

# Talend 사용자 매뉴얼

2014년 1월  
솔루션팀

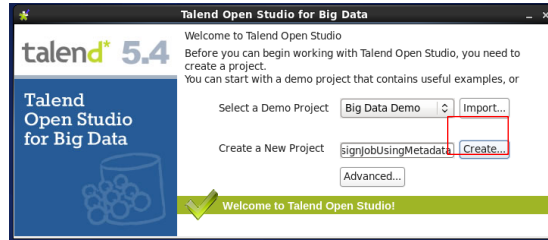
1. 프로젝트 생성
2. Main Windows 구성 설명
3. Job Design 실습
  - 3.1 CSV -> XML 생성
  - 3.2 정렬하기
  - 3.3 맵핑하기
4. Oracle 접속 실습
  - 4.1 Oracle DB 접속하기
  - 4.2 AmazonRDS 접속하기
  - 4.3 Amazon RDS - Multi thread execution 사용하기
5. 컴포넌트 사용 실습
  - 5.1 데이터 집계와 정렬하기
  - 5.2 데이터 타입 변환
  - 5.3 한 필드가 콤마로 구분된 파일 추출하기
  - 5.4 Error messages 출력하기
  - 5.5 소셜 네트워크에서 데이터 추출하기
  - 5.6 이름으로 필터링하기
  - 5.7 두개의 파일 조인하여 엑셀과 리젝트 파일로 출력하기
  - 5.8 데이터 표준화하기
  - 5.9 데이터 매핑과 Reject하기
  - 5.10 Inner join rejection

# 1. 프로젝트 생성

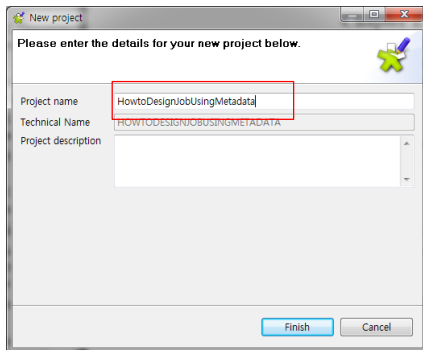
①



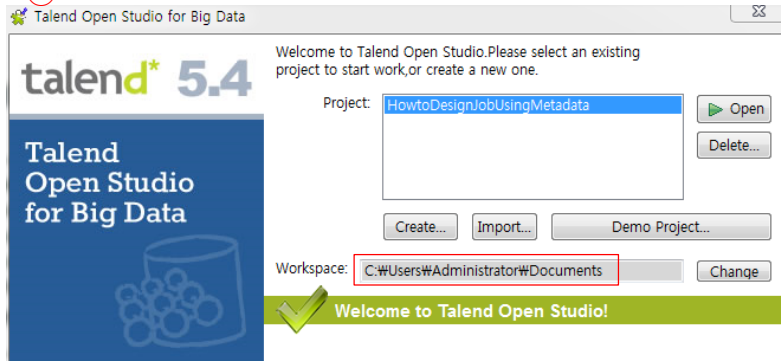
②



③



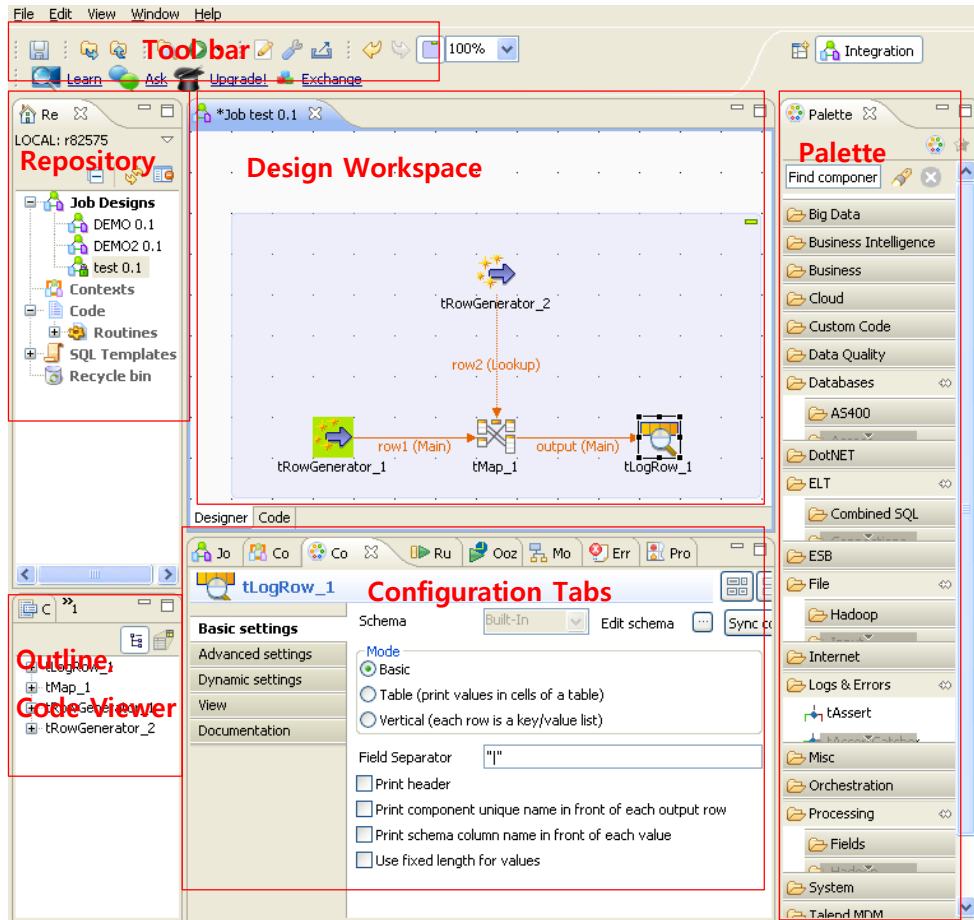
④



1. 탈렌드 BigData 프로그램을 실행하면 로고 화면이 표시된다.
2. 데모 프로젝트를 열거나 새로운 프로젝트를 열수 있는 창이 나타난다. 이 창에서 프로젝트를 생성하기 위하여 Create 버튼을 클릭한다.
3. 새로운 프로젝트의 이름과 프로젝트 설명을 입력한 후 Finish버튼을 클릭한다. Workspace 폴더 안에 프로젝트 이름으로 폴더가 생성된다.
4. 새로 만든 프로젝트가 표시되고, 열거나 삭제할 수 있다. 이 화면에서 추가로 새로운 프로젝트를 만들거나 다른 프로젝트를 Import할 수 있고, 데모 프로젝트를 불러올 수 있다. Workspace를 변경하고자 하는 경우에는 Change버튼을 클릭하여 다른 폴더를 지정해 주면 된다. 새로 생성한 프로젝트를 선택한 후 Open을 클릭하면 Job을 Design할 수 있는 화면이 나타난다.

## 2. Main Windows 구성 설명

### 2.1 Main Windows 구성













#### Main Windows 구성

- **menu bar**
  - 기본적인 이클립스의 기능과 함께 다양한 Talend Studio의 메뉴를 제공한다.
- **toolbar**
  - Save, export items, Run job과 같은 일반적으로 사용하는 작업을 빠르게 액세스할 수 있는 아이콘을 제공한다.
- **Repository tree view**
  - 잡을 디자인할 때 사용할 수 있는 모든 기술 항목, JobDesigns, Contexts, 재사용 가능한 루틴 등을 모아 보여준다.
- **Design workspace**
  - Designer tab에서 수행할 잡을 디자인할 수 있다. Code tab에서는 java code를 보여준다.
- **Palette**
  - Design workspace에 가져올 수 있는 컴포넌트들을 보여준다.
- **Configuration views**
  - Design workspace에서 디자인한 컴포넌트의 여러가지 값을 설정할 수 있다.
- **Outline, Code Viewer**
  - Outline tab은 열려져 있는 Job에서 사용되고 있는 컴포넌트들을 리스트한다.
  - Code Viewer tab은 Design workspace에서 선택한 컴포넌트에 해당하는 코드를 보여준다.

## 2. Main Windows 구성 설명

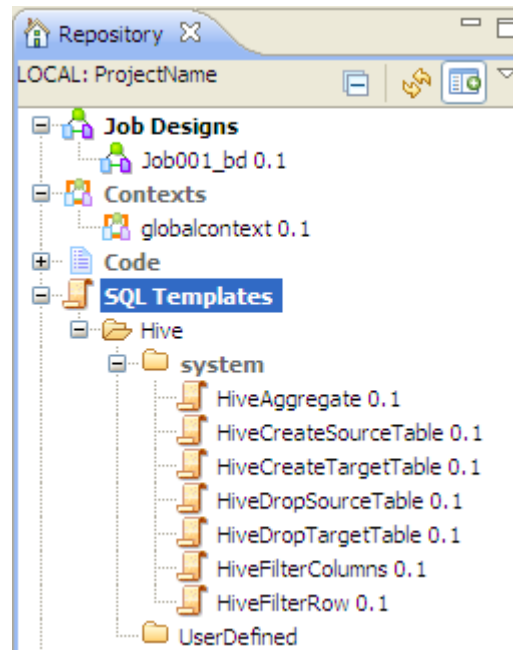
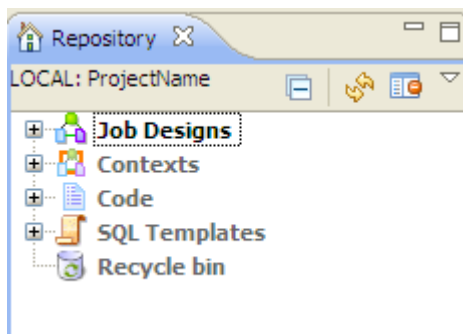
### 2.2 Toolbar 설명

메뉴명	아이콘	설명
Save		디자인한 잡을 저장한다.
Save as		새로운 잡으로 저장한다.
Export items		레파지토리 아이템을 아카이브 파일로 export 한다. 새로운 버전의 Talend Studio나 다른 workstation에서 import하여 사용할 경우 소스파일이 포함되는지 확인해야 한다.
Import items		아카이브 파일로부터 import한다.
Find a specific job		리포지토리 트리 뷰에 나열된 모든 작업을 열 수 있는 대화 상자를 표시한다.
Run job		현재 디자인한 내용을 실행한다.
Create		생성 마법사를 시작한다. 이 메뉴를 통해, 잡 디자인, 컨텍스트 및 루틴을 포함한 모든 레파지토리 항목을 만들 수 있다.
Project settings		[프로젝트 설정] 대화 상자를 실행하여 현재 프로젝트에 대한 설명을 추가하고 팔레트의 표시를 사용자 정의 할 수 있다.
Detect and update all jobs		작업에 사용할 수 있는 모든 업데이트를 검색한다.
Export Talend projects		[Export Talend projects] 마법사를 시작한다.

## 2. Main Windows 구성 설명

### 2.3 Repository 설명

노드명	설명
Job Designs	잡 디자인 폴더는 현재 프로젝트에서 디자인 되어 있는 잡을 트리형식으로 보여준다.
Contexts	파일경로나 DB접속정보 같은 여러 잡에서 재사용할 변수들을 저장한다.
Code	이 프로젝트와 다른 프로젝트에서 사용할 수 있는 코드를 그룹화한 라이브러리이다.
SQL Templates	표준화한 사용자 정의 SQL templates를 생성할 수 있는 폴더이다.
Recycle bin	Repository에서 삭제한 내용을 보관한다. Recycle bin icon을 오른쪽 클릭해서 recycle bin 을 비우기 전까지는 파일 시스템에 삭제된 내용이 존재하고, Recycle bin 에서 바로 삭제된 내용을 복구하거나 완전히 삭제가 가능하다.



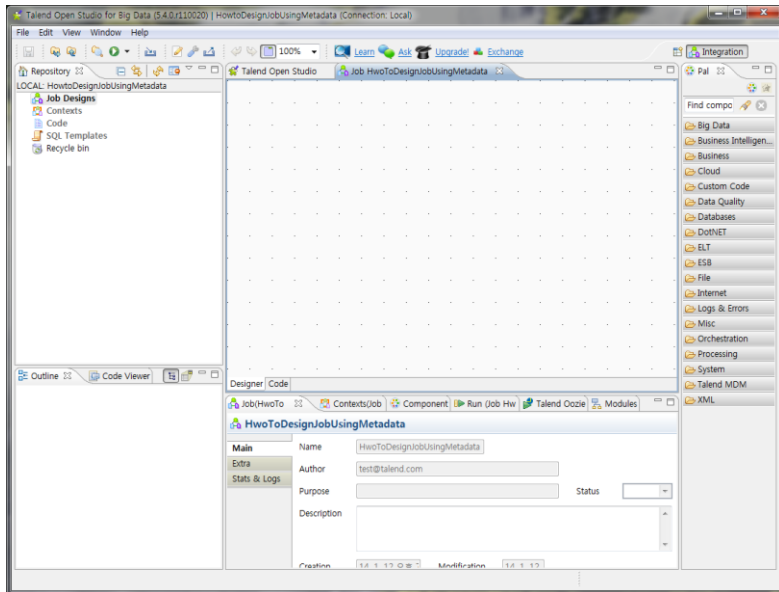
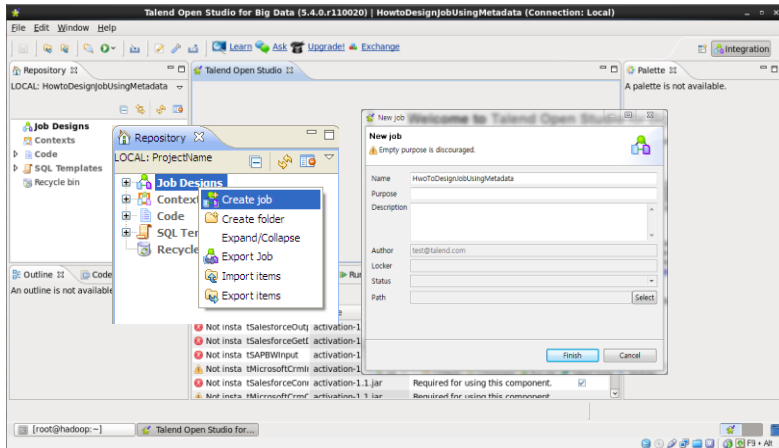
## 2. Main Windows 구성 설명

### 2.4 Configuration 설명

View	설명
Component	Component 뷰는 각 구성 요소에 대한 특정 매개 변수의 자세한 사항들이다. 작업을 구축하기 위해 각각의 컴포넌트가 필요한 필드를 작성해야 한다.
Run Job	Run Job 뷰는 현재 작업의 실행을 보여주고, 실행 후에는 로그 결과를 보여준다.
Oozie scheduler	현재 작업을 실행하거나 원격 HDFS 서버에 주기적으로 실행되도록 예약할 수 있다.
Error Log	Error Log 뷰는 주로 작업 실행 오류에 사용되고 경고나 작업 실행 중에 발생하는 오류의 history를 보여준다. 로그 탭은 또한 Java 구성 요소 운영 진행에 대한 유익한 기능을 가지고 있다. 예를 들어 오류 로그 탭은 기본적으로 숨겨져 있고 Window > Show views 창에서 General 폴더를 확장 후 Error Log를 선택하면 탭이 보여진다.
Modules	작업을 원활하게 수행하기 위하여 컴포넌트에서 필요하거나 참조되는 모듈을 보여준다. 필요한 외부 모듈을 다운받아 설치할 수 있다.
Job view	Job view는 design workspace에 열려 있는 작업과 관련된 다양한 정보를 표시한다. 여기에는 다음과 같은 Tab이 있다.
	<b>Main tab</b> 이 탭은 design workspace에 열려있는 작업에 대한 기본 정보 예를들면 이름, 저자, 버전 번호등을 표시하고, 읽기 전용이고, 이 탭을 편집하기 위해서는 Job을 닫은 후 저장소 트리 뷰에 레이블을 마우스 오른쪽 버튼으로 클릭하고 드롭다운 목록에서 속성 편집을 클릭한다.
	<b>Extra tab</b> 이 탭은 멀티 스레드 실행과 암시적 컨텍스트 로딩 기능을 포함하여 추가 매개 변수를 표시한다.
	<b>Stats/Log tab</b> 이 탭은 모든 작업에 대한 통계 및 로그를 활성화 또는 비활성화 할 수 있고 통계, Log, Volumetrics를 사용여부를 간단히 설정할 수 있다. 또는 Catcher 컴포넌트를 사용하지 않고 한 번에 전체 활성 작업에 대해 이 기능을 설정할 수 있는데, 모든 컴포넌트는 설정에 따라 추적하고 로그를 파일이나 데이터베이스 테이블에 로그를 남길 수 있다. 또한 프로젝트 세팅을 다시 읽어오거나 저장할 수 있다.
	<b>Version tab</b> design workspace에 열려져 있는 잡의 다른 버전에 대하여 생성일자와 변경일자를 보여준다.
Problems	Problems 뷰는 설정의 일부가 누락 된 것과 같은 경우 컴포넌트에 도킹 아이콘에 링크된 메시지를 표시한다. 아이콘 / 메시지는 Error, Warning and Infos 세 가지 유형이 있다.
Job Hierarchy	선택한 잡의 하위 잡을 보여준다. 이 뷰를 표시하려면, Repository 트리 뷰에서 상위 작업을 오른쪽 클릭하고 드롭 다운 목록에서 Open Job Hierarchy 을 선택하거나, Window > Show view 창에서 Talend > Job Hierarchy를 선택하면 된다. tRunJob 컴포넌트를 통해 상위 잡 및 하나 이상의 하위 잡을 만들 경우에만 Job Hierarchy를 볼 수 있다.

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환



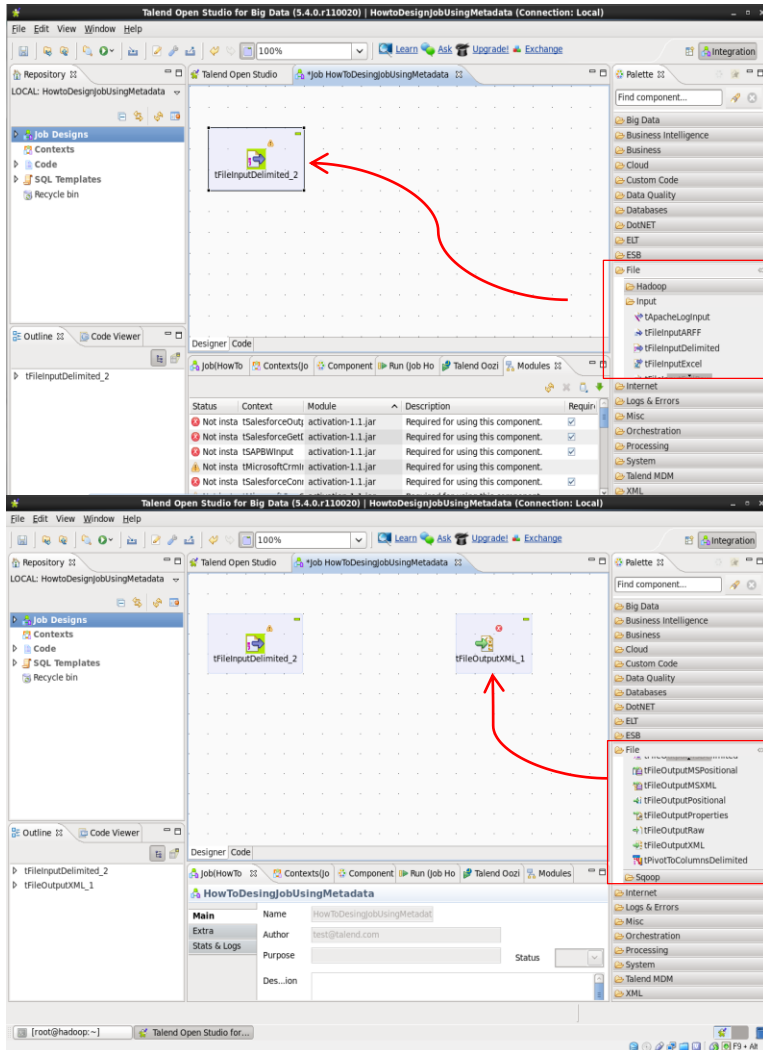
개요 : CSV 파일을 읽어서 XML파일로 변환한다.

1. 새로운 잡을 만들기 위해서 Job Designs에서 마우스 오른쪽 버튼 클릭 후 Create Job을 클릭한다.
2. Name에 새로운 Job 이름, Purpose, Description을 입력한다.  
Workspace 지정 경로 밑에  
HOWTODESIGNJOBUSINGMETADATA#process에  
HwoToDesignJobUsingMetadata\_0.1.item  
HwoToDesignJobUsingMetadata\_0.1.properties  
HwoToDesignJobUsingMetadata\_0.1.screenshot  
3개의 파일이 생성된다
3. Job을 Design할 수 있는 메인 화면이 표시된다



### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환



1. 콤마로 구분되어 있는 Input 파일을 사용하기 위하여 Palette-File-Input-tFileInputDelimited Component 선택 후 Job Designer에 드롭한다.
2. XML 형식의 출력파일을 만들기 위하여 Palette-File-Output-tFileOutputXML Component 선택 후 Job Designer에 드롭한다

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환

The screenshot displays the Talend Open Studio interface. The main workspace shows a job design with a flow from **tFileInputDelimited\_1** to **row1 (Main)** and then to **tFileOutputXML\_1**. The left sidebar shows the 'tFileInputDelimited\_1' component settings. The 'Basic settings' tab is selected, showing the 'File Name/Input Stream' as 'D:/TOS\_BD Test/001/customer.csv'. The 'Header' is set to 0 and 'Footer' to 0. The 'Schema' is set to 'Built-in'. The 'Edit schema' button is highlighted. The 'Schema of tFileInputDelimited\_1' dialog is open, showing a table of columns: id, CustomerName, CustomerAddr, idState, id2, RegTime, RegisterTime, Sum1, Sum2. The 'id' column is highlighted.

Column	Key	Type	N...	Date Patt...	Len...	Prec...	De...	Co...
id	<input checked="" type="checkbox"/>	int	<input type="checkbox"/>		9			
CustomerName	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
CustomerAddr	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
idState	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
id2	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
RegTime	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
RegisterTime	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
Sum1	<input type="checkbox"/>	float	<input type="checkbox"/>		10			
Sum2	<input type="checkbox"/>	float	<input type="checkbox"/>		10			

1. 두 컴포넌트를 연결하기 위하여 tFileInputDelimited에서 마우스 오른쪽 버튼을 클릭한 채로 tFileOutputXML로 드래그한다.
2. 버튼을 놓으면 row1(Main) 연결선이 생긴다.  
tFileInputDelimited에서 마우스 오른쪽 버튼을 클릭 후 row-main을 선택하고 tFileOutputXML을 선택하여도 연결선이 생긴다.
3. Job Designer에서 tFileInputDelimited를 더블클릭하면 아래 부분에Component 판넬이 보인다.
4. Basic settings의 File name/Stream의 [...]버튼에서 input file로 사용할 파일을 선택한다.
5. Input File의 헤더가 6라인이므로 Header에 6을 입력한다.
6. Edit Schema [...] 버튼을 클릭한다
7. tFileInputDelimited 스키마 창의 아래쪽의 + 키를 눌러서 빈 라인을 생성한 후 Input File에서 사용하는 컬럼 정의 내역을 그림과 같이 입력한다

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환

The screenshot shows the Talend Open Studio interface. At the top, a job design is visible with a flow from **tFileInputDelimited\_1** to **tFileOutputXML\_1** labeled **row1 (Main)**. Below the design, the **Designer** tab is active, and the **tFileOutputXML\_1** component is selected. The **Basic settings** panel shows the following configuration:

- File Name:** "D:/TOS\_BD Test/001/customer.xml"
- Incoming record is a document:** ☐
- Row tag:** "row"
- Schema:** Built-In (with an **Edit schema** button highlighted by a red box)
- Sync columns:** ☐

Below the settings, the **Schema of tFileOutputXML\_1** dialog is open, showing two tables for comparison:

Column	Key	Type	zNullab	Date Pat	Leng	Preci	Def	Com
id	<input checked="" type="checkbox"/>	int	<input type="checkbox"/>		9			
CustName	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
CustAddr	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
idState	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
id2	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
RegTime	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
RegisterTim	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
Sum1	<input type="checkbox"/>	float	<input type="checkbox"/>		10			
Sum2	<input type="checkbox"/>	float	<input type="checkbox"/>		10			

Column	Key	Type	zNullab	Date Pat	Leng	Preci	Def	Com
id	<input checked="" type="checkbox"/>	int	<input type="checkbox"/>		9			
CustName	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
CustAddr	<input type="checkbox"/>	String	<input type="checkbox"/>		255			
idState	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
id2	<input type="checkbox"/>	int	<input type="checkbox"/>		2			
RegTime	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
RegisterTim	<input type="checkbox"/>	String	<input type="checkbox"/>		30			
Sum1	<input type="checkbox"/>	float	<input type="checkbox"/>		10			
Sum2	<input type="checkbox"/>	float	<input type="checkbox"/>		10			

1. tFileOutputXML을 더블클릭한 후 아래의 Component 패널에서 생성할 XML 파일의 경로와 이름을 입력한다
2. Edit schema의 [...] 을 클릭하여 input / output column 을 확인한다

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환

The screenshot shows the Talend Open Studio interface. The top part displays a job design with a flow from **tFileInputDelimited\_1** to **tFileOutputXML\_1**. The flow is labeled with performance metrics: "5000 rows in 0.54s", "9345.79 rows/s", and "row1 (Main)".

The bottom part shows the configuration for the **tFileOutputXML\_1** component. The **Advanced settings** tab is selected. The **Root tags** section shows a list of tags with "customer" selected. The **Output format** section shows a table with columns: Column, As attribute, Use sche..., and Label.

Column	As attribute	Use sche...	Label
id	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
CustomerName	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
CustomerAddr	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
idState	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
id2	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
RegTime	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
RegisterTime	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"
Sum1	<input type="checkbox"/>	<input checked="" type="checkbox"/>	"label"

1. Advanced settings에서 1000 Rows마다 파일을 분리하기 위하여 "Split output in several files"를 체크한다.
2. Root tags를 입력하기 위하여 [+] 추가 후 "customer"를 입력한다.  
"를 입력하지 않으면 실행시 오류가 발생하므로 주의해서 입력한다.
3. Output format에서 id 컬럼의 As attribute를 체크한다.

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환

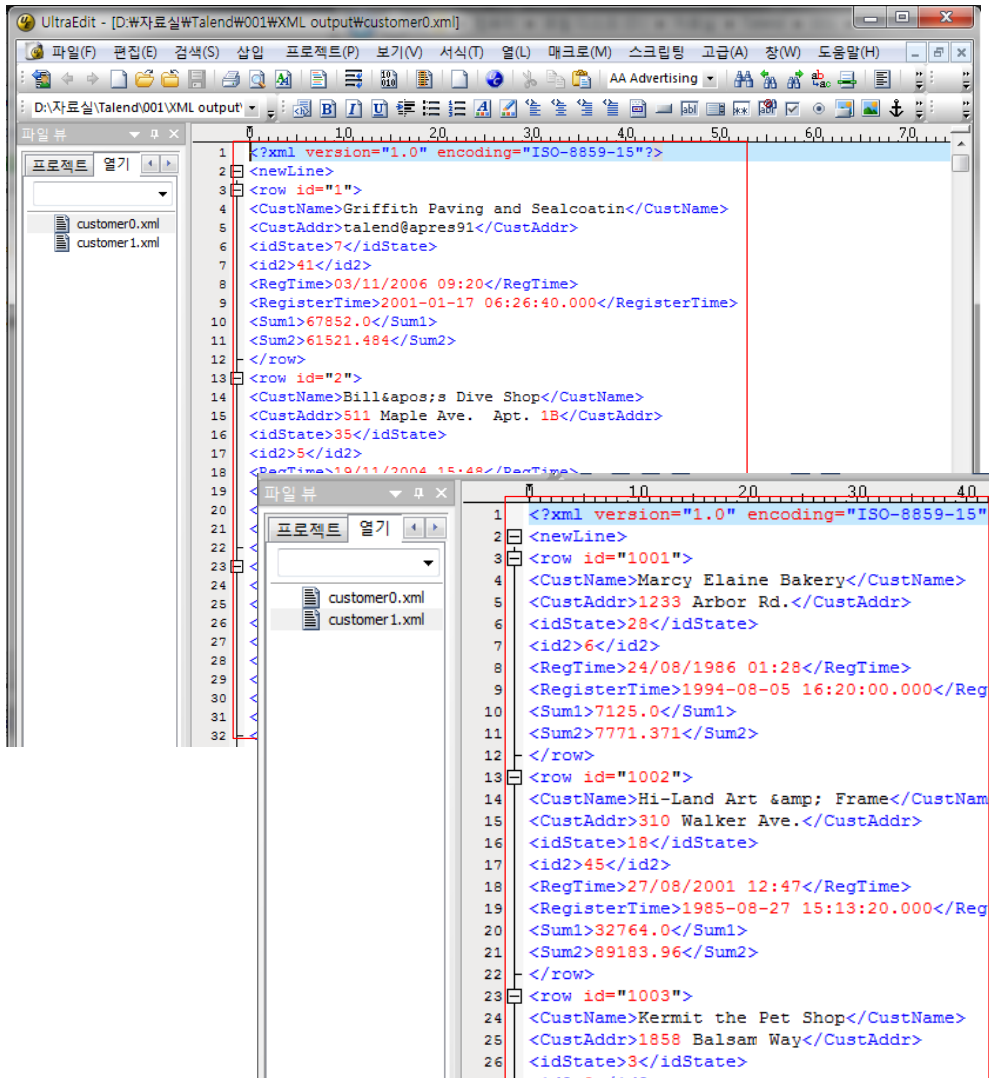
The screenshot displays the Talend Open Studio interface. The main workspace shows a job design with two components: `tFileInputDelimited_1` and `tFileOutputXML_1`, connected by a data flow arrow. Above the arrow, statistics are shown: "5000 rows in 0.49s", "10121.46 rows/s", and "row1 (Main)". Below the workspace, the "Run" button is highlighted with a red rectangle. The "Execution" panel shows the job's execution log, which includes the start time, statistics, and end time. A "Find Errors in Jobs" dialog box is open in the bottom left corner, displaying a table of errors.

Resource	Description
HwoToDesignJobUsingMetadata	
General	
	customer cannot be resolved to a variable
	customer cannot be resolved to a variable

1. Run 패널의 Run 버튼을 클릭하여 정상적으로 실행되는지 확인한다. 정상적으로 실행되면 화면과 같은 결과가 표시되고, 오류가 발생한 경우 오류 팝업창이 나타난다.

### 3. Job Design 실습

#### 3.1 CSV -> XML 파일로 변환

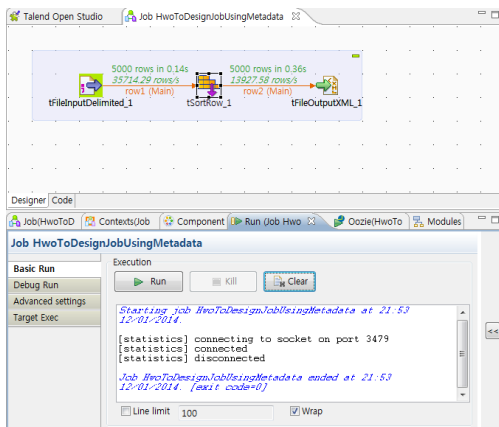
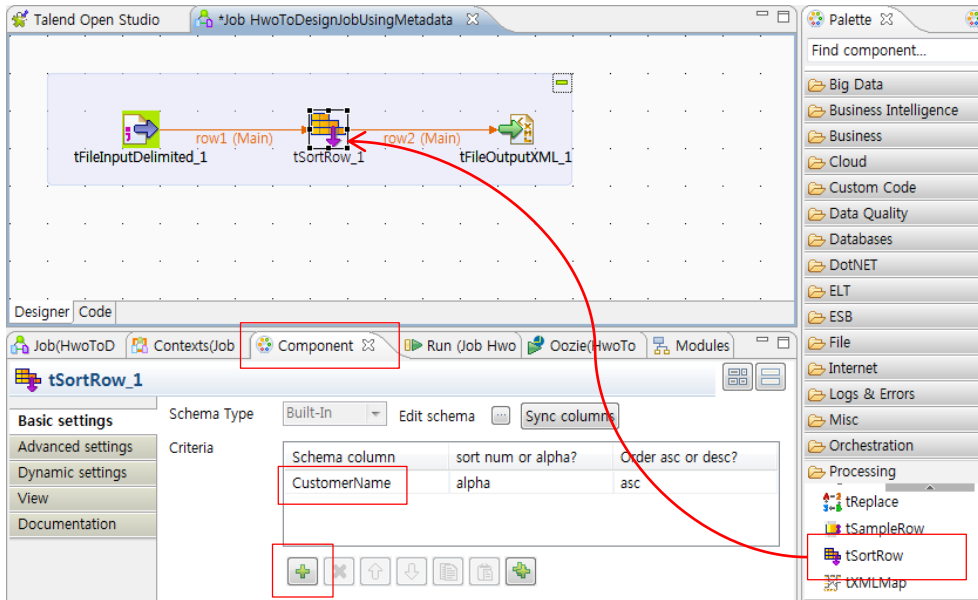


```
<?xml version="1.0" encoding="ISO-8859-15"?>
<newLine>
<row id="1">
  <CustName>Griffith Paving and Sealcoat</CustName>
  <CustAddr>talend@apres91</CustAddr>
  <idState>7</idState>
  <id2>41</id2>
  <RegTime>03/11/2006 09:20</RegTime>
  <RegisterTime>2001-01-17 06:26:40.000</RegisterTime>
  <Sum1>67852.0</Sum1>
  <Sum2>61521.484</Sum2>
</row>
<row id="2">
  <CustName>Bill's Dive Shop</CustName>
  <CustAddr>511 Maple Ave. Apt. 1B</CustAddr>
  <idState>35</idState>
  <id2>5</id2>
  <RegTime>19/11/2004 15:48</RegTime>
  <RegisterTime>1994-08-05 16:20:00.000</RegisterTime>
  <Sum1>7125.0</Sum1>
  <Sum2>7771.371</Sum2>
</row>
<row id="1001">
  <CustName>Marcy Elaine Bakery</CustName>
  <CustAddr>1233 Arbor Rd.</CustAddr>
  <idState>28</idState>
  <id2>6</id2>
  <RegTime>24/08/1986 01:28</RegTime>
  <RegisterTime>1994-08-05 16:20:00.000</RegisterTime>
  <Sum1>7125.0</Sum1>
  <Sum2>7771.371</Sum2>
</row>
<row id="1002">
  <CustName>Hi-Land Art & Frame</CustName>
  <CustAddr>310 Walker Ave.</CustAddr>
  <idState>18</idState>
  <id2>45</id2>
  <RegTime>27/08/2001 12:47</RegTime>
  <RegisterTime>1985-08-27 15:13:20.000</RegisterTime>
  <Sum1>32764.0</Sum1>
  <Sum2>89183.96</Sum2>
</row>
<row id="1003">
  <CustName>Kermit the Pet Shop</CustName>
  <CustAddr>1858 Balsam Way</CustAddr>
  <idState>3</idState>
```

1. tFileOutputXML에서 지정한 저장경로에 가면 customer0.xml ~ customer4.xml 5개의 xml 파일이 생성되어 있고 내용을 확인해 보면 하나의 파일에 id가 1000개씩 저장되어 있다.

### 3. Job Design 실습

#### 3.2 CSV -> XML 파일로 변환시 정렬하기



개요 : XML파일로 생성하기 전에 고객명으로 정렬한 후 XML 파일로 생성한다.

1. CSV 파일을 XML 파일로 변환시 고객명으로 정렬하여 저장한다
2. Palette/Processing/tSortRow를 드래그하여 row1에 드롭한다
3. tSortRow 이 추가되고 Component에서 "+"를 눌러서 정렬할 컬럼을 추가하고 콤보박스에서 정렬할 컬럼인 CustName과 숫자, 문자 중 alpha를 선택한다.
4. 오름차순은 asc, 내림차순은 desc를 선택하면 된다.
5. Run 패널의 Run 버튼을 클릭하여 정상적으로 실행되는지 확인한다.

### 3. Job Design 실습

#### 3.2 CSV -> XML 파일로 변환시 정렬하기

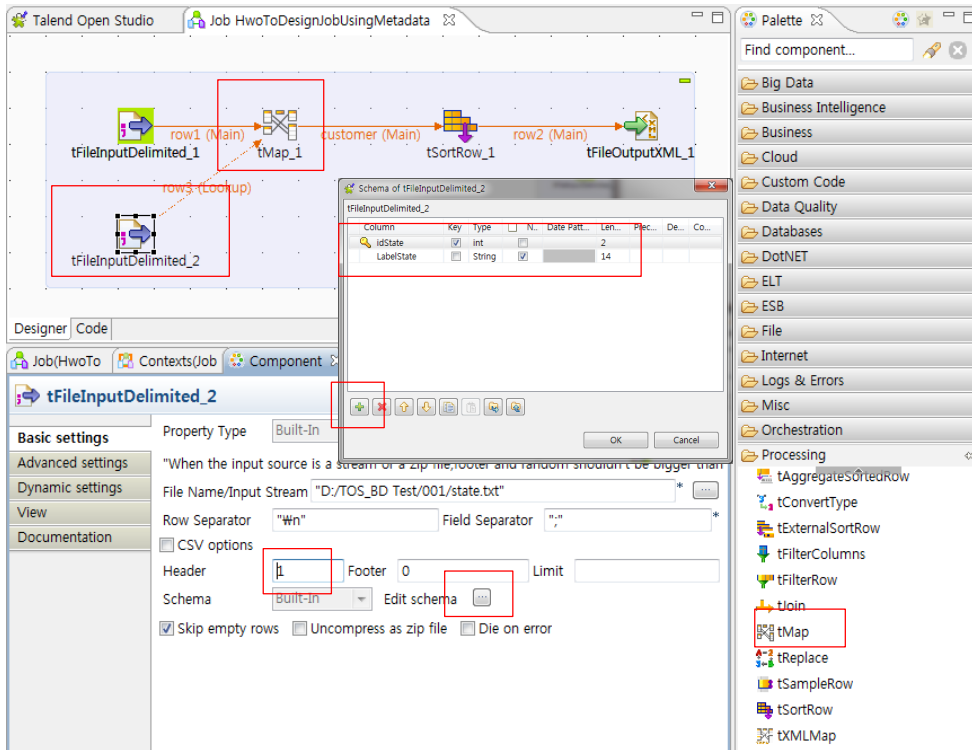
```
<?xml version="1.0" encoding="ISO-8859-15"?>
<customer>
<row id="224">
<CustomerName>AA Advertising</CustomerName>
<CustomerAddr>35 Glenlake Dr.</CustomerAddr>
<idState>13</idState>
<id2>30</id2>
<RegTime>18/05/2004 04:47</RegTime>
<RegisterTime>2005-03-03 16:28:16.000</RegisterTime>
<Sum1>75136.0</Sum1>
<Sum2>13600.888</Sum2>
</row>
<row id="313">
<CustomerName>AA Advertising</CustomerName>
<CustomerAddr>245 Oakridge Ave.</CustomerAddr>
<idState>7</idState>
<id2>1</id2>
<RegTime>04/09/2001 16:08</RegTime>
<RegisterTime>2002-10-18 12:06:40.000</RegisterTime>
<Sum1>43120.0</Sum1>
<Sum2>78881.51</Sum2>
</row>
<row id="393">
<CustomerName>AA Advertising</CustomerName>
<CustomerAddr>899 Harvard Ct</CustomerAddr>
<idState>38</idState>
<id2>27</id2>
<RegTime>14/03/2005 02:28</RegTime>
<RegisterTime>1987-09-22 10:06:40.000</RegisterTime>
<Sum1>22565.0</Sum1>
<Sum2>94625.33</Sum2>
</row>
<row id="578">
<CustomerName>AA Advertising</CustomerName>
<CustomerAddr>673 Calais Drive</CustomerAddr>
<idState>25</idState>
<id2>34</id2>
<RegTime>01/10/1997 11:00</RegTime>
<RegisterTime>2005-09-14 15:54:56.000</RegisterTime>
<Sum1>24902.0</Sum1>
<Sum2>62266.684</Sum2>
</row>
```

1. tFileOutputXML에서 지정한 저장경로에 가면 customer0.xml ~ customer4.xml 5개의 xml 파일이 생성되어 있고 내용을 확인해 보면 하나의 파일에 id가 1000개씩 저장되어 있고, 정렬순서로 지정한 CustomerName 순서로 정렬되어 있다



### 3. Job Design 실습

#### 3.3 CSV -> XML 파일로 맵핑 후 정렬하기



개요 : State id가 있는 컬럼을 State이름으로 맵핑한 후 고객명으로 정렬하여 xml파일로 생성한다.

1. State 이름이 있는 Input 파일을 사용하기 위하여 Palette-File-Input-tFileInputDelimited Component 선택 후 Job Designer에 드롭하면 tFileInputDelimited\_2가 생성된다. tFileInputDelimited\_2를 클릭한 후 Component Tab으로 이동한다. File Name의 ... 버튼을 클릭하여 state.txt 파일이 들어있는 폴더를 지정한다. Header 필드에는 헤더라인이 1라인이므로 1을 입력한다. Edit schema 의 ... 버튼을 클릭하여 + 버튼을 2번 클릭하여 새로운 라인을 2개 추가 후 State 이름 파일의 스키마를 그림과 같이 입력한다.
2. 고객파일과 State이름을 매핑하기 위하여 Processing-tMap Component를 선택하여 row1 위에 드롭하고, new Output name을 지정하라는 팝업창이 나오면 customer를 입력한다.
3. State 이름이 있는 Input 파일 Component인 tFileInputDelimited\_2를 클릭한 후 tMap\_1에 드래그&드롭하면 row3(Lookup) 라인이 생성된다.



### 3. Job Design 실습

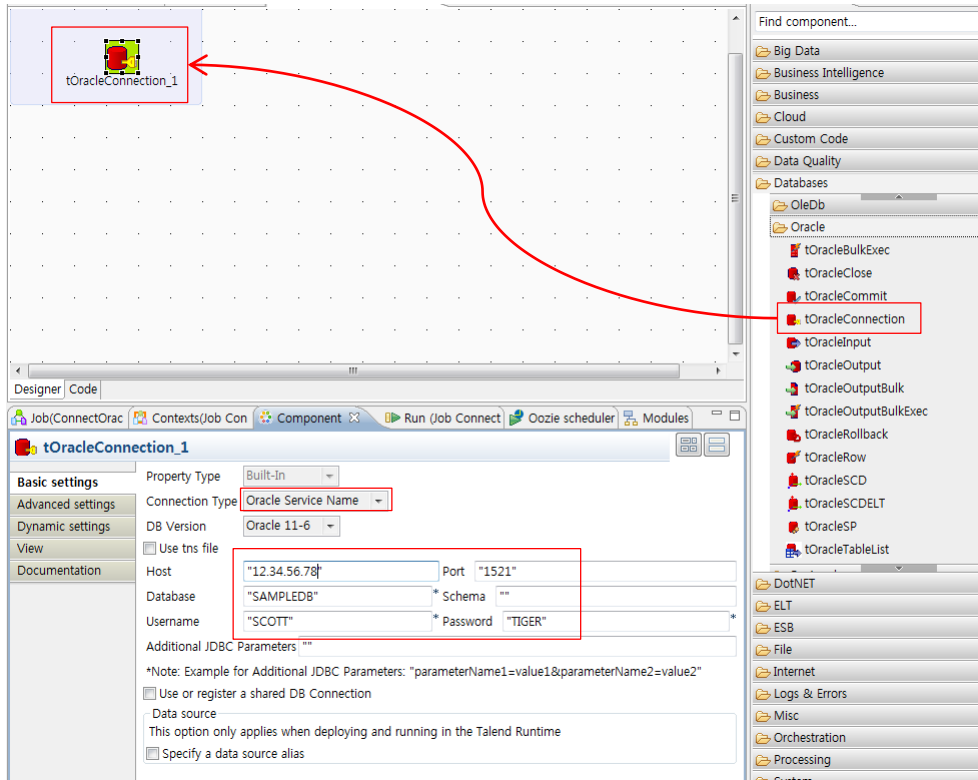
#### 3.3 CSV -> XML 파일로 맵핑 후 정렬하기

```
<?xml version="1.0" encoding="ISO-8859-15"?>
<customer>
<row id="1354">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>750 Homewood Ave.</CustomerAddr>
<LabelState>Hawaii</LabelState>
</row>
<row id="1491">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>43078 Kennedy Dr.</CustomerAddr>
<LabelState>Pennsylvania</LabelState>
</row>
<row id="1722">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>1882 St. Johns</CustomerAddr>
<LabelState>Wisconsin</LabelState>
</row>
<row id="1837">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>1126 Hillcrest Ave</CustomerAddr>
<LabelState>Rhode Island</LabelState>
</row>
<row id="2067">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>14836 W. River Oaks Dr.</CustomerAddr>
<LabelState>Iowa</LabelState>
</row>
<row id="2099">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>1077 Court Ave.</CustomerAddr>
<LabelState>New Hampshire</LabelState>
</row>
<row id="2261">
<CustomerName>Pac Builders</CustomerName>
<CustomerAddr>225 Jeffrys Place</CustomerAddr>
<LabelState>Indiana</LabelState>
</row>
```

1. tFileOutputXML에서 지정한 저장경로에 가면 customer0.xml ~ customer4.xml 5개의 xml 파일이 생성되어 있고 내용을 확인해 보면 하나의 파일에 id가 1000개씩 저장되어 있고, LabelState에 State이름이 추가되어 있는지 확인한다.

## 4. Oracle 접속 실습

### 4.1 Oracle DB 접속하기



개요 : 오라클 A테이블에서 B테이블로 데이터를 추가한다.

1. 오라클에 접속하기 위해서는 팔레트의 Databases/Oracle/tOracleConnection component를 사용한다. 그림과 같이 tOracleConnection을 Job Designer에 드롭한다.
2. 오라클 접속 방식에 따라서 Connection Type을 선택한다. 여기에서는 Oracle Service Name을 선택한 후 Host, Port, Database, Schema, UserName, Password등 오라클 접속 정보를 입력한다.

## 4. Oracle 접속 실습

### 4.1 Oracle DB 접속하기

The screenshot shows the Job Designer interface with the 'tOracleInput\_1' component selected. The 'Basic settings' tab is active, showing the 'Use an existing connection' checkbox checked and the 'Component List' dropdown set to 'tOracleConnection\_1'. The 'Schema' dropdown is set to 'Built-In', and the 'Table Name' is 'AQU\_BAT\_LOG'. The 'Query' field contains a SQL statement. A 'Schema of tOracleInput\_1' dialog box is open at the bottom, displaying a table of column details.

Column	Db Column	Key	Type	DB Type	N.	Date Pattern (Ctrl+Space avail...)	Length	Precis...
LOG_ID	LOG_ID	<input checked="" type="checkbox"/>	Integer	INT	<input checked="" type="checkbox"/>		12	
BAT_ID	BAT_ID	<input type="checkbox"/>	String		<input checked="" type="checkbox"/>		10	
BAT_TP	BAT_TP	<input type="checkbox"/>	String		<input checked="" type="checkbox"/>		10	
BAT_DT	BAT_DT	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
STRT_DTTM	STRT_DTTM	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
END_DTTM	END_DTTM	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
SCC_YN	SCC_YN	<input type="checkbox"/>	String		<input checked="" type="checkbox"/>		1	

1. 테이블의 데이터를 가져오기 위해서 팔레트의 Databases/Oracle/tOracleInput component를 선택한 후 Job Designer에 드롭한다.
2. tOracleInput \_1의 Component 탭에서 tOracleConnection\_1 의 오라클 접속정보를 사용하기 위해서 "Use an existing connect"을 체크한다.
3. 조회할 테이블 정보의 스키마를 편집하기 위해 Edit schema 를 누른 후 팝업창에서 테이블의 정보를 입력한다.
4. Guest Query 버튼을 누르면 해당 Query가 표시된다.

## 4. Oracle 접속 실습

### 4.1 Oracle DB 접속하기

The screenshot shows the Job Designer interface with the following components and settings:

- Component List:** tOracleConnection\_1, tOracleInput\_1, tOracleOutput\_1.
- Basic settings:** ☒ Use an existing connection, Component List: tOracleConnection\_1.
- Advanced settings:** Table: TEST\_LOG.
- Dynamic settings:** Action on table: None, Action on data: Insert.
- View:** Schema: Built-in, Edit schema (highlighted).
- Documentation:** ☐ Die on error.

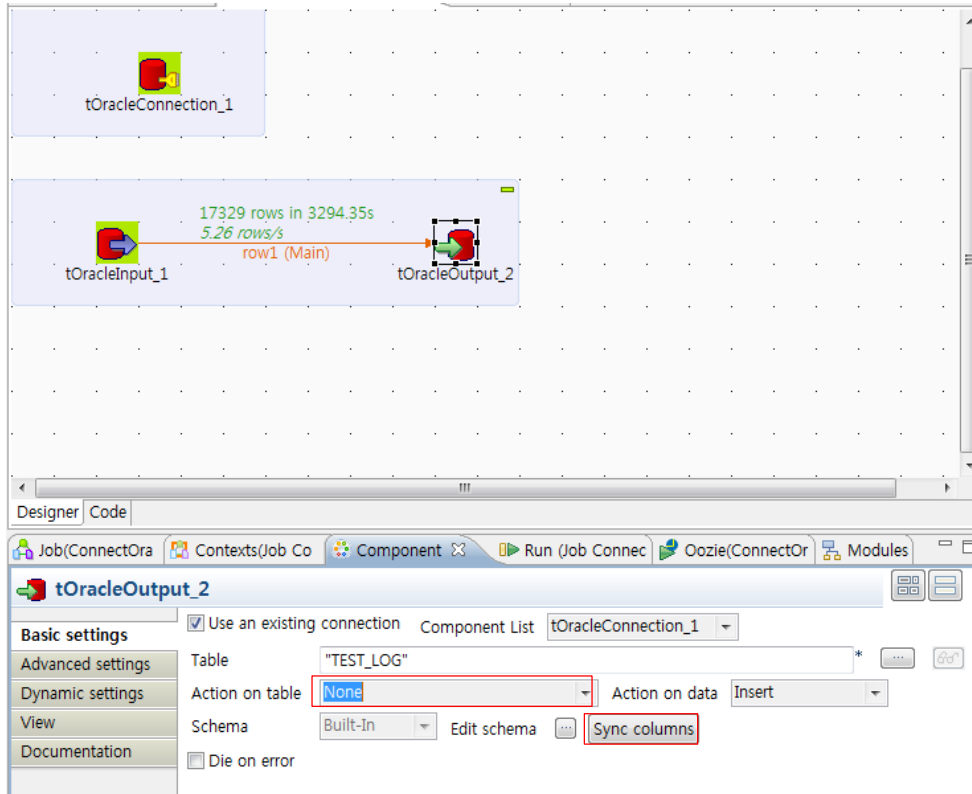
The 'Schema of tOracleOutput\_1' dialog box is open, showing the following table structure:

Column	Db Column	Key	Type	DB Type	N...	Date Pattern (Ctrl+Space ...	Len...	Pre
LOG_ID	LOG_ID	<input checked="" type="checkbox"/>	Int...	INT	<input checked="" type="checkbox"/>		12	
BAT_ID	BAT_ID	<input type="checkbox"/>	Stri...		<input checked="" type="checkbox"/>		10	
BAT_TP	BAT_TP	<input type="checkbox"/>	Stri...		<input checked="" type="checkbox"/>		10	
BAT_DT	BAT_DT	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
STRT_DTTM	STRT_DTTM	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
END_DTTM	END_DTTM	<input type="checkbox"/>	Date	DATE	<input checked="" type="checkbox"/>	"yyyy/MM/dd hh:mm:ss"		
SCC_YN	SCC_YN	<input type="checkbox"/>	Stri...		<input checked="" type="checkbox"/>		1	

1. 데이터를 저장할 테이블을 위해서 팔레트의 Databases/Oracle/tOracleOutput component를 선택한 후 Job Designer에 드롭한다.
2. tOracleOutput\_1의 Component 탭에서 tOracleConnection\_1의 오라클 접속정보를 사용하기 위해서 "Use an existing connect"을 체크한다.
3. 테이블 정보의 스키마를 편집하기 위해 Edit schema 를 누른 후 팝업창에서 테이블의 정보를 입력한다.

## 4. Oracle 접속 실습

### 4.1 Oracle DB 접속하기



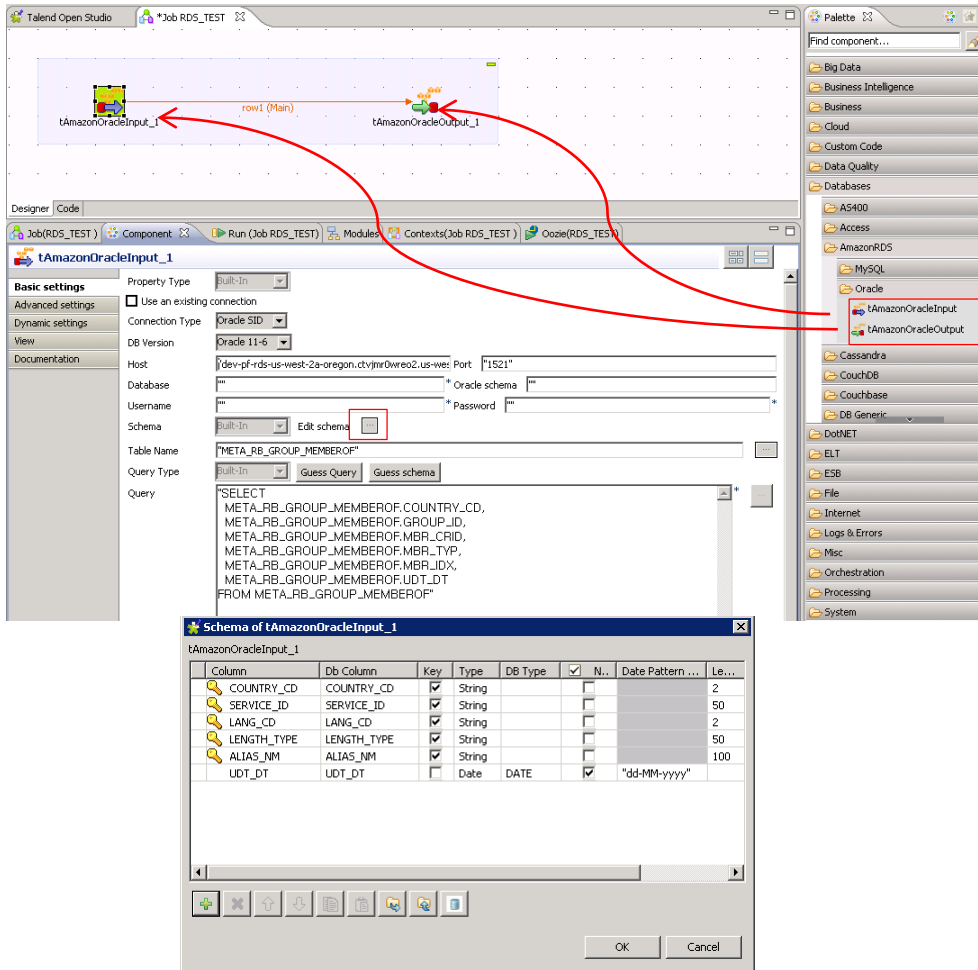
1. 두 컴포넌트를 연결하기 위하여 tOracleInput \_1 에서 마우스 오른쪽 버튼을 클릭한 채로 tOracleOutput \_1 로 드래그한다.
2. 버튼을 놓으면 row1(Main) 연결선이 생긴다. tOracleInput \_1 에서 마우스 오른쪽 버튼을 클릭 후 row-main 을 선택하고 tOracleOutput \_1 을 선택하여도 연결선이 생긴다.
3. Input data와 Output data를 매핑하기 위하여 Sync columns을 누르면 자동으로 매핑하여 준다.
4. Action on table 종류 : Drop and Create table, Create table, Create table if not exists, Drop table if exists and create, Clear table, Truncate table, Truncate table with reuse storage
5. Action on data 종류 : Insert, Update, Insert or update, Update or Insert, Delete
6. Run Tab에서 실행

\* Input Table schema db에서 가져오는 방법 찾기

\* 데이터 추가속도는 초당 5.26건 : 너무 느림!!!

## 4. Oracle 접속 실습

### 4.2 AmazonRDS 접속하기



1. AmazonRDS 접속 컴포넌트 tAmazonOracleInput, tAmazonOracleOutput를 드래그하여 Job Designer에 드롭한다.
2. 두 컴포넌트를 연결하기 위하여 tAmazonOracleInput\_1 에서 마우스 오른쪽 버튼을 클릭한 채로 tAmazonOracleOutput\_1 로 드래그한다.
3. AmazonRDS에 접속하기 위하여 tnsname.ora 에 있는 접속정보를 tAmazonOracleInput\_1, tAmazonOracleOutput\_1 각 필드에 입력한다.
4. Edit schema 를 누른 후 팝업창에서 테이블의 정보를 입력한 후 Guest Query 버튼을 누르면 해당 Query가 표시된다. 일부 데이터만 가져오려는 경우 Query에서 조건절을 추가로 입력한다.

\* 데이터 추가속도는 179,524건 추가시 초당 12,974건 INSERT됨



## 4. Oracle 접속 실습

### 4.3 Amazon RDS - Multi thread execution 사용하기

The screenshot displays the Oozie Job Designer interface. At the top, a workflow diagram shows three components: CM\_STG\_OLD\_BAT, tAmazonOracleInput\_1, and INNO\_TEST. CM\_STG\_OLD\_BAT is connected to tAmazonOracleInput\_1 via a green arrow labeled 'OnSubjobOk'. tAmazonOracleInput\_1 is connected to INNO\_TEST via a red arrow labeled 'row1 (Main)'. Below the workflow, the 'Amazon\_RDS\_Oracle' component settings are visible. The 'Main' tab is selected, showing 'Reload from project settings', 'Save to project settings', and 'Use Project Settings' checked. The 'Extra' tab is also visible, with 'Multi thread execution' checked. The 'Basic settings' tab is selected, showing 'Oracle SID' and 'Oracle 11-g' as the database connection details. The 'Advanced settings' tab is also visible, showing 'Use tns file' checked and 'TNS File' set to 'D:/oracle/product/11.2.0/cdi'. The 'Host' field is set to '1521\*'. The 'Database' field is set to 'Schema'. The '사용자' (User) field is set to '비밀번호'. The 'Additional JDBC Parameters' field is set to 'parameterName1=value1&parameterName2=value2'. The 'Use or register a shared DB Connection' checkbox is checked, and the 'Shared DB Connection Name' is set to 'TIVYSCNT\*'. The 'Documentation' tab is also visible.

1. AmazonRDS 접속 컴포넌트 tAmazonOracleConnection, tAmazonOracleInput, tAmazonOracleOutput를 드래그하여 Job Designer에 드롭한다.
2. 두 컴포넌트를 연결하기 위하여 tAmazonOracleConnection\_1 에서 마우스 오른쪽 버튼을 클릭한 후 trigger/OnSubjobOK 을 클릭 후 tAmazonOracleInput\_1 을 클릭하면 OnSubjobOK 선이 연결된다.
3. Multi thread execution을 사용하기 위하여 Job tab의 Extra화면에서 Multi thread execution을 체크한다.
4. tAmazonOracleConnection\_1 에 접속하기 위하여 tnsname.ora 에 있는 접속정보를 tAmazonOracleConnection\_1 각 필드에 입력한다.

## 4. Oracle 접속 실습

### 4.3 Amazon RDS - Multi thread execution 사용하기

The screenshot displays the Oozie Designer interface for configuring Amazon RDS components. The top diagram shows a workflow: CM\_STG\_OLD\_BAT → OnSubjobOk → tAmazonOracleInput\_1 → row1 (Main) → INNO\_TEST.

**tAmazonOracleInput\_1 Configuration:**

- Basic settings:** ☒ Use an existing connection. Connection: AmazