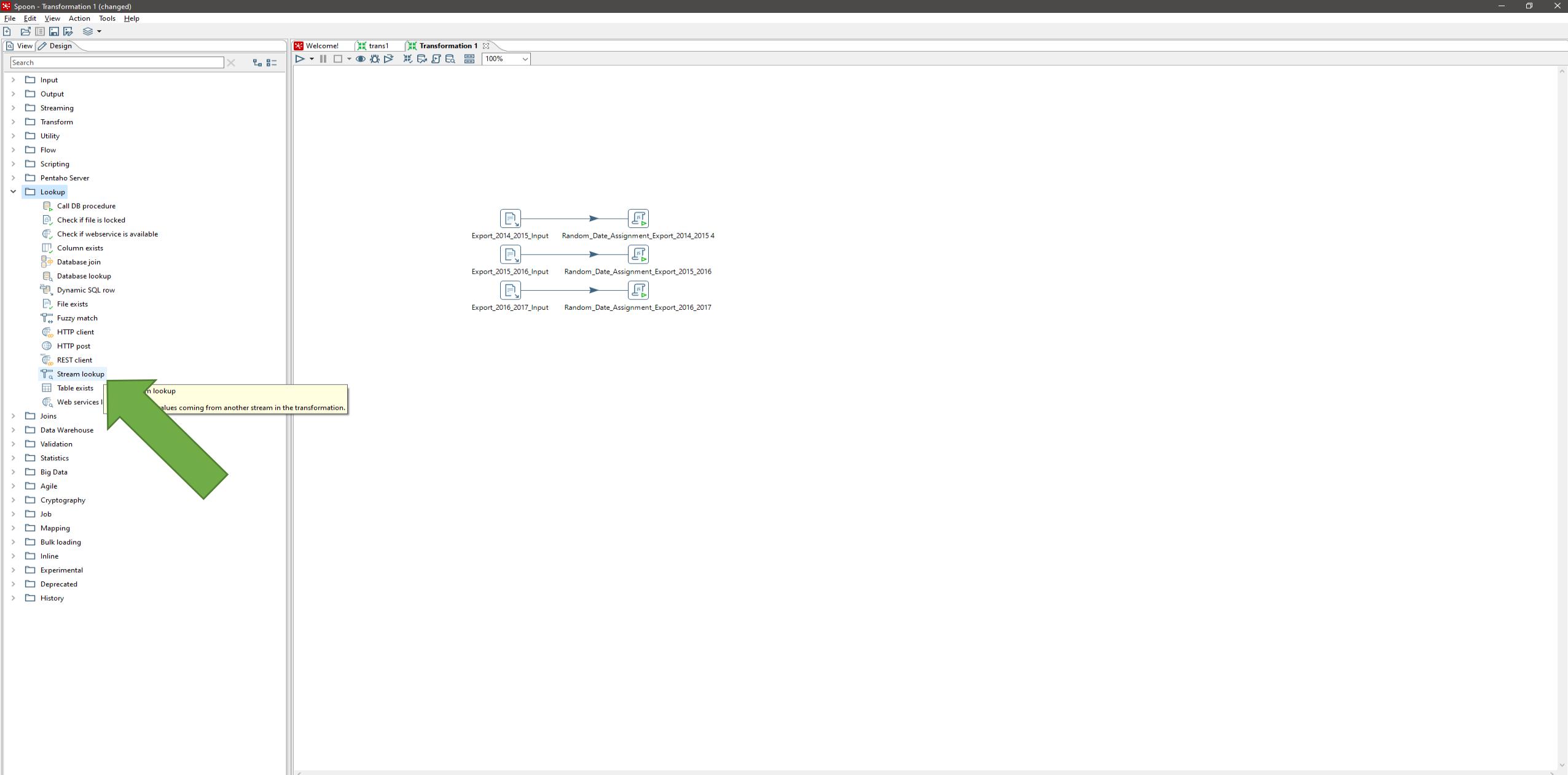
```

    Script 1
    // Script for random assignation of month and year
    // d1 creates present date, sets its Day to the 1st, then randomizes its month
    // if month is less than April, assign 2015; if month is more April and higher assign 2014
    // reason is CSV-entries are aggregated from range April to March
    // d3 finally creates a string representing the adequate year + randomized month: "YYYYMM" <- Month_ID for DIM_Zeit
    var d1 = new Date();
    d1.setDate(1);
    d1.setMonth(Math.random()*12);
    var d2 = date2str(d1, "MM");
    if (str2num(d2) >= 4) {
        d1.setFullYear(2014);
    } else {
        d1.setFullYear(2015);
    }
    var d3 = date2str(d1, "yyyy") + d2;
    
```

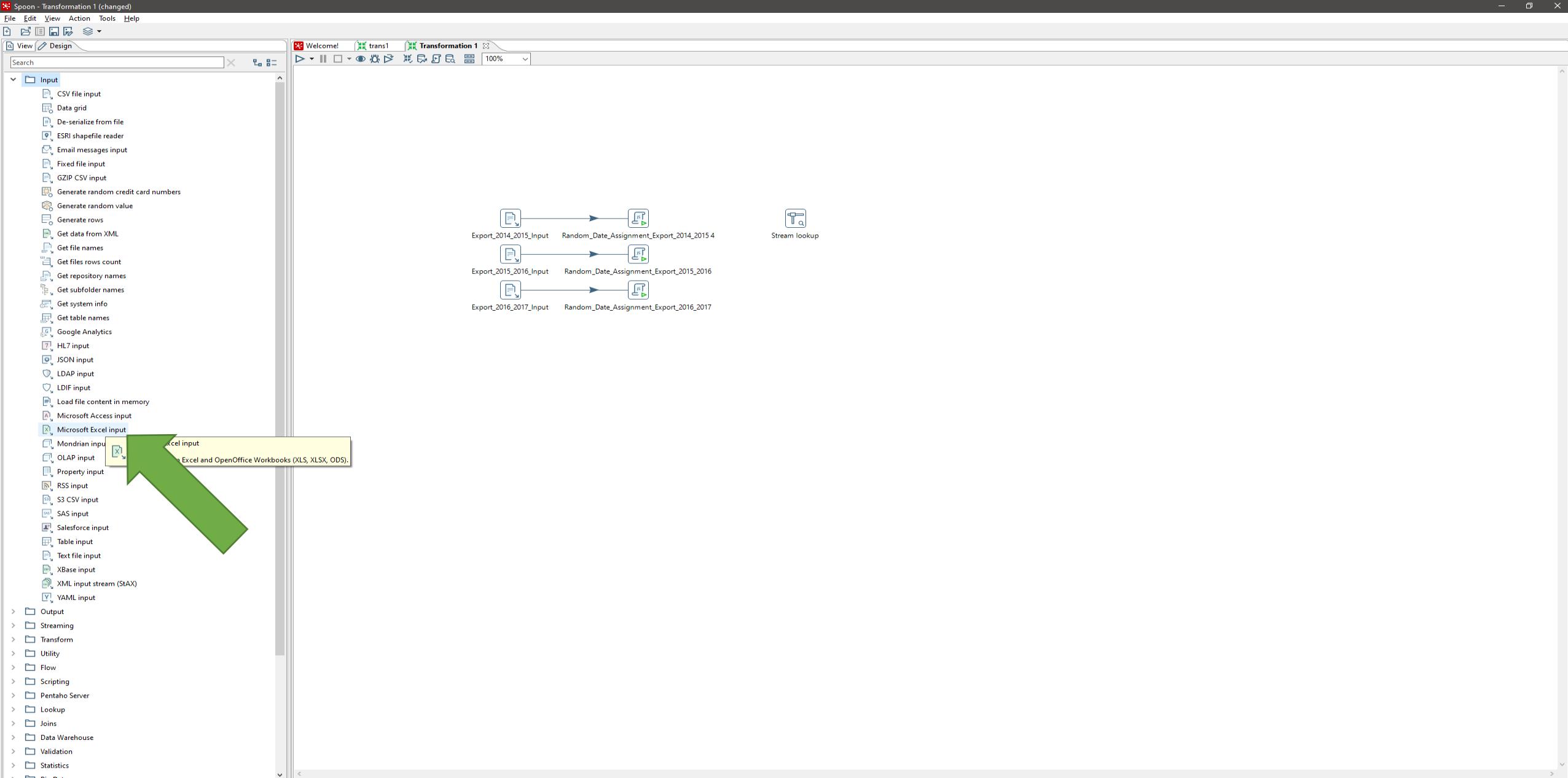
Execution Results  
 Logging Execution History Step Metrics Performance Graph Metrics Preview data  
 Stepname Copynr Read Written Input Output  
 1 DIM\_Product 0 0 168 168  
 2 Export\_2015\_2016\_Input 0 0 20771 20772  
 3 Export\_2014\_2015\_Input 0 0 20830 20831  
 4 Export\_2016\_2017\_Input 0 0 20830 20831  
 5 Random\_Date\_Assignment\_Export\_2014\_2015 0 20830 20830 0  
 6 Random\_Date\_Assignment\_Export\_2015\_2016 0 20771 20771 0  
 7 Random\_Date\_Assignment\_Export\_2016\_2017 0 20830 20830 0  
 8 DIM\_Region 0 0 233 233  
 9 DIM\_Time 0 0 1460 1460  
 10 Stream\_Lookup\_DIM\_Time\_Exports 0 63891 62431 0  
 11 Stream\_Lookup\_DIM\_Product\_Exports 0 62599 62431 0  
 12 Stream\_Lookup\_DIM\_Region\_Exports 0 62664 62431 0  
 13 Cleansing\_Exports 0 62431 62431 0  
 14 Output\_Exports 0 60370 60369 0 60369

Fields  
 Position: 17, 32  
 Compatibility mode? Optimization level: 9  
 # Fieldname Rename to Type Length Precision Replace value 'Fieldname' or 'Rename to'  
 1 d3 Integer N  
 OK Cancel Get variables Test script

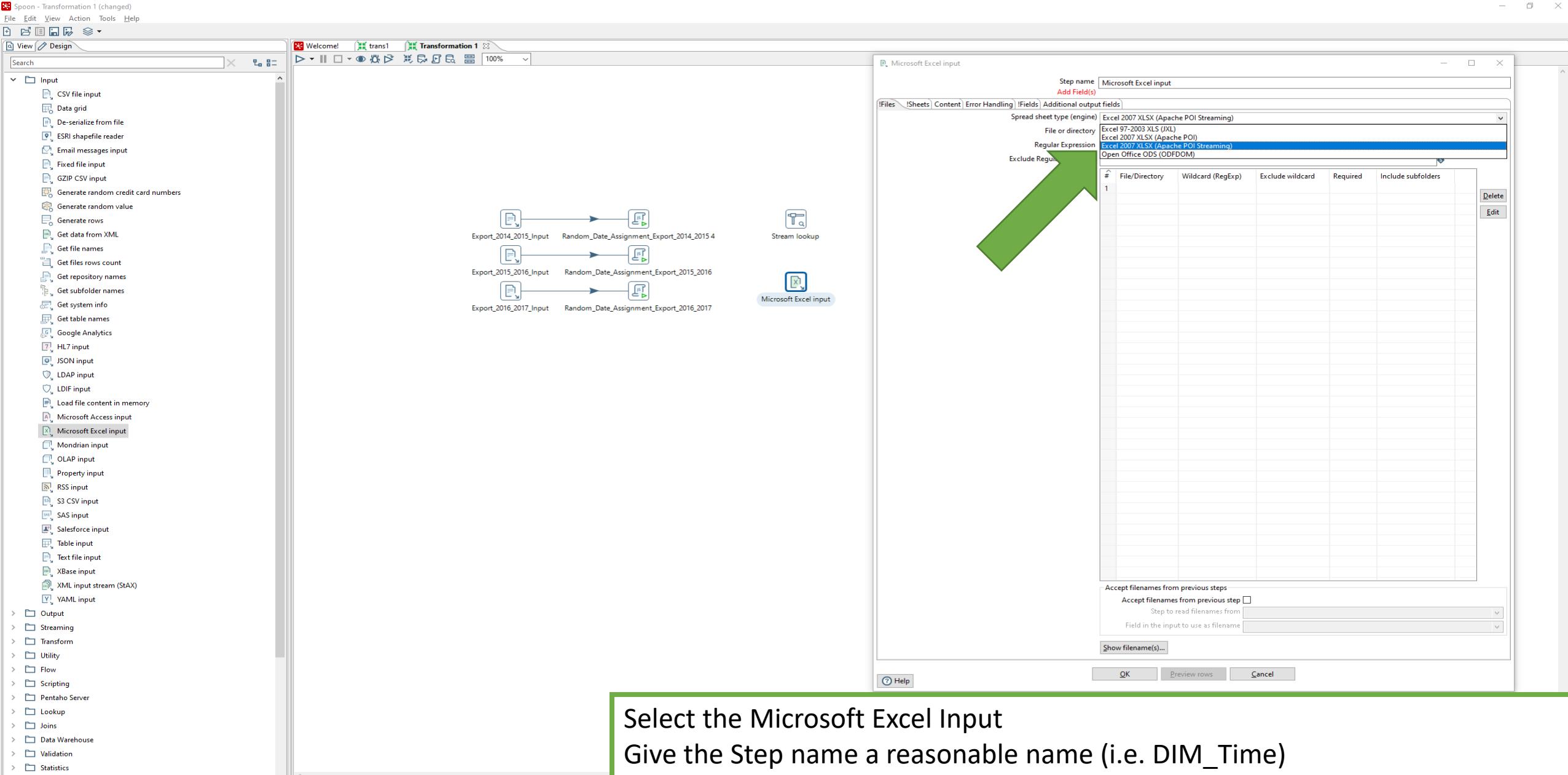
Adjust d1.setFullYear as described above and press OK  
 Do the same for „Random\_Date\_Assignment\_Export\_2016\_2017“



We will now create the first Lookup for our inputs to enrich dimensions  
Find Stream Lookup under Lookup and drag it onto the main screen



Since our Time dimension DIM\_Time is an Excel file, we additionally need a MS Excel Input Step  
Drag the MS Input Step onto the main screen



Select the Microsoft Excel Input  
Give the Step name a reasonable name (i.e. DIM\_Time)  
Select Spread sheet type and choose Excel 2007 XLSX (Apache POI Streaming)  
Click on „Browse...“ , search and choose DIM\_Time.xlsx  
Finally press „Add“

Spoon - Transformation 1 (changed)  
 File Edit View Action Tools Help

View Design  
 Search

**Input**

- CSV file input
- Data grid
- De-serialize from file
- ESRI shapefile reader
- Email messages input
- Fixed file input
- GZIP CSV input
- Generate random credit card numbers
- Generate random value
- Generate rows
- Get data from XML
- Get file names
- Get files rows count
- Get repository names
- Get subfolder names
- Get system info
- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input**
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input
- SAS input
- Salesforce input
- Table input
- Text file input
- XBase input
- XML input stream (StAX)
- YAML input

**Output**

**Streaming**

**Transform**

**Utility**

**Flow**

**Scripting**

Pentaho Server

Lookup

Joins

Data Warehouse

Validation

Statistics

DIM Data

Welcome! trans1 Transformation 1

Transformation 1

Stream lookup

Microsoft Excel input

Step name Microsoft Excel input  
**Add Field(s)**

Files Sheets Content Error Handling Fields Additional output fields

Spread sheet type (engine) Excel 2007 XLSX (Apache POI Streaming)

File or directory  Add Browse...

