


# Transformation of BHL Europe data using Pentaho Data Integration

**BHL**  
Biodiversity Heritage Library  
**EUROPE**

 co-funded by the  
Community programme  
eContentplus



2

## Outline:

- Introduction to Pentaho Data Integration
- Transforming Data from Relational Databases
- Transforming Data from XML Input



**BHL**  
Biodiversity Heritage Library  
**EUROPE**

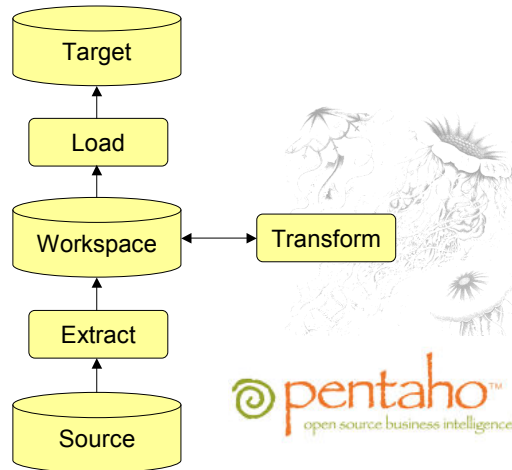
2009, AIT Forschungsgesellschaft mbH

 co-funded by the  
Community programme  
eContentplus



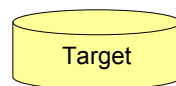
### Introduction to Pentaho Data Integration (1)

- Part of Pentaho Business Intelligence Suite
- ETL Tool
- Open Source



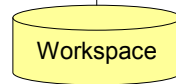
### Introduction to Pentaho Data Integration (2)

Output steps



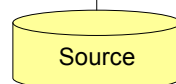
Load

Transform steps



Transform

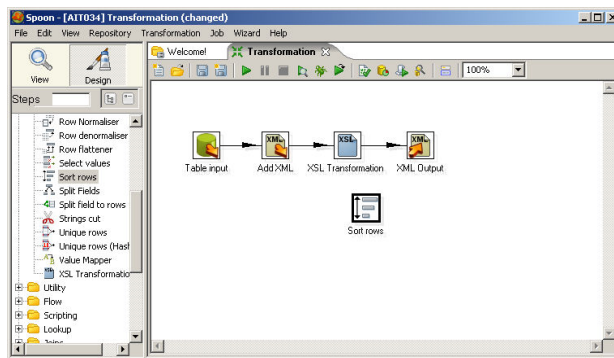
Input steps



Extract



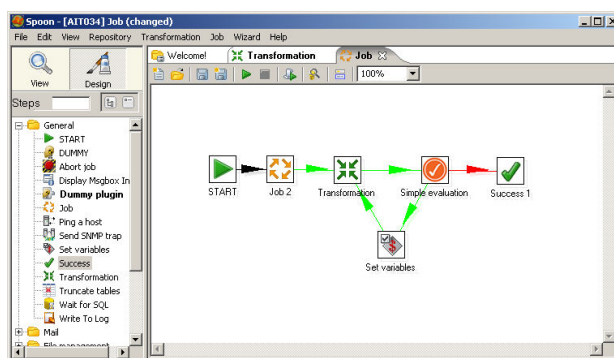
## Introduction to Pentaho Data Integration (3)



Steps are created and connected via Drag and Drop.



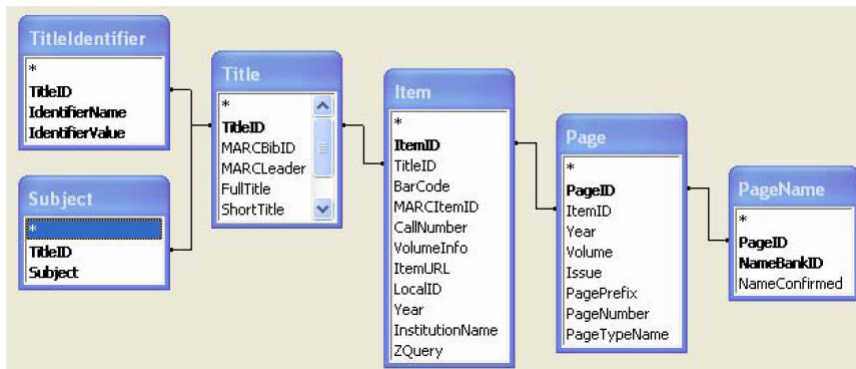
## Introduction to Pentaho Data Integration (4)



- Jobs control the workflow
- Other jobs and transformations can be called
- Conditional evaluations are possible (loops)

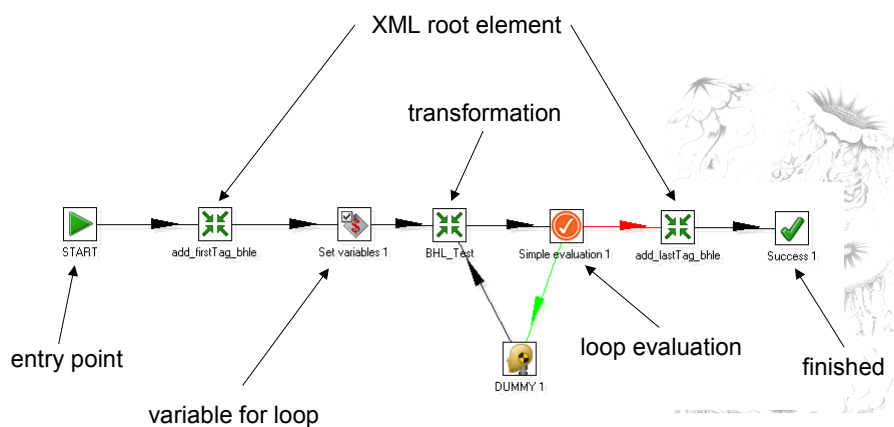


### Transforming Data from Relational Databases

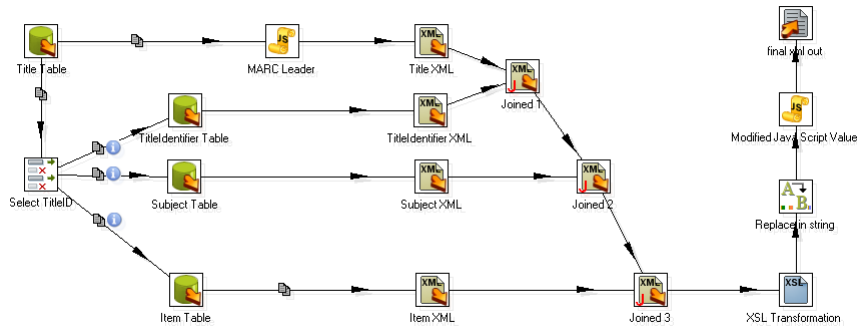


Database schema (Export) of Missouri Botanical Garden (MOBOT)

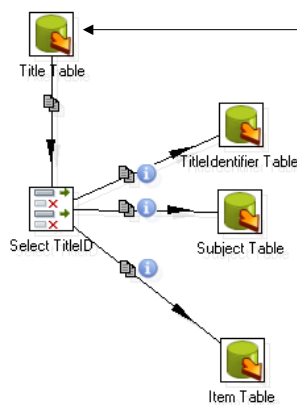
### Missouri Botanical Garden (MOBOT) - Job



## Missouri Botanical Garden (MOBOT) - Transformation



## Getting Data from a Relational Database (1)



**Table input**

Step name: Title Table

Connection: BHL\_Export

SQL: `SELECT * FROM Title  
LIMIT 1 OFFSET ${titleid}`

Line 1 Column 0

Enable lazy conversion:

Replace variables in:

Insert data from step: [dropdown]

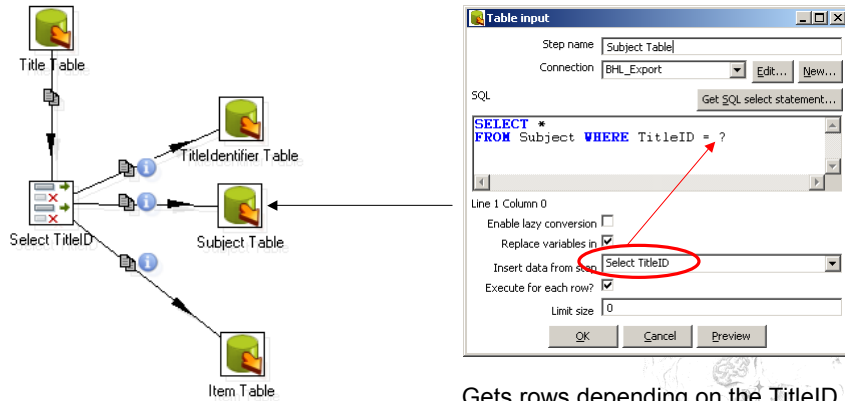
Execute for each row?:

Limit size: 1

OK Cancel Preview

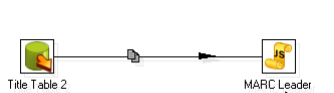
Gets one row depending on loop variable of the job.

## Getting Data from a Relational Database (2)



Gets rows depending on the TitleID of the current title.

## Special Treatment of Fields



JavaScript can handle special fields, e.g. extracting information from the MARCLeader.

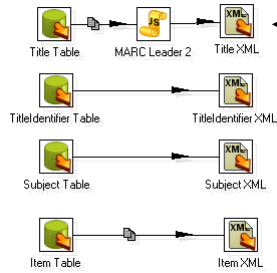
```
Script 1
//Script here
var marcleader= MARCLeader;
var resourceValue;
var manuscript;
var collection;

var pos_six = MARCLeader.charAt(6);
var pos_seven = MARCLeader.charAt(7);

switch (pos_six) {
case "a": resourceValue = "text";
          break;
case "t": resourceValue = "text";
          break;

case "e": resourceValue = "cartographic";
          break;
case "f": resourceValue = "cartographic";
          break;
case "c": resourceValue = "notated music";
          break;
case "d": resourceValue = "notated music";
          break;
}
```

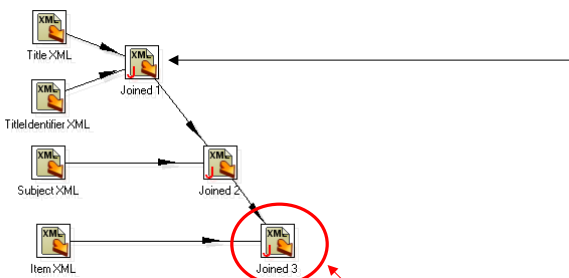
### Creating XML Parts



XML data is created with information from each table.

#	Fieldname	Element name
1	TitleID	
2	MARCBibID	
3	MARCLeader	
4	resourceValue	typeOfResource
5	manuscript	
6	collection	
7	FullTitle	
8	ShortTitle	
9	PublicationDetails	
10	CallNumber	
11	StartYear	
12	EndYear	
13	LanguageCode	
14	TL2Author	
15	TitleURL	

### Joining XML Parts



These XMLs are merged into a single XML field.

Step name: Joined 1

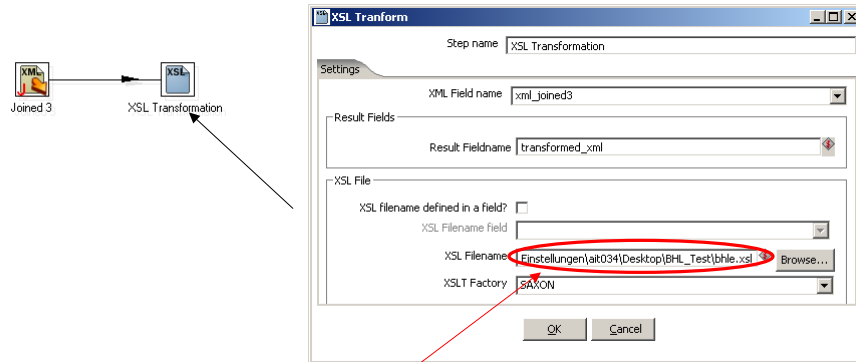
Target stream properties:  
 Target XML step: Title XML  
 Target XML field: title\_xml

Source stream properties:  
 Source XML step: Title Identifier XML  
 Source XML field: titleidentifier\_xml

Join condition properties:  
 XPath Statement: //Title  
 Complex join?:   
 Join comparison:

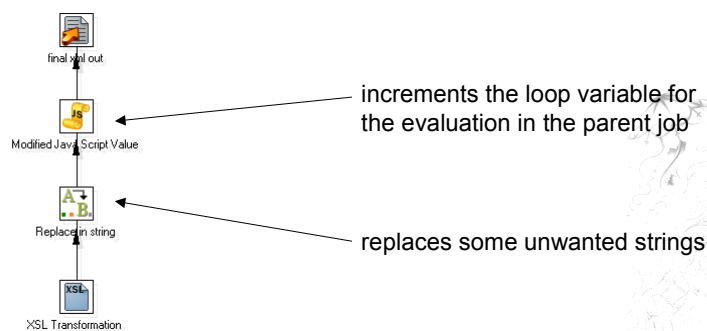
Result stream properties:  
 Result XML field: xml\_joined1  
 Encoding: UTF-8  
 Omit XML header:

## Applying a XSL Transformation



A XSL file can be specified to perform a transformation on the incoming XML field in order to adapt it to a specific XML schema, e.g. MODS.

## Finalizing the Transformation (1)



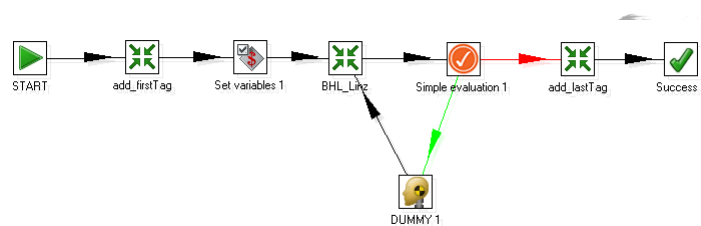


## Finalizing the Transformation (2)



The text file output step is used to store the transformed information in a XML file.

## Landesmuseum Oberösterreich (LANDOE) - Job

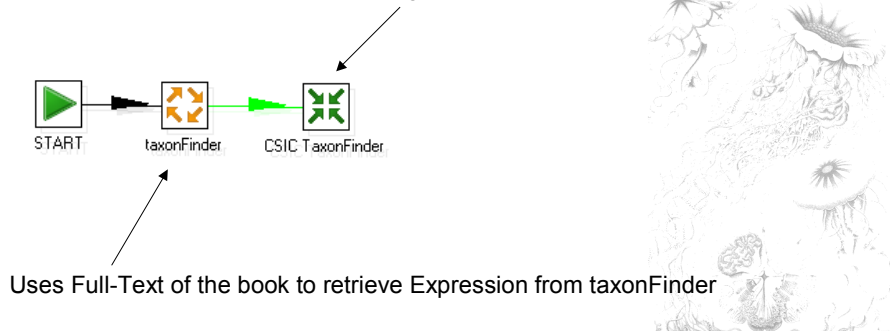


This job is equal to the one used for the MOBOT data.

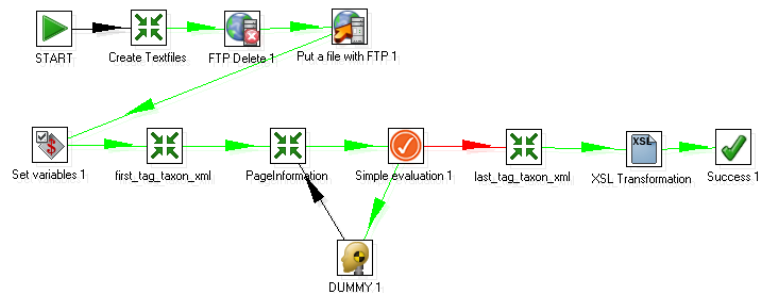


## Consejo Superior de Investigaciones Cientificas (CSIC)

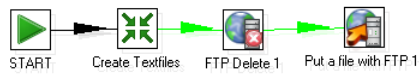
Creates CSIC record and merges it with the Information from taxonFinder



## Job: taxonFinder

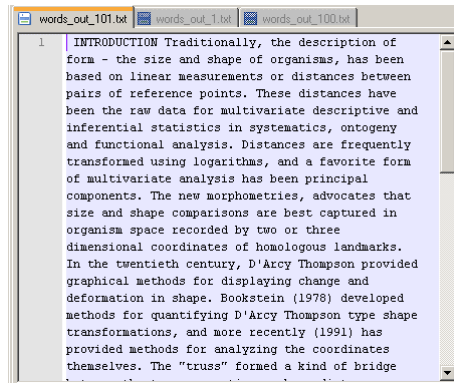


## Create Files and Transfer via FTP

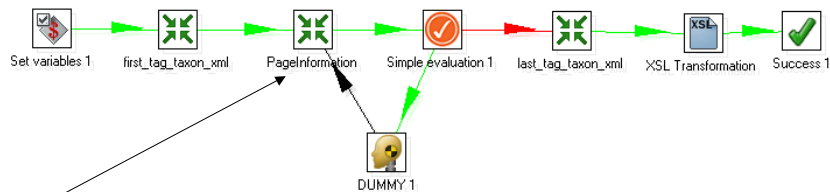


Text files are created for each page, containing all its words space separated.

These files are then transferred via FTP in order to pass URL to the taxonFinder service.

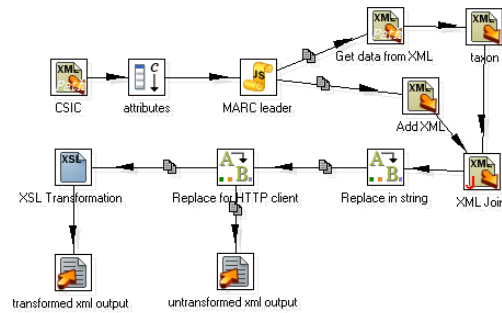


## Create XML containing taxonFinder words (1)



In the loop taxonFinder is called for every single page and the words occurring are saved with the appropriate link to the page in the Internet Archive.

## Merging the record with taxonFinder Information (1)



## Getting Data from a XML Input File



Fields can be specified using XPath statements. These fields can then be transformed as before.

Get XML Data

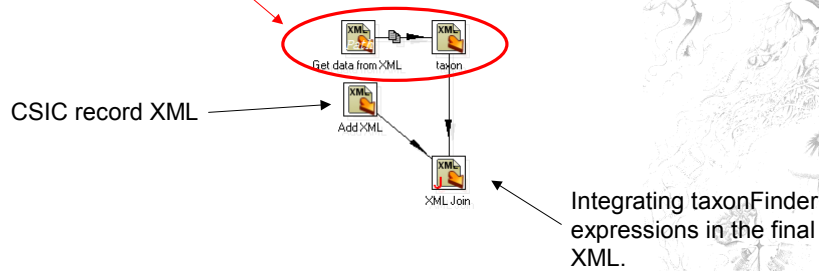
Step name: CSIC

#	Name	XPath	Element
1	MARLeader	record/leader	Node
2	controlfield	record/controlfield[@tag="001"]	Node
3	controlfield2	record/controlfield[@tag="005"]	Node
4	datafield_019_a	record/datafield[@tag="019"]/subfield[@code="a"]	Node
5	datafield_020_a	record/datafield[@tag="020"]/subfield[@code="a"]	Node
6	datafield_040_a	record/datafield[@tag="040"]/subfield[@code="a"]	Node
7	datafield_040_b	record/datafield[@tag="040"]/subfield[@code="b"]	Node
8	datafield_040_c	record/datafield[@tag="040"]/subfield[@code="c"]	Node
9	datafield_080_a-1	record/datafield[@tag="080"][1]/subfield[@code="a"]	Node
10	datafield_080_a-2	record/datafield[@tag="080"][2]/subfield[@code="a"]	Node
11	datafield_080_a-3	record/datafield[@tag="080"][3]/subfield[@code="a"]	Node
12	datafield_100_a	record/datafield[@tag="100"]/subfield[@code="a"]	Node
13	datafield_245_a	record/datafield[@tag="245"]/subfield[@code="a"]	Node
14	datafield_245_c	record/datafield[@tag="245"]/subfield[@code="c"]	Node
15	datafield_260_a	record/datafield[@tag="260"]/subfield[@code="a"]	Node
16	datafield_260_b	record/datafield[@tag="260"]/subfield[@code="b"]	Node
17	datafield_260_c	record/datafield[@tag="260"]/subfield[@code="c"]	Node

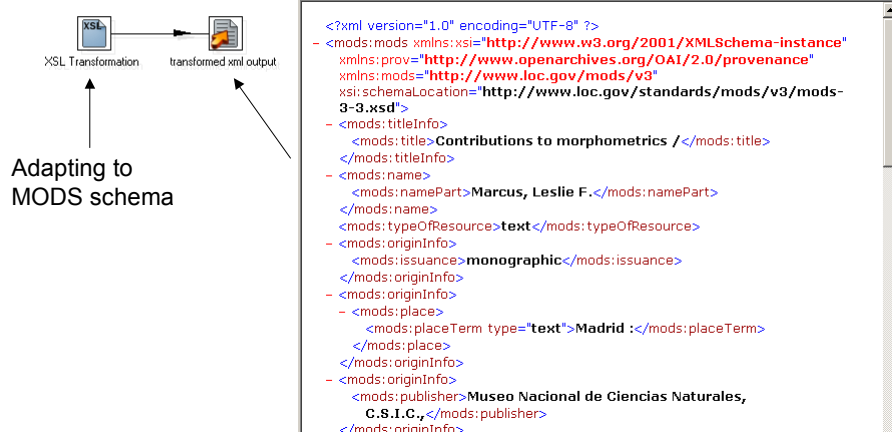
Buttons: Get fields, OK, Preview rows, Cancel

## Merging the record with taxonFinder Information (2)

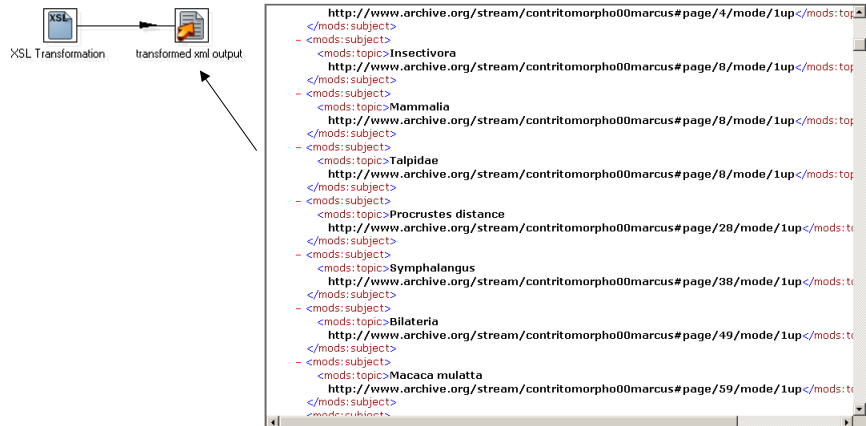
Reading XML file created with taxonFinder



## XML Output in MODS Schema (1)



## XML Output in MODS Schema (2)



## Useful Links/References

- Official Website of the Business Intelligence Suite:  
<http://www.pentaho.org>
- Website of Pentaho Data Integration:  
<http://kettle.pentaho.org>
- A beginners tutorial:  
<http://etl-tools.info/en/pentaho/kettle-spoon.htm>
- Pentaho Data Integration User Guide:  
<http://wiki.pentaho.com/display/EAI/Spoon+User+Guide>