Regular Expression

Exclude Regular Expression

Selected files:

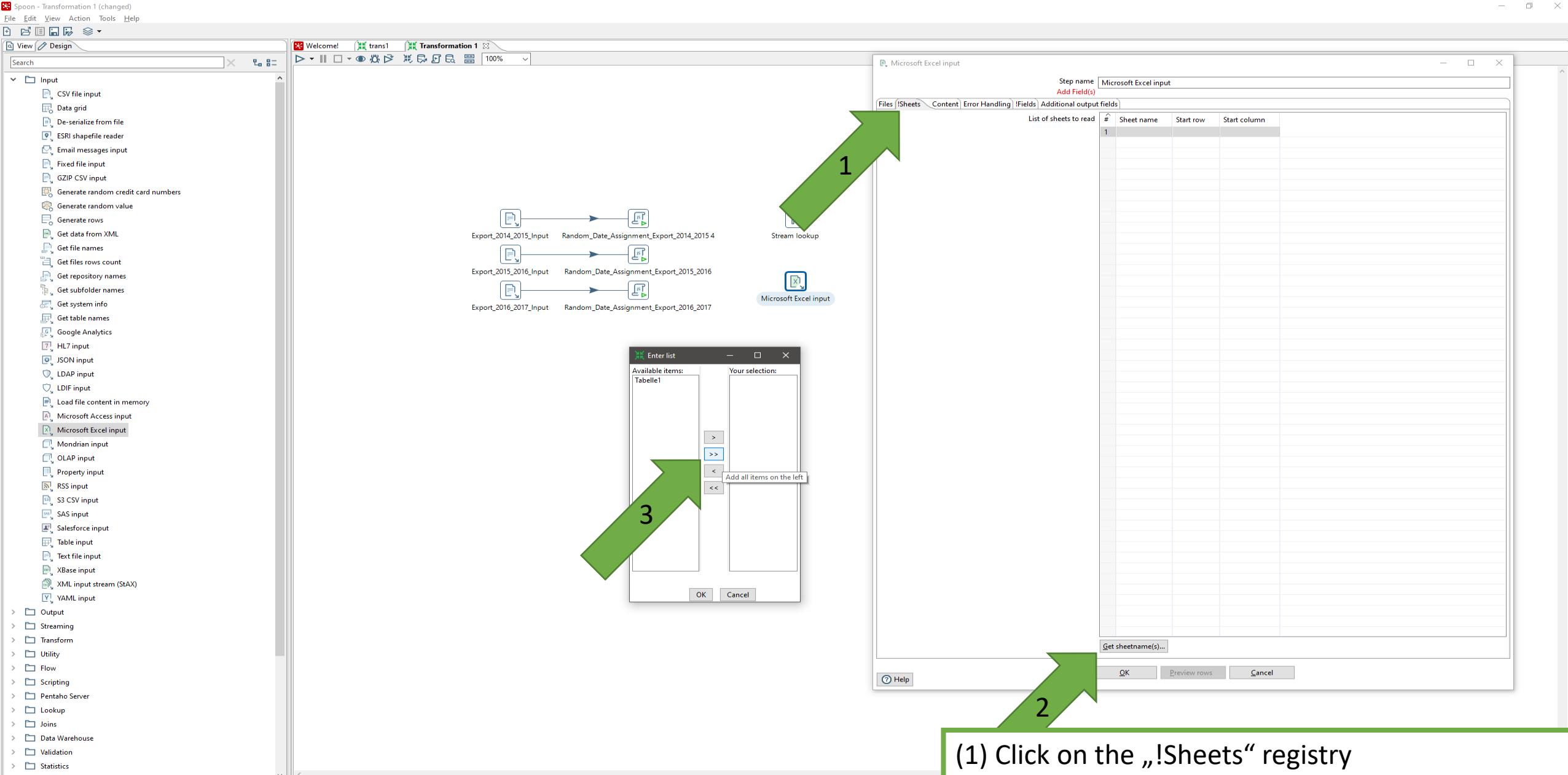
#	File/Directory	Wildcard
1	C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\DIM_Time.xlsx	

Delete Edit

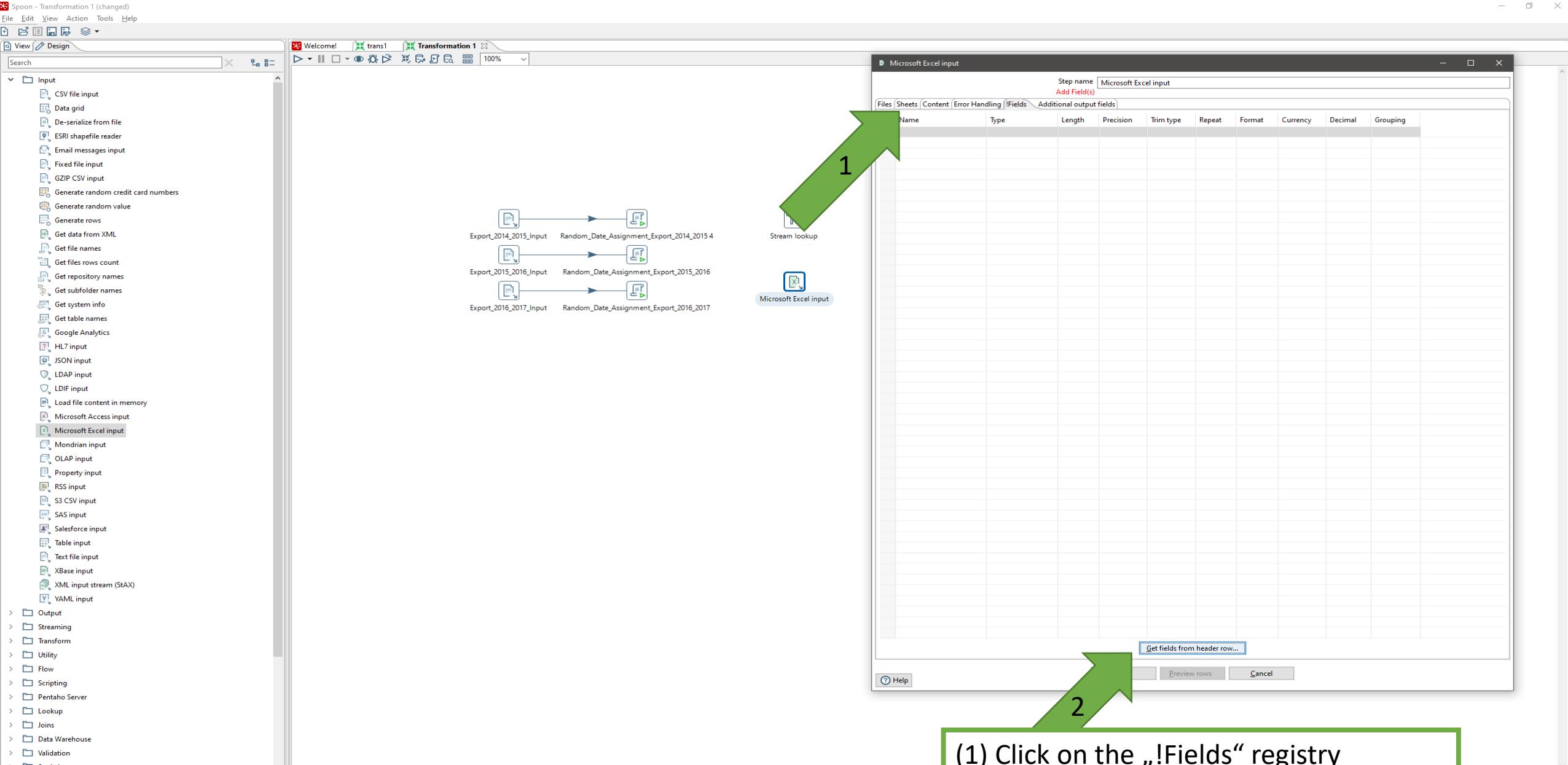
Accept filenames from previous steps   
 Step to read filenames from   
 Field in the input to use as filename

Show filename(s)... OK Preview rows Cancel

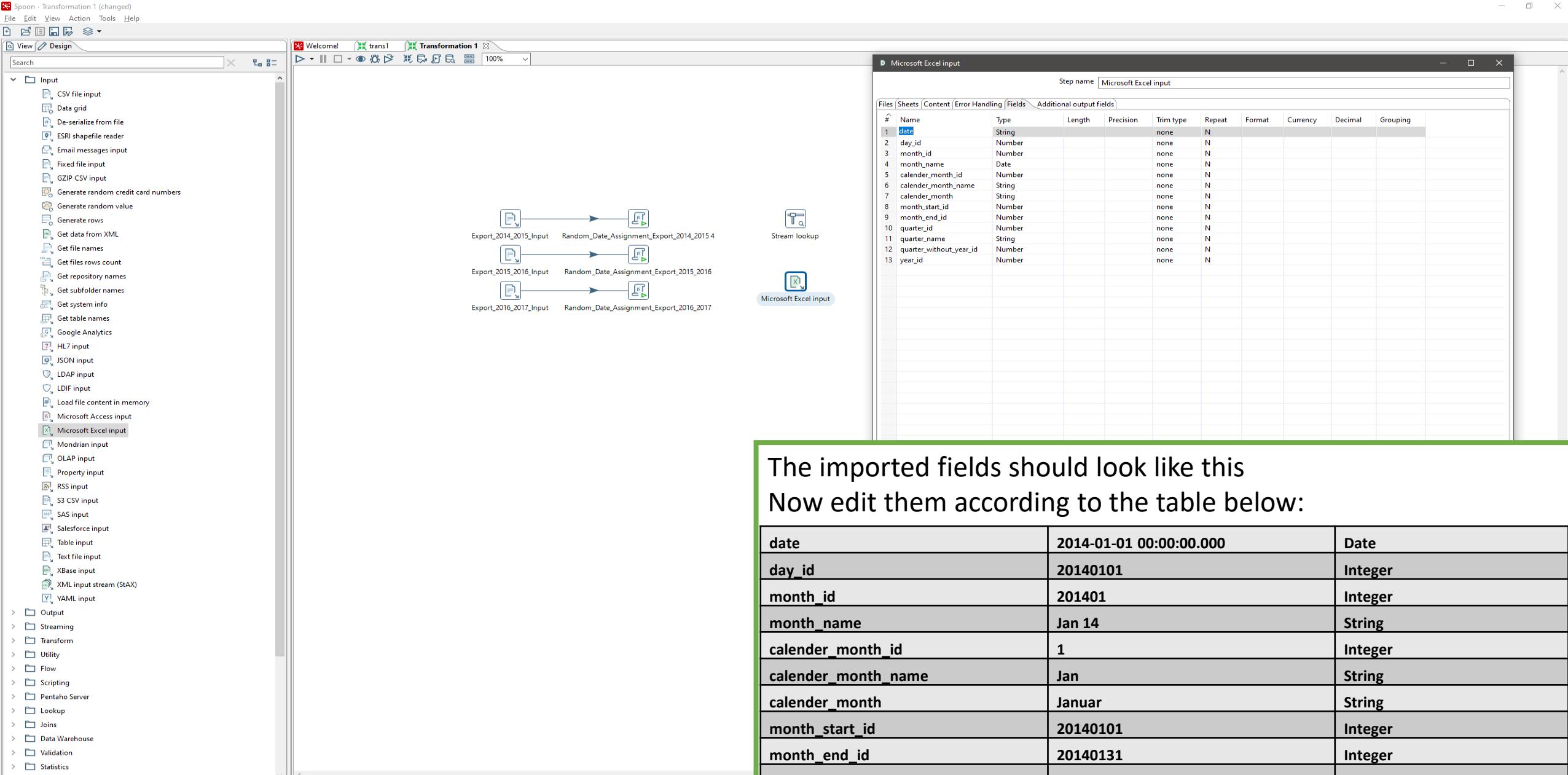
DIM\_Time.xlsx should be added to the selected files  
 Screen should look like stated above



- (1) Click on the „!Sheets“ registry  
 (2) Select „Get sheetnames...“  
 (3) Press on „>>“ to add all items to your selection  
 Press OK



- (1) Click on the „!Fields“ registry  
(2) Select „Get fields from header row...“



The imported fields should look like this  
Now edit them according to the table below:

date	2014-01-01 00:00:00.000	Date
day_id	20140101	Integer
month_id	201401	Integer
month_name	Jan 14	String
calender_month_id	1	Integer
calender_month_name	Jan	String
calender_month	Januar	String
month_start_id	20140101	Integer
month_end_id	20140131	Integer
quarter_id	20141	Integer
quarter_name	Q1 14	String
quarter_without_year_ID	1	Integer
year_id	2014	Integer

Spoon - Transformation 1 (changed)  
File Edit View Action Tools Help

Welcome! trans1 Transformation 1

View Design Search

Input

- CSV file input
- Data grid
- De-serialize from file
- ESRI shapefile reader
- Email messages input
- Fixed file input
- GZIP CSV input
- Generate random credit card numbers
- Generate random value
- Generate rows
- Get data from XML
- Get file names
- Get files rows count
- Get repository names
- Get subfolder names
- Get system info
- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input
- SAS input
- Salesforce input
- Table input
- Text file input
- XBase input
- XML input stream (StAX)
- YAML input

Output

Streaming

Transform

Utility

Flow

Scripting

Pentaho Server

Lookup

Join

Data Warehouse

Validation

Statistics

Bin Data

Transformation 1

Export\_2014\_2015\_Input → Random\_Date\_Assignment\_Export\_2014\_2015 4

Export\_2015\_2016\_Input → Random\_Date\_Assignment\_Export\_2015\_2016

Export\_2016\_2017\_Input → Random\_Date\_Assignment\_Export\_2016\_2017

Stream lookup

Microsoft Excel input

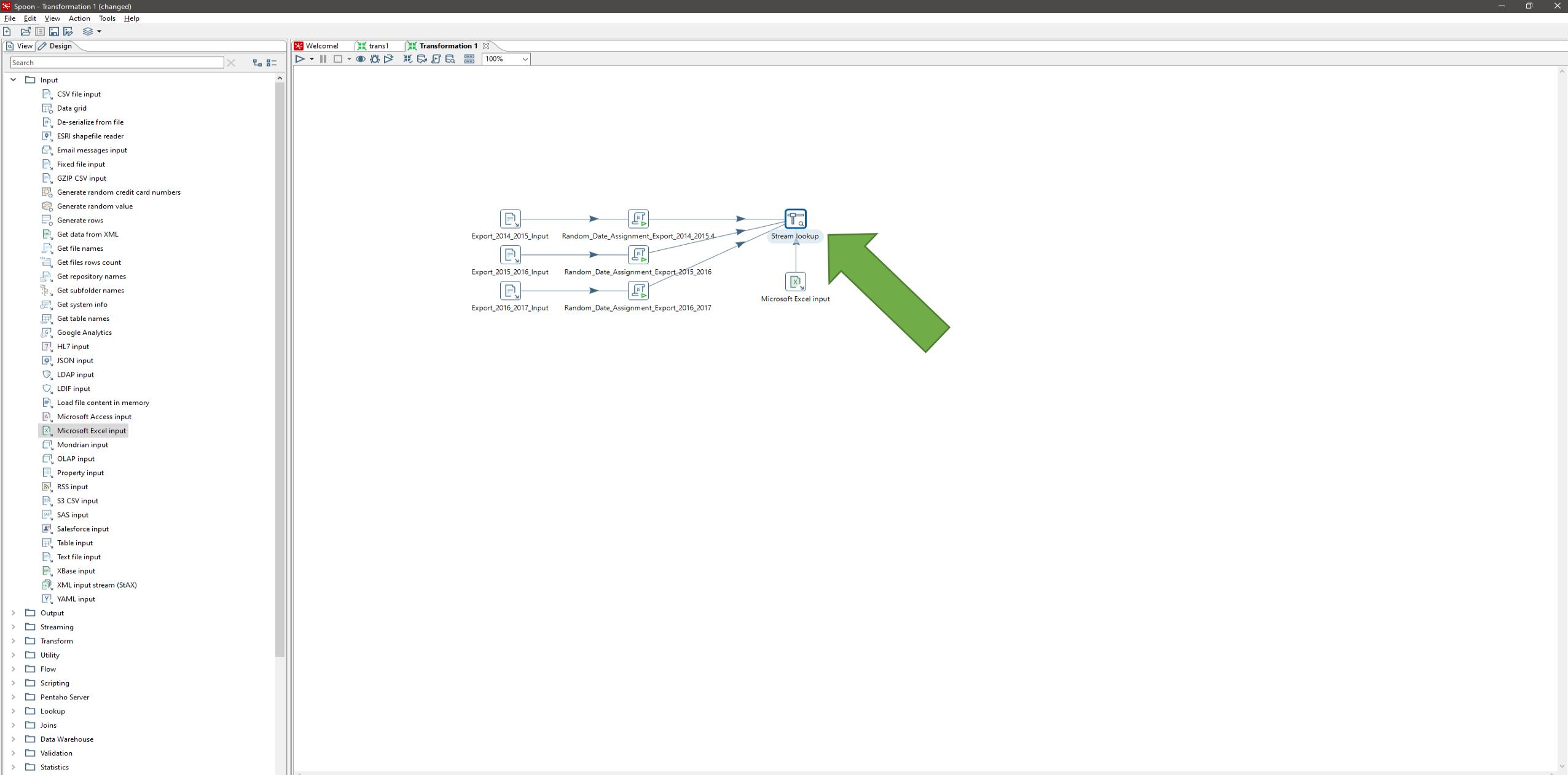
Step name Microsoft Excel input

Files Sheets Content Error Handling Fields Additional output fields

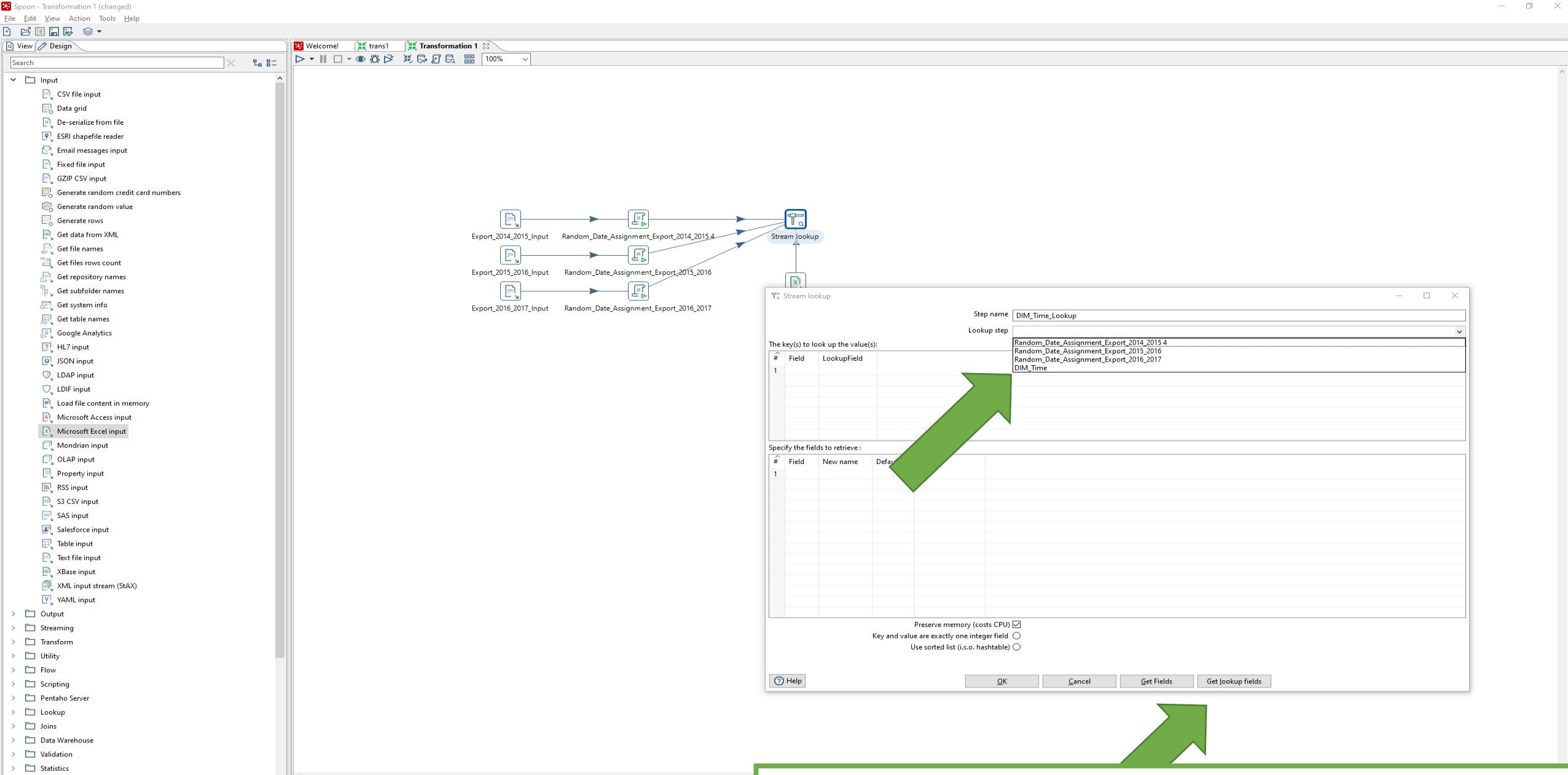
#	Name	Type	Length	Precision	Trim type	Repeat	Format	Currency	Decimal	Grouping
1	date	Date			none	N	yyyy-M...			
2	day_id	Integer			none	N				
3	month_id	Integer			none	N				
4	month_name	String			none	N				
5	calender_month_id	Integer			none	N				
6	calender_month_name	String			none	N				
7	calender_month	String			none	N				
8	month_start_id	Integer			none	N				
9	month_end_id	Integer			none	N				
10	quarter_id	Integer			none	N				
11	quarter_name	String			none	N				
12	quarter_without_year_id	Integer			none	N				
13	year_id	Integer			none	N				

Get fields from header row... OK Preview rows Cancel

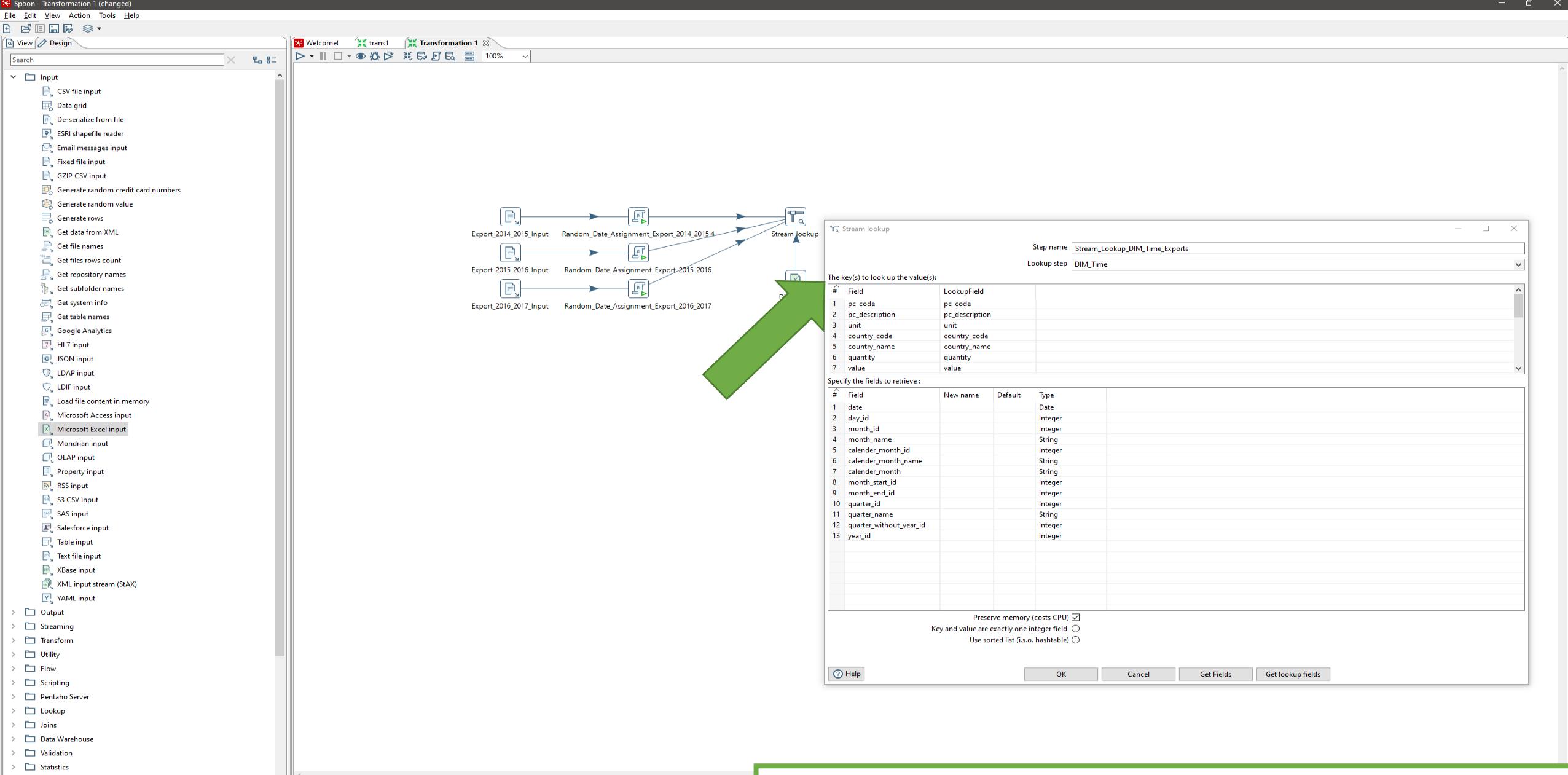
Do not forget to adjust the format for the date to:  
„yyyy-MM-dd HH:mm:ss.SSS“  
Press OK



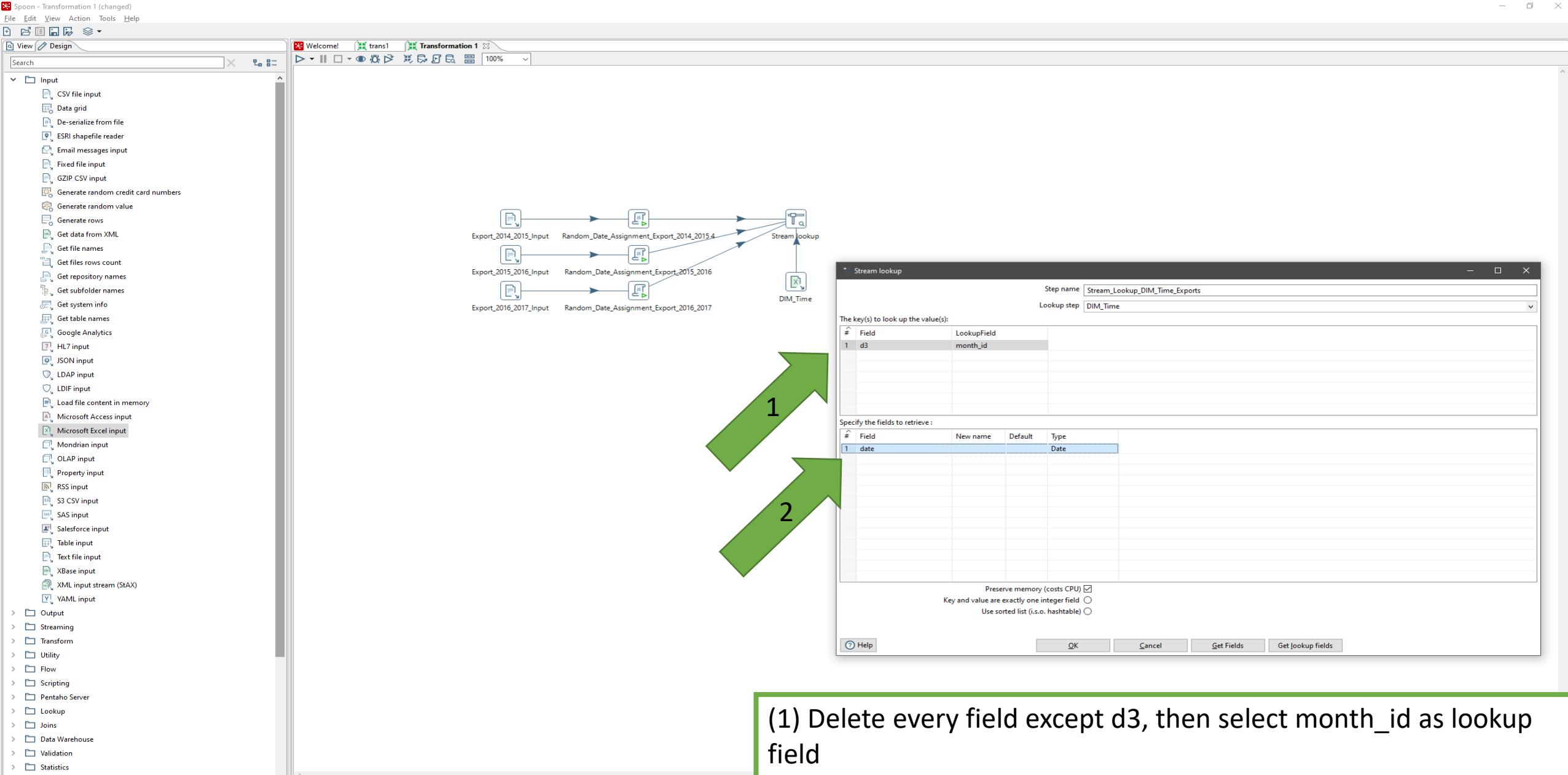
Now connect the steps as shown in the picture above  
Then double-click „Stream lookup“ to edit it



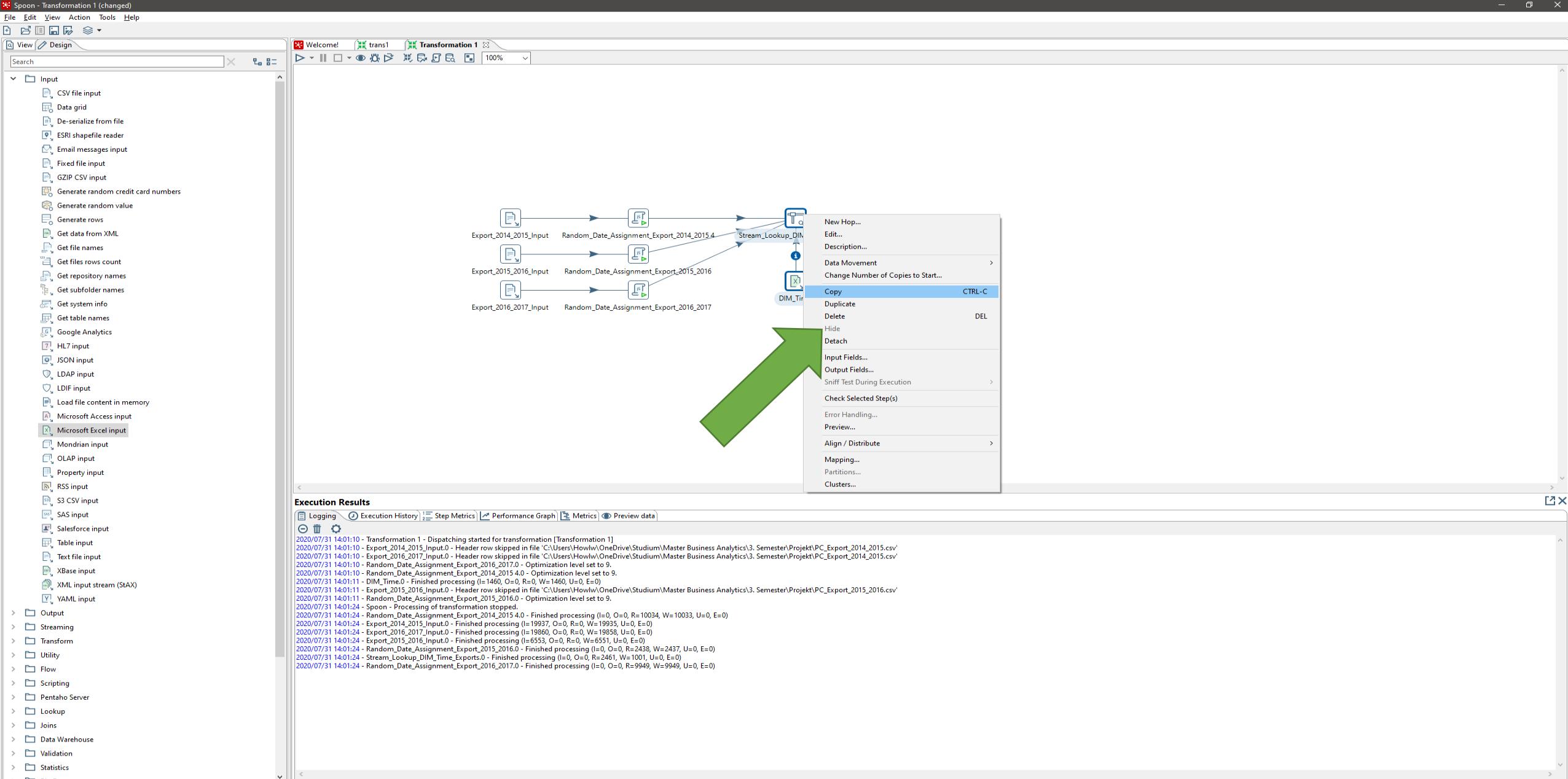
Name the step name properly  
Select the Lookup step which is „DIM\_Time“ in our case  
Select „Get Fields“ and „Get lookup fields“



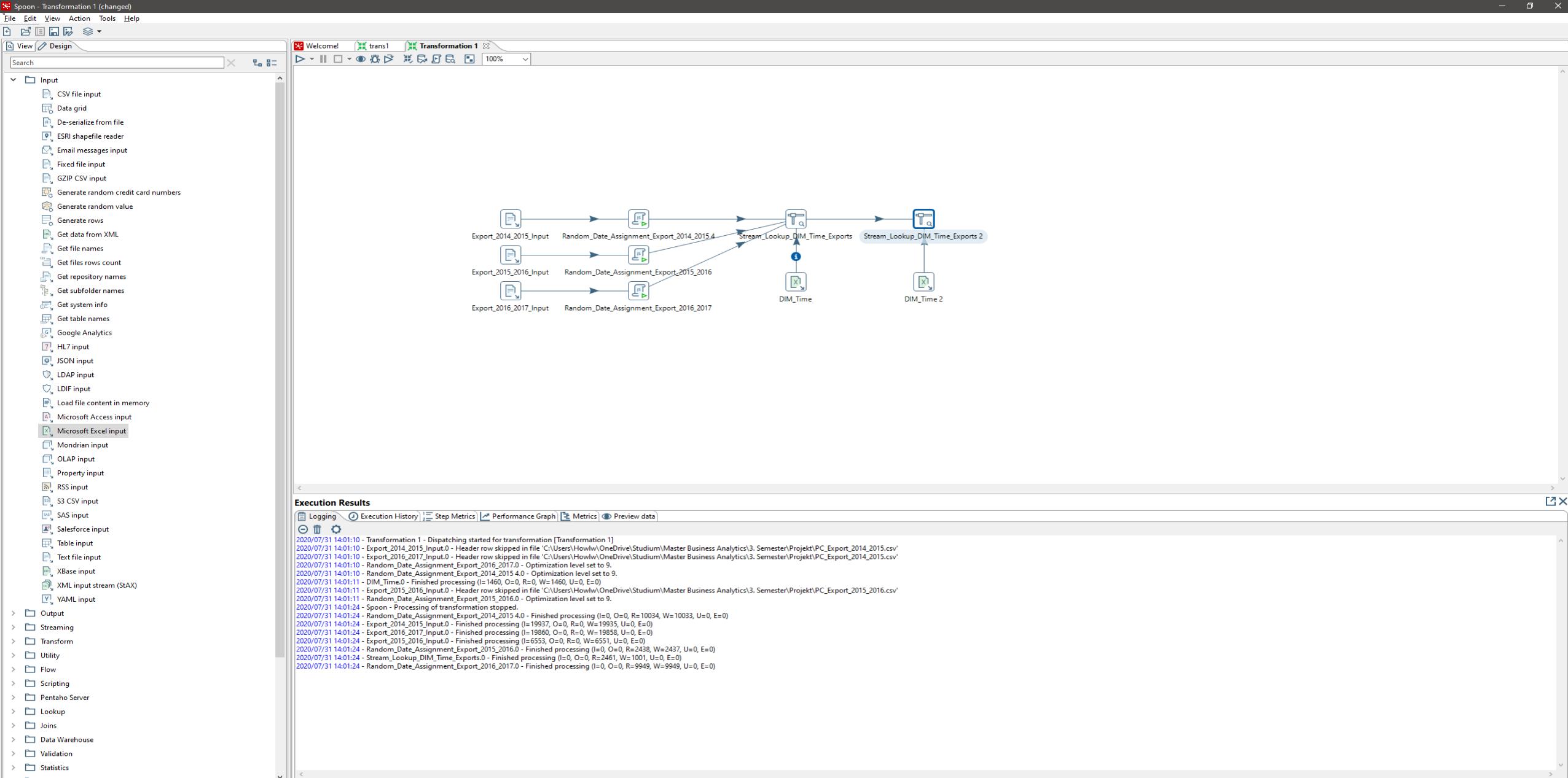
The screen should now look like this  
Since we look up the columns of DIM\_Time, we already know that d3 is the primary key to connect with DIM\_Time



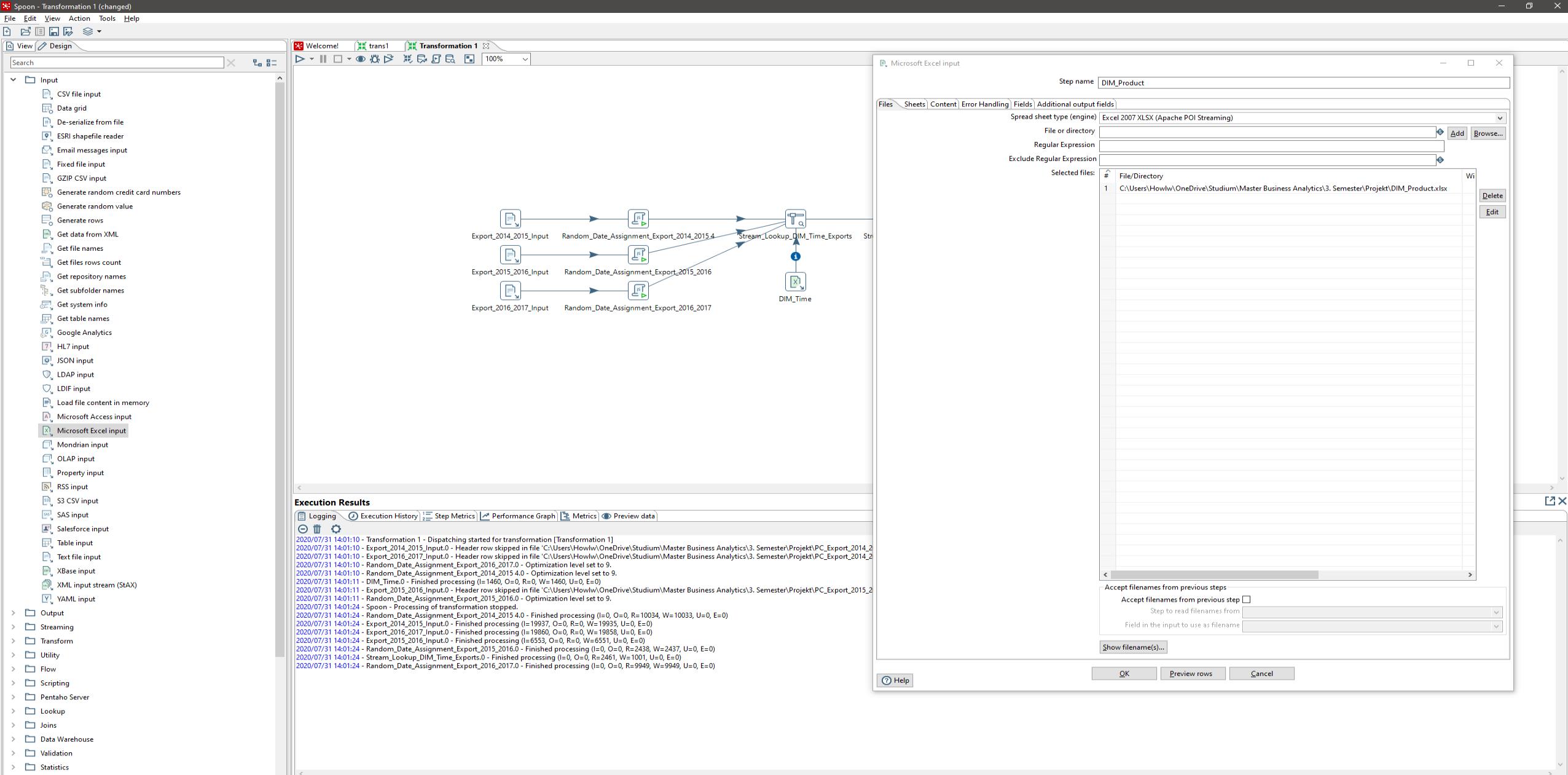
- (1) Delete every field except d3, then select month\_id as lookup field
- (2) Since we need dates to visualize timelines in Tableau, this field is the only one we are interested in: Delete every field except date in the fields to retrieve



Copy and Paste both the lookup step and the Excel Input Step right next to the first ones



Connect the lookup steps and double-click the newly created Excel Input Step to set up DIM\_Product



Rename it properly  
Select „Browse...“ to select DIM\_Product.xlsx and add it

Spoon - Transformation 1 (changed)

File Edit View Action Tools Help

View Design

Search

**Input**

- CSV file input
- Data grid
- De-serialize from file
- ESRI shapefile reader
- Email messages input
- Fixed file input
- GZIP CSV input
- Generate random credit card numbers
- Generate random value
- Generate rows
- Get data from XML
- Get file names
- Get files rows count
- Get repository names
- Get subfolder names
- Get system info
- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input**
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input
- SAS input
- Salesforce input
- Table input
- Text file input
- XBase input
- XML input stream (StAX)
- YAML input

**Output**

Streaming

Transform

Utility

Flow

Scripting

Pentaho Server

Lookup

Joins

Data Warehouse

Validation

Statistics

Bin Data

Welcome! trans1 Transformation 1

100%

Stream

Time\_Exports Stream

DIM\_Time

Export\_2014\_2015\_Input Random\_Date\_Assignment\_Export\_2014\_2015\_4

Export\_2015\_2016\_Input Random\_Date\_Assignment\_Export\_2015\_2016

Export\_2016\_2017\_Input Random\_Date\_Assignment\_Export\_2016\_2017

**Microsoft Excel input**

Step name: DIM\_Product

Files Sheets Content Error Handling Fields Additional output fields

#	Name	Type	Length	Precision	Trim type	Repeat	Format	Currency	Decimal	Grouping
1	pc_code	String			none	N				
2	pc_description	String			none	N				
3	pc_category	String			none	N				

Execution Results

Logging Execution History Step Metrics Performance Graph Metrics Preview data

```

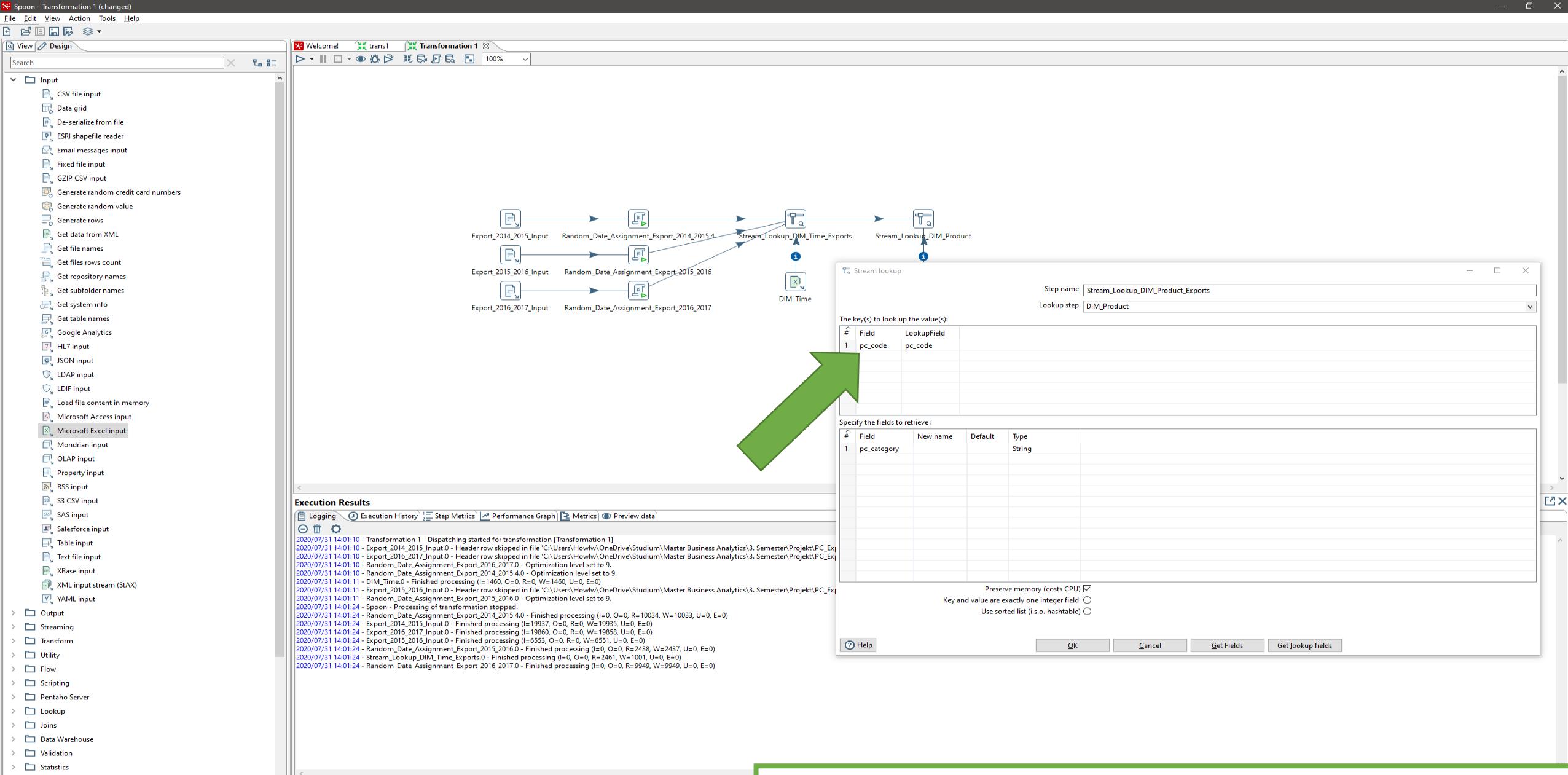
2020/07/31 14:01:10 - Transformation 1 - Dispatching started for transformation [Transformation 1]
2020/07/31 14:01:10 - Export_2014_2015_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_Export_2014_2015.xlsx'
2020/07/31 14:01:10 - Export_2016_2017_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_export_2014_2016.xlsx'
2020/07/31 14:01:10 - Random_Date_Assignment_Export_2014_2015_4.0 - Optimization level set to 9.
2020/07/31 14:01:10 - Random_Date_Assignment_Export_2015_2016.0 - Optimization level set to 9.
2020/07/31 14:01:10 - Random_Date_Assignment_Export_2016_2017.0 - Optimization level set to 9.
2020/07/31 14:01:11 - DIM_Time.0 - Finished processing (I=1460, O=0, R=0, W=1460, U=0, E=0)
2020/07/31 14:01:11 - Export_2015_2016_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_Export_2015_2016.xlsx'
2020/07/31 14:01:11 - Random_Date_Assignment_Export_2015_2016_0.0 - Optimization level set to 9.
2020/07/31 14:01:24 - Spoon - Processing of transformation stopped.
2020/07/31 14:01:24 - Random_Date_Assignment_Export_2014_2015_4.0 - Finished processing (I=0, O=0, R=10034, W=10033, U=0, E=0)
2020/07/31 14:01:24 - Export_2014_2015_Input.0 - Finished processing (I=19937, O=0, R=0, W=19935, U=0, E=0)
2020/07/31 14:01:24 - Export_2016_2017_Input.0 - Finished processing (I=19860, O=0, R=0, W=19858, U=0, E=0)
2020/07/31 14:01:24 - Random_Date_Assignment_Export_2015_2016_0.0 - Finished processing (I=6553, O=0, R=0, W=6551, U=0, E=0)
2020/07/31 14:01:24 - Random_Date_Assignment_Export_2015_2016_0.0 - Finished processing (I=0, O=0, R=2438, W=2437, U=0, E=0)
2020/07/31 14:01:24 - Stream_Lookup_DIM_Time_Exports.0 - Finished processing (I=0, O=0, R=2461, W=1001, U=0, E=0)
2020/07/31 14:01:24 - Random_Date_Assignment_Export_2016_2017_0.0 - Finished processing (I=0, O=0, R=9949, W=9949, U=0, E=0)

```

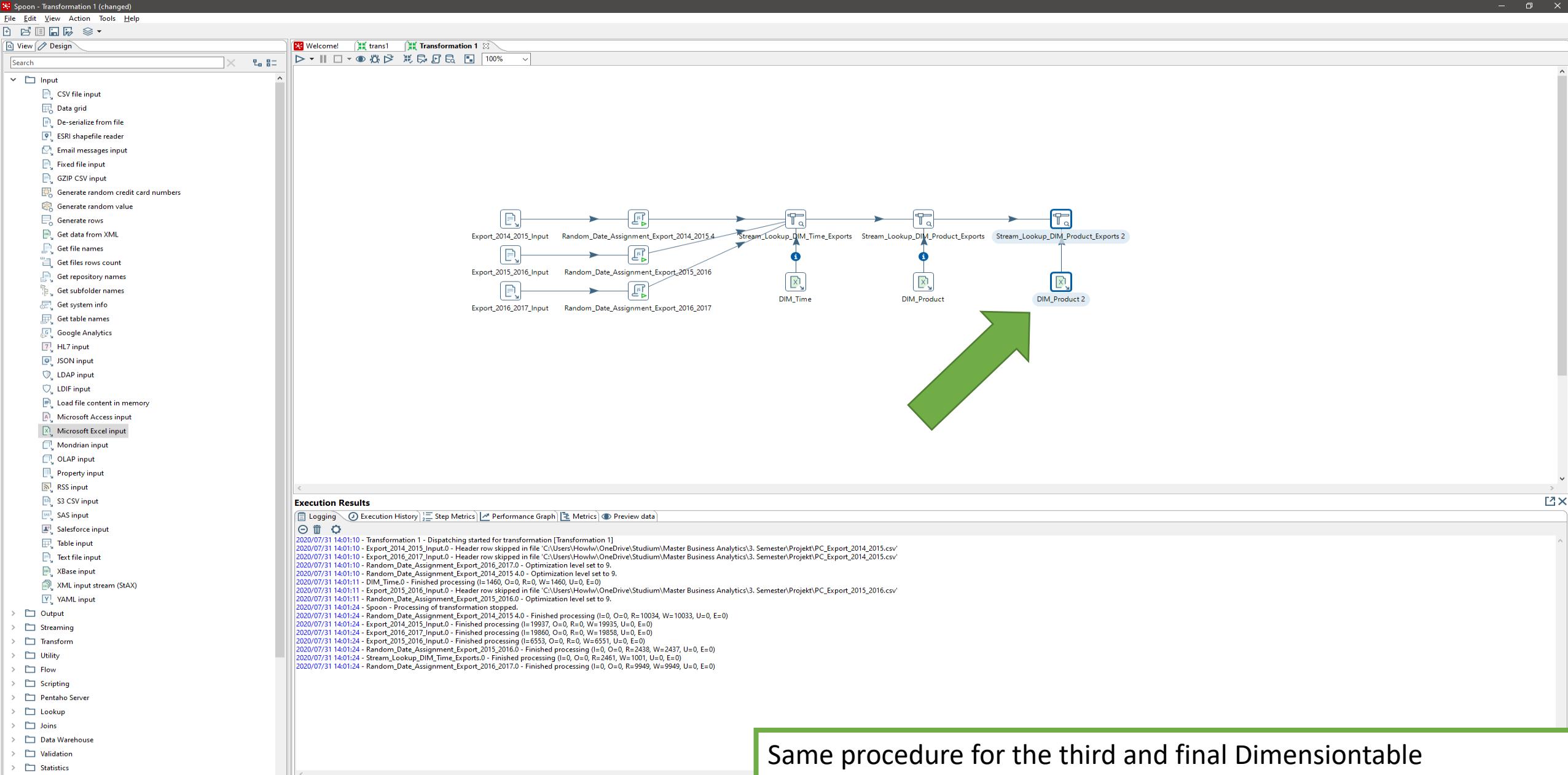
Get fields from header row...

OK Preview rows Cancel

Select Fields in the registry and choose „Get fields from header row...“  
The columns of DIM\_Product.xlsx should now be loaded  
Press OK



Double-click on the Lookup Step for DIM\_Product and change it according to the picture above  
Press OK



Same procedure for the third and final Dimensiontable integration  
 Again copy and paste both Stream Lookup and Microsoft Excel Input, connect the steps and then edit them

Spoon - Transformation 1 (changed)

View Design

Search

Input

- CSV file input
- Data grid
- De-serialize from file
- ESRI shapefile reader
- Email messages input
- Fixed file input
- GZIP CSV input
- Generate random credit card numbers
- Generate random value
- Generate rows
- Get data from XML
- Get file names
- Get files rows count
- Get repository names
- Get subfolder names
- Get system info
- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input
- SAS input
- Salesforce input
- Table input
- Text file input
- XBase input
- XML input stream (StAX)
- YAML input

Output

Streaming

Transform

Utility

Flow

Scripting

Pentaho Server

Lookup

Joins

Data Warehouse

Validation

Statistics

Bin Data

Welcome! trans1 Transformation 1

Transformation 1

100%

Execution Results

Logging Execution History Step Metrics Performance Graph Metrics Preview data

2020/07/31 14:01:10 - Transformation 1 - Dispatching started for transformation [Transformation 1]  
2020/07/31 14:01:10 - Export\_2014\_2015\_Input0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_Export\_2014\_2015.xlsx'  
2020/07/31 14:01:10 - Export\_2016\_2017\_Input0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_export\_2016\_2017.xlsx'  
2020/07/31 14:01:10 - Random\_Date\_Assignment\_Export\_2014\_2015\_4 - Optimization level set to 9.  
2020/07/31 14:01:10 - Random\_Date\_Assignment\_Export\_2015\_2016\_0 - Optimization level set to 9.  
2020/07/31 14:01:10 - Random\_Date\_Assignment\_Export\_2016\_2017\_0 - Optimization level set to 9.  
2020/07/31 14:01:11 - Stream\_Lookup\_DIM\_Time\_Exports0 - Finished processing (I=1460, O=0, R=0, W=1460, U=0, E=0)  
2020/07/31 14:01:11 - Export\_2015\_2016\_Input0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_Export\_2015\_2016.xlsx'  
2020/07/31 14:01:11 - Random\_Date\_Assignment\_Export\_2015\_2016\_0 - Optimization level set to 9.  
2020/07/31 14:01:24 - Spoon - Processing of transformation stopped.  
2020/07/31 14:01:24 - Random\_Date\_Assignment\_Export\_2014\_2015\_4 - Finished processing (I=0, O=0, R=10034, W=10033, U=0, E=0)  
2020/07/31 14:01:24 - Export\_2014\_2015\_Input0 - Finished processing (I=19937, O=0, R=0, W=19935, U=0, E=0)  
2020/07/31 14:01:24 - Export\_2016\_2017\_Input0 - Finished processing (I=19860, O=0, R=0, W=19858, U=0, E=0)  
2020/07/31 14:01:24 - Export\_2015\_2016\_Input0 - Finished processing (I=6553, O=0, R=0, W=6551, U=0, E=0)  
2020/07/31 14:01:24 - Random\_Date\_Assignment\_Export\_2015\_2016\_0 - Finished processing (I=0, O=0, R=2438, W=2437, U=0, E=0)  
2020/07/31 14:01:24 - Stream\_Lookup\_DIM\_Time\_Exports0 - Finished processing (I=0, O=0, R=2461, W=1001, U=0, E=0)  
2020/07/31 14:01:24 - Random\_Date\_Assignment\_Export\_2016\_2017\_0 - Finished processing (I=0, O=0, R=9949, W=9949, U=0, E=0)

Microsoft Excel input

Step name: DIM\_Region

Spread sheet type (engine): Excel 2007 XLSX (Apache POI Streaming)

File or directory:

Regular Expression:

Exclude Regular Expression:

Selected files:

#	File/Directory
1	C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\DIM_Region.xlsx

Accept filenames from previous steps   
Step to read filenames from   
Field in the input to use as filename

Show filename(s)...

OK Preview rows Cancel

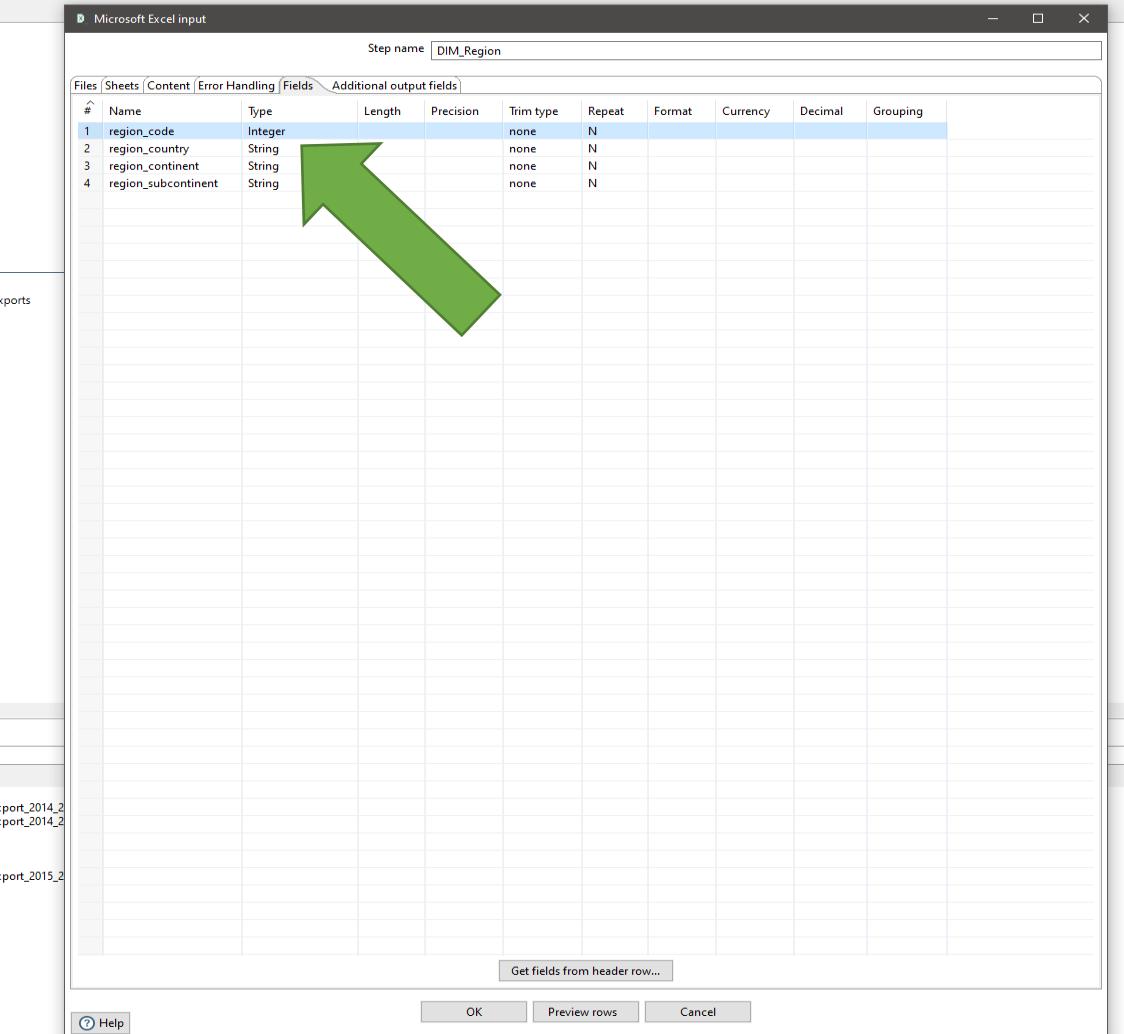
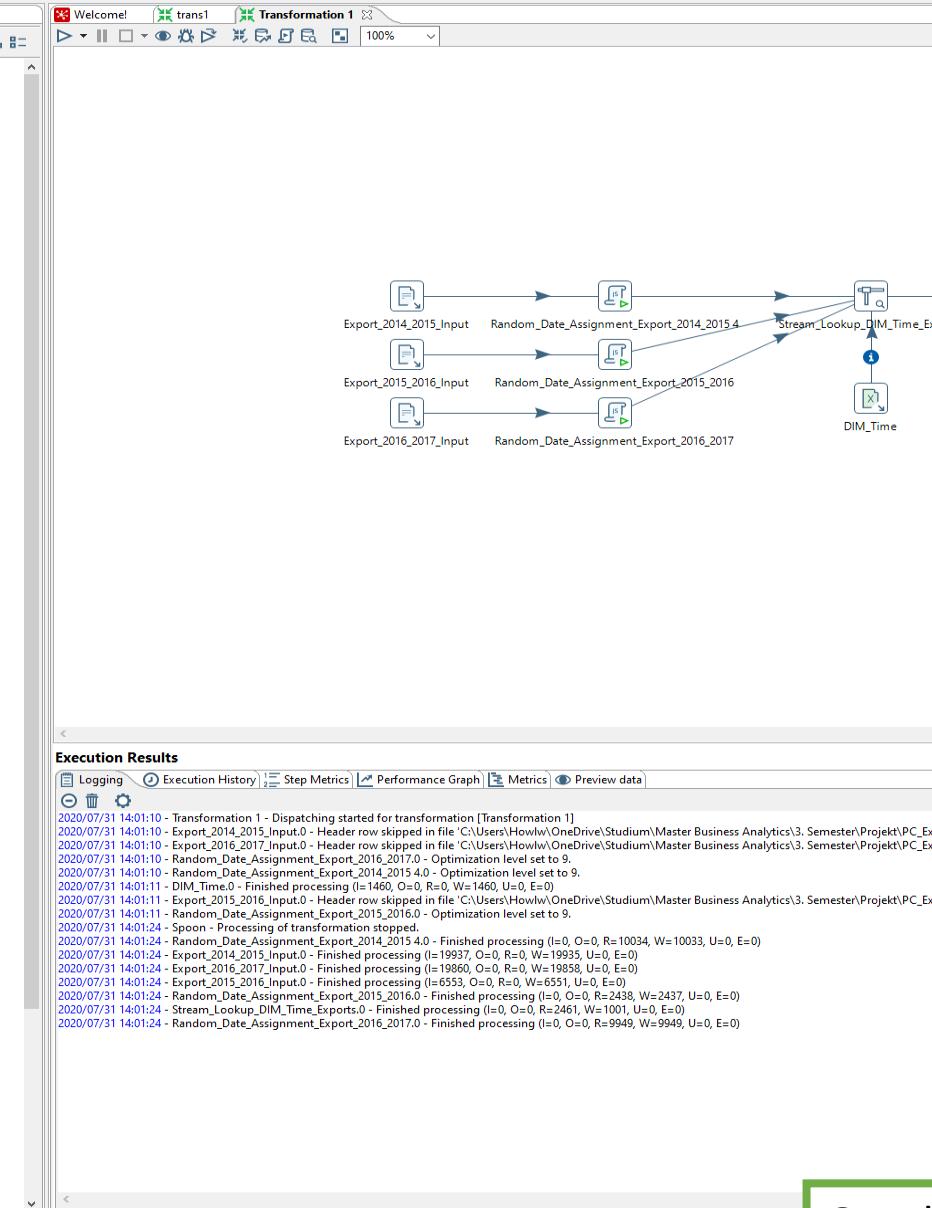
```
graph LR; A[Export_2014_2015_Input] --> B[Random_Date_Assignment_Export_2014_2015_4]; C[Export_2015_2016_Input] --> D[Random_Date_Assignment_Export_2015_2016_0]; E[Export_2016_2017_Input] --> F[Random_Date_Assignment_Export_2016_2017_0]; G[Stream_Lookup_DIM_Time_Exports] --> H[DIM_Time]
```

Name the Step properly  
Select the DIM\_Region.xlsx file via „Browse...“ and add it  
Click on Fields

**View Design**

Search

- Input**
  - CSV file input
  - Data grid
  - De-serialize from file
  - ESRI shapefile reader
  - Email messages input
  - Fixed file input
  - GZIP CSV input
  - Generate random credit card numbers
  - Generate random value
  - Generate rows
  - Get data from XML
  - Get file names
  - Get files rows count
  - Get repository names
  - Get subfolder names
  - Get system info
  - Get table names
  - Google Analytics
  - HL7 input
  - JSON input
  - LDAP input
  - LDIF input
  - Load file content in memory
  - Microsoft Access input
  - Microsoft Excel input**
  - Mondrian input
  - OLAP input
  - Property input
  - RSS input
  - S3 CSV input
  - SAS input
  - Salesforce input
  - Table input
  - Text file input
  - XBase input
  - XML input stream (StAX)
  - YAML input
- Output**
- Streaming**
- Transform**
- Utility**
- Flow**
- Scripting**
- Pentaho Server
- Lookup
- Join
- Data Warehouse
- Validation
- Statistics
- Bin Data



Get the fields from the header row  
Do not forget to change the datatype of region\_code to integer since this was defined in the operational data

Spoon - Transformation 1 (changed)

File Edit View Action Tools Help

View Design

Search

> Input  
Output  
Streaming  
Transform  
  Add XML  
  Add a checksum  
  Add constants  
  Add sequence  
  Add value fields changing sequence  
  Calculator  
  Closure generator  
  Concat fields  
  Get ID from slave server  
  Number range  
  Replace in string  
  Row denormaliser  
  Row flattener  
  Row normaliser  
  Select values  
  Set field value  
  Set field value to a constant  
  Sort rows  
  Split field to rows  
  Split fields  
  String operations  
  Strings cut  
  Unique rows  
  Unique rows (HashSet)  
  Value mapper  
  XSL transformation

> Utility  
Flow  
Scripting  
Pentaho Server  
  Lookup  
  Joins  
Data Warehouse  
Validation  
Statistics  
Big Data  
Agile  
Cryptography  
Job  
Mapping  
Bulk loading  
Inline  
Experimental  
Deprecated  
History

Welcome! trans1 Transformation 1

100%

Stream lookup

Step name: Stream\_Lookup\_DIM\_Region\_Exports  
Lookup step: DIM\_Region

The key(s) to look up the value(s):

#	Field	LookupField
1	country_code	region_code

Specify the fields to retrieve:

#	Field	New name	Default	Type
1	region_country			String
2	region_continent			String
3	region_subcontinent			String

Preserve memory (costs CPU)   
Key and value are exactly one integer field   
Use sorted list (i.s.o. hashtable)

OK Cancel Get Fields Get lookup fields

Execution Results

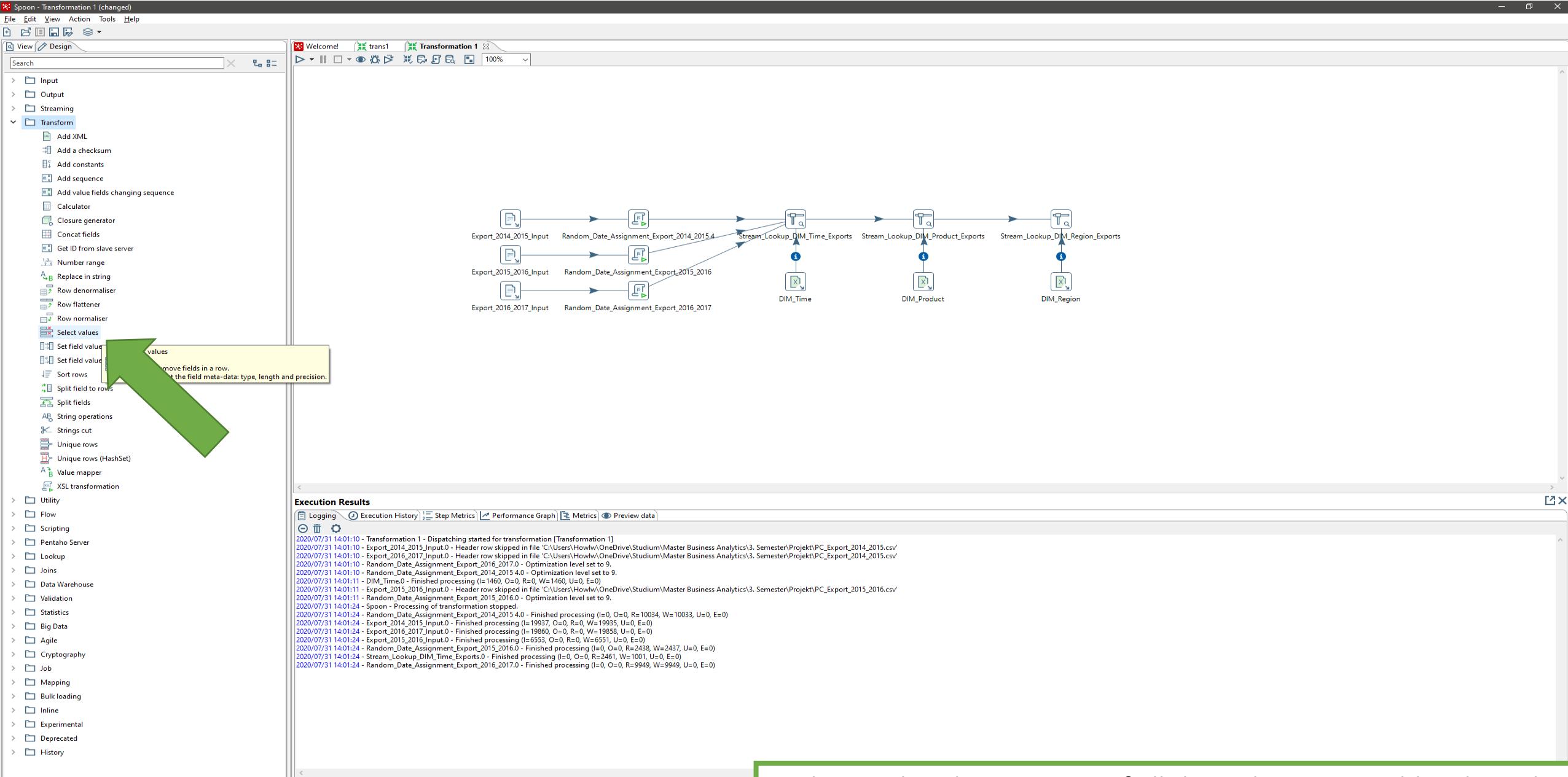
Logging Execution History Step Metrics Performance Graph Metrics Preview data

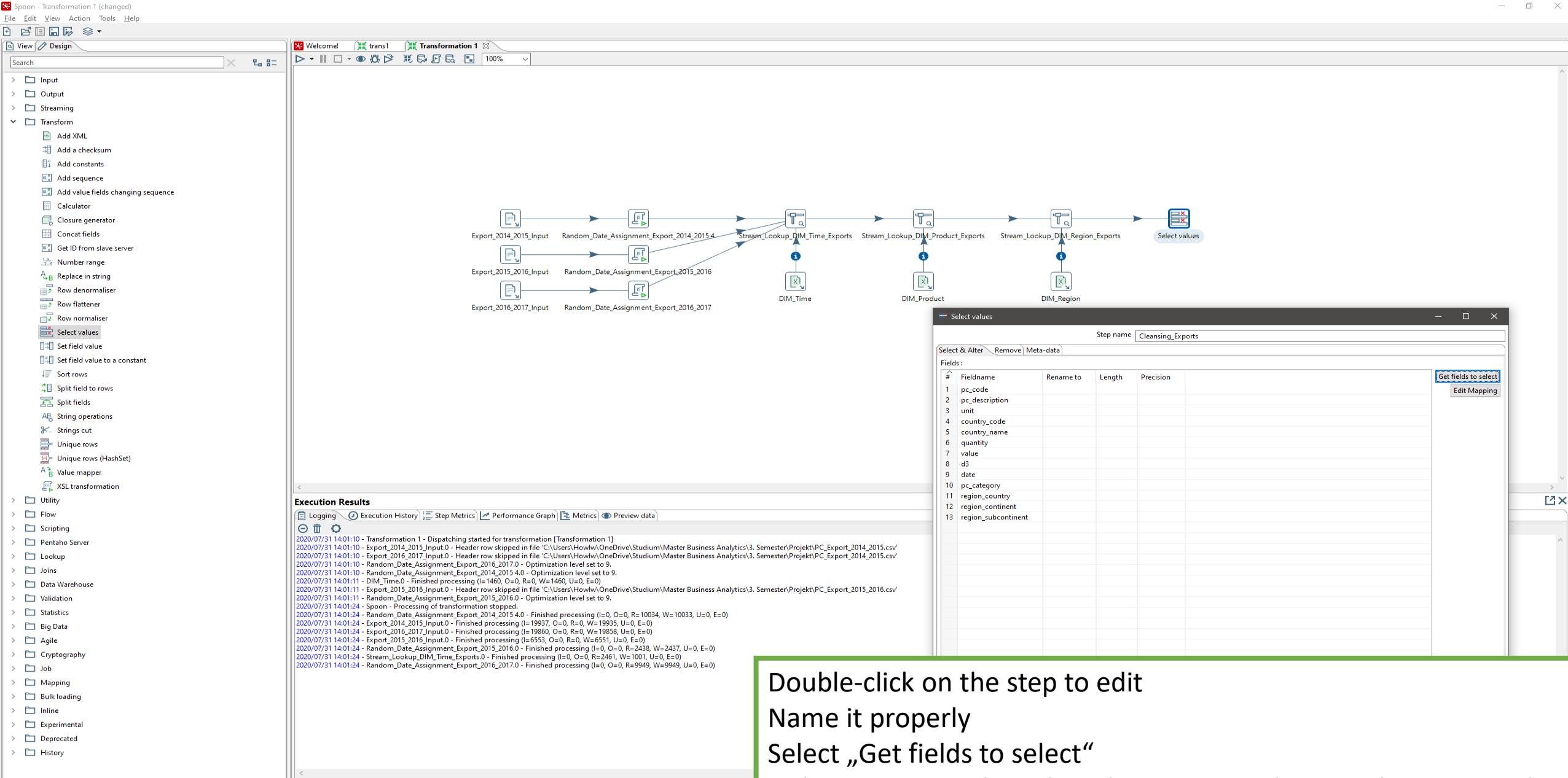
```

2020/07/31 14:01:10 - Transformation 1 - Dispatching started for transformation [Transformation 1]
2020/07/31 14:01:10 - Export_2014_2015_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_1\PC_1\Export_2014_2015\Input.0'
2020/07/31 14:01:10 - Export_2015_2016_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_1\PC_1\Export_2015_2016\Input.0'
2020/07/31 14:01:10 - Random_Date_Assignment_Export_2014_2015.0 - Optimization level set to 9.
2020/07/31 14:01:10 - Random_Date_Assignment_Export_2014_2015.0 - Optimization level set to 9.
2020/07/31 14:01:11 - DIM_Time.0 - Finished processing (I=1460, O=0, R=0, W=1460, U=0, E=0)
2020/07/31 14:01:11 - Random_Date_Assignment_Export_2015_2016.0 - Finished processing (I=19860, O=0, R=0, W=19858, U=0, E=0)
2020/07/31 14:01:12 - Export_2016_2017_Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC_Export_2015_2016.csv'
2020/07/31 14:01:11 - Random_Date_Assignment_Export_2015_2016.0 - Optimization level set to 9.
2020/07/31 14:01:12 - Random_Date_Assignment_Export_2015_2016.0 - Finished processing (I=0, O=0, R=0, W=10034, U=0, E=0)
2020/07/31 14:01:12 - Stream_Lookup_DIM_Time_Exports.0 - Finished processing (I=0, O=0, R=2438, W=2437, U=0, E=0)
2020/07/31 14:01:12 - Stream_Lookup_DIM_Time_Exports.0 - Finished processing (I=0, O=0, R=2461, W=1001, U=0, E=0)
2020/07/31 14:01:12 - Stream_Lookup_DIM_Time_Exports.0 - Finished processing (I=0, O=0, R=9949, W=9949, U=0, E=0)
2020/07/31 14:01:12 - Random_Date_Assignment_Export_2016_2017.0 - Finished processing (I=0, O=0, R=9949, W=9949, U=0, E=0)
2020/07/31 14:01:12 - Spoon - Processing of transformation stopped.

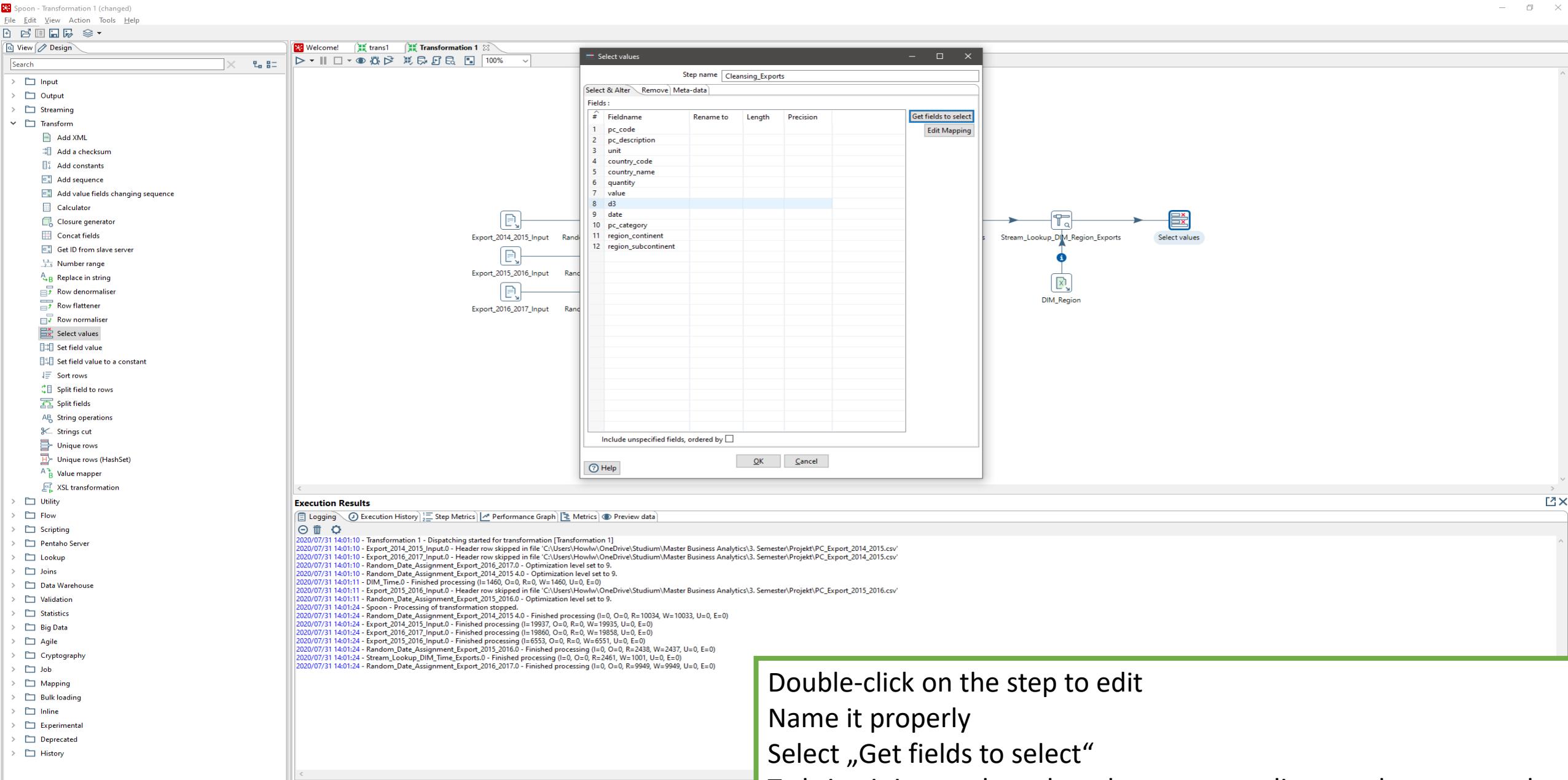
```

Edit the Lookup Step as stated on the picture above  
Press OK

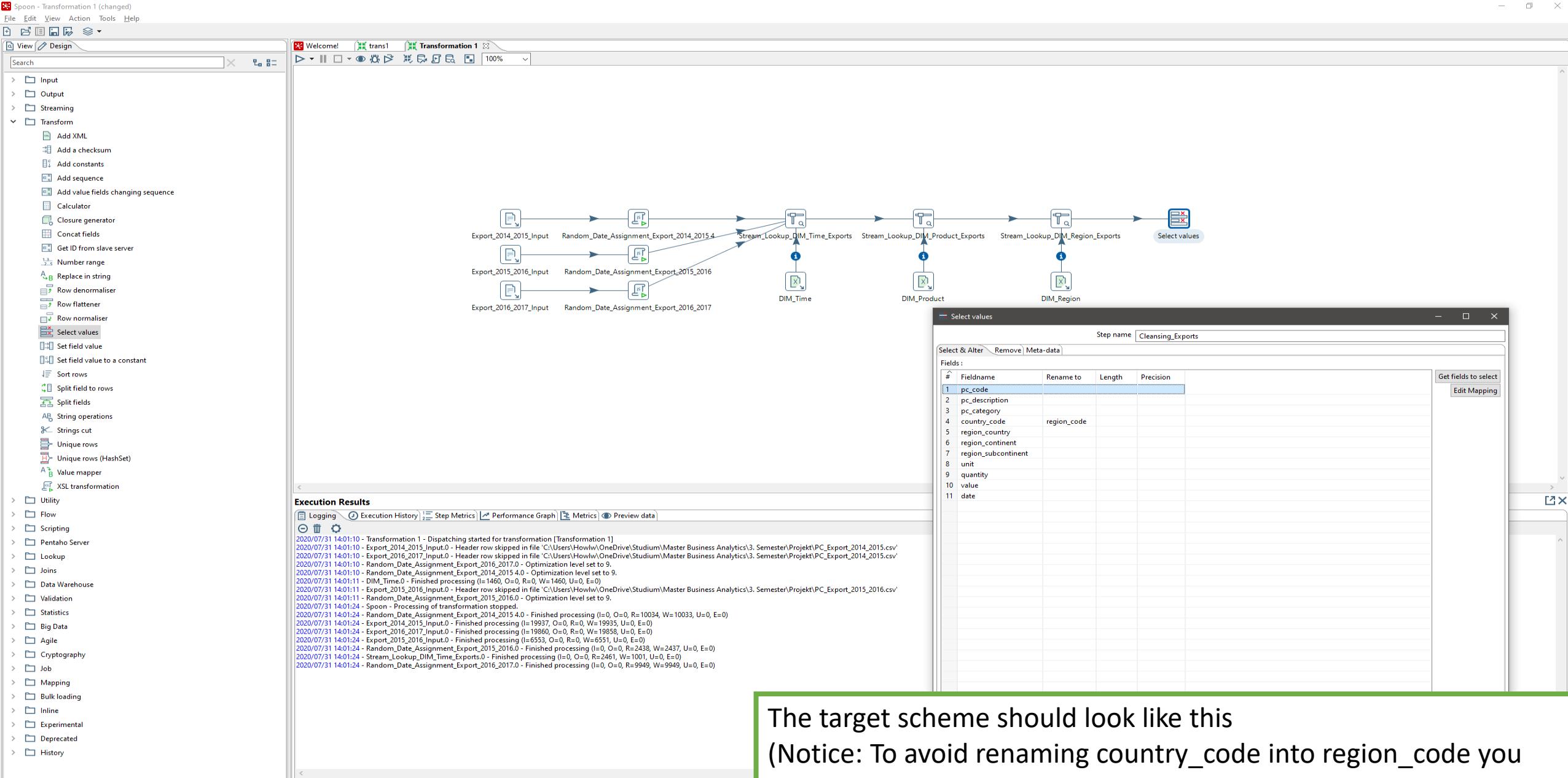




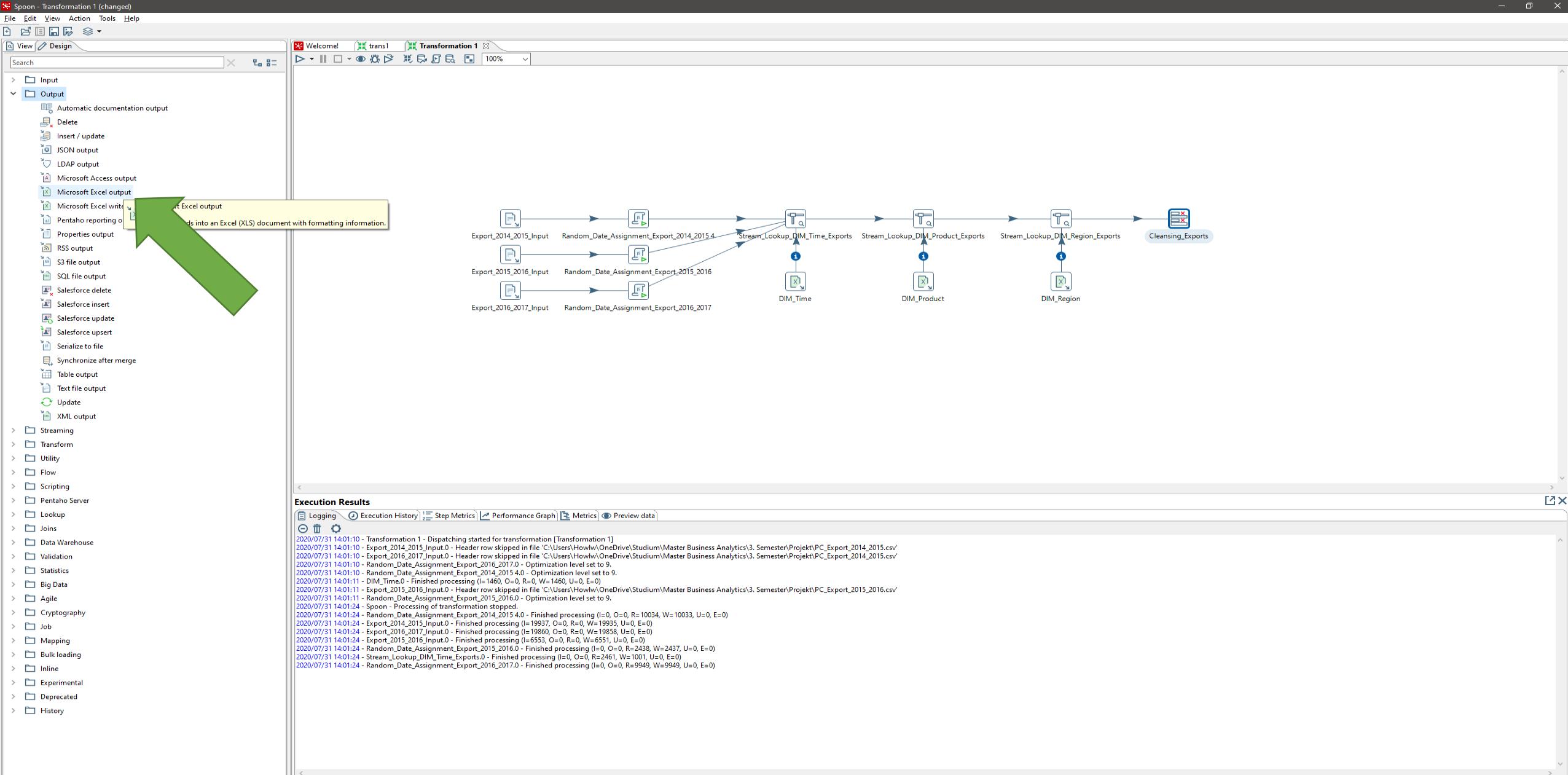
Double-click on the step to edit  
Name it properly  
Select „Get fields to select“  
To bring it into order select the corresponding number next to the filename and press STRG+Arrow Up to move rows upwards and STRG+Arrow Down to move rows downwards



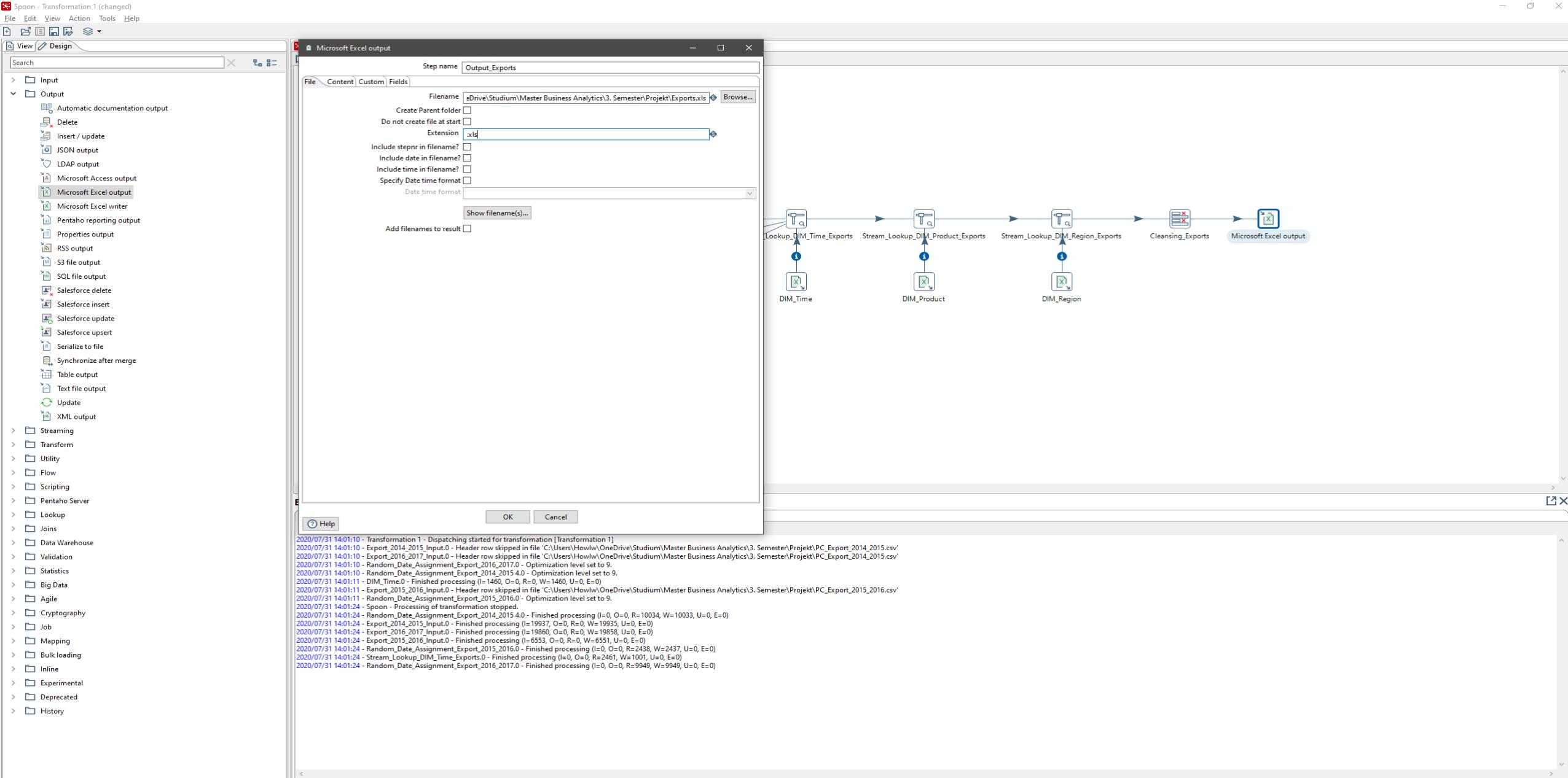
Double-click on the step to edit  
 Name it properly  
 Select „Get fields to select“  
 To bring it into order select the corresponding number next to the filename and press STRG+Arrow Up to move rows upwards and STRG+Arrow Down to move rows downwards



The target schema should look like this  
(Notice: To avoid renaming country\_code into region\_code you can retrieve the region\_code field from DIM\_region and finally delete country\_code in this step)  
Press OK



Finally for the Output Step drag Microsoft Excel output onto the main screen and double-click it



Name it properly  
Click „Browse...“ to choose the destination for saving the file

Spoon - Tuttrans

File Edit View Action Tools Help

View Design

Search

- > Input
- > Output
- > Streaming
- > Transform
  - Add XML
  - Add a checksum
  - Add constants
  - Add sequence
  - Add value fields changing sequence
  - Calculator
  - Closure generator
  - Concat fields
  - Get ID from slave server
  - Number range
  - Replace in string
  - Row denormaliser
  - Row flattener
  - Row normaliser
  - Select values
  - Set field value
  - Set field value to a constant
  - Sort rows
  - Split field to rows
  - Split fields
  - String operations
  - Strings cut
  - Unique rows
  - Unique rows (HashSet)
  - Value mapper
  - XSL transformation
- > Utility
- > Flow
- > Scripting
- > Pentaho Server
- > Lookup
- > Joins
- > Data Warehouse
- > Validation
- > Statistics
- > Big Data
- > Agile
- > Cryptography
- > Job
- > Mapping
- > Bulk loading
- > Inline
- > Experimental
- > Deprecated
- > History

Welcome! trans1 Tuttrans

100%

Execution Results

2020/07/31 21:36:16 - Spoon - Using legacy execution engine  
 2020/07/31 21:36:16 - Spoon - Transformation opened.  
 2020/07/31 21:36:16 - Spoon - Launching transformation [Tuttrans]...  
 2020/07/31 21:36:16 - Spoon - Started the transformation execution.  
 2020/07/31 21:36:16 - Tuttrans - Dimension started for transformation [Tuttrans]  
 2020/07/31 21:36:16 - Export\_2015\_2016.Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_Export\_2015\_2016.csv'  
 2020/07/31 21:36:16 - DIM\_Product.0 - Finished processing (I=168, O=0, R=0, W=168, U=0, E=0)  
 2020/07/31 21:36:16 - Random\_Date\_Assignment\_Export\_2015\_2016.0 - Optimization level set to 9.  
 2020/07/31 21:36:16 - Export\_2016\_2017.Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_Export\_2014\_2015.csv'  
 2020/07/31 21:36:16 - Random\_Data\_Assignment\_Export\_2016\_2017.0 - Optimization level set to 9.  
 2020/07/31 21:36:16 - DIM\_Region.0 - Finished processing (I=233, O=0, R=0, W=233, U=0, E=0)  
 2020/07/31 21:36:16 - Export\_2014\_2015.Input.0 - Header row skipped in file 'C:\Users\Howlw\OneDrive\Studium\Master Business Analytics\3. Semester\Projekt\PC\_Export\_2014\_2015.csv'  
 2020/07/31 21:36:16 - Random\_Date\_Assignment\_Export\_2014\_2015.0 - Optimization level set to 9.  
 2020/07/31 21:36:16 - DIM\_Time.0 - Finished processing (I=146, O=0, R=0, W=146, U=0, E=0)  
 2020/07/31 21:36:17 - Export\_2014\_2015.Input.0 - Finished processing (I=20831, O=0, R=0, W=20830, U=0, E=0)  
 2020/07/31 21:36:17 - Export\_2016\_2017.Input.0 - Finished processing (I=20831, O=0, R=0, W=20830, U=0, E=0)  
 2020/07/31 21:36:17 - Export\_2015\_2016.Input.0 - Finished processing (I=20772, O=0, R=0, W=20771, U=0, E=0)  
 2020/07/31 21:36:17 - Random\_Data\_Assignment\_Export\_2014\_2015.0 - Finished processing (I=0, O=0, R=20830, W=20830, U=0, E=0)  
 2020/07/31 21:36:17 - Random\_Date\_Assignment\_Export\_2015\_2016.0 - Finished processing (I=0, O=0, R=20771, W=20771, U=0, E=0)  
 2020/07/31 21:36:32 - Stream\_Lookup\_DIM\_Time\_Exports.0 - linen 50000  
 2020/07/31 21:37:46 - Stream\_Lookup\_DIM\_Product\_Exports.0 - linen 50000  
 2020/07/31 21:38:16 - Stream\_Lookup\_DIM\_Exports.0 - Finished processing (I=0, O=0, R=63891, W=62431, U=0, E=0)

Microsoft Excel output

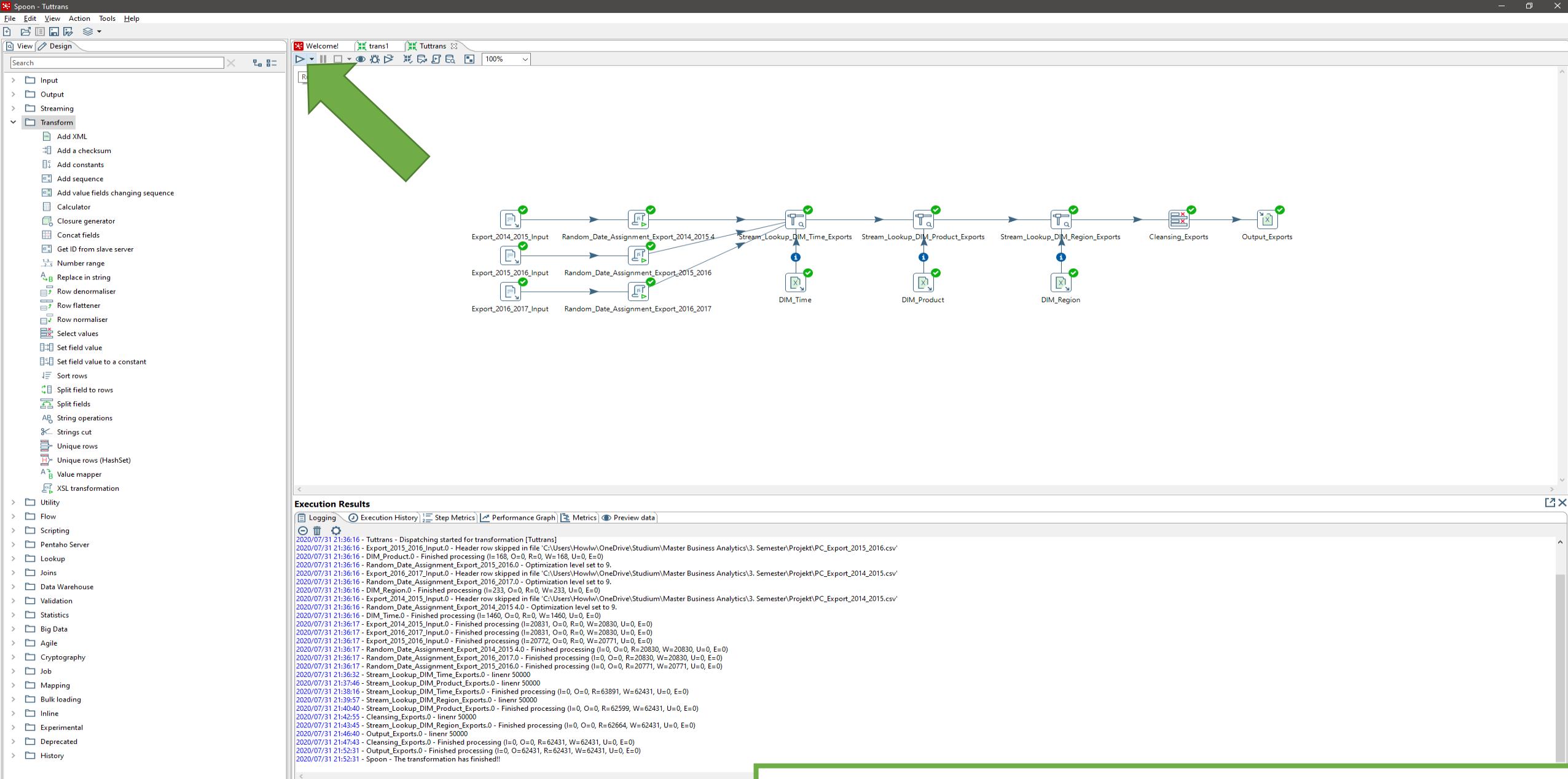
Step name: Output\_Exports

File Content Custom Fields

#	Name	Type	Format
1	pc_code	String	
2	pc_description	String	
3	pc_category	String	
4	region_code	Integer	
5	region_country	String	
6	region_continent	String	
7	region_subcontinent	String	
8	unit	String	
9	quantity	Integer	
10	value	Number	000000000000.00
11	date	Date	

Get Fields Minimal width OK Cancel Help

Switch to Fields and edit the Fields accordingly  
 Press OK



Finally press „Run“ and then again „Run“ without changing anything  
The transformation should start and finish successfully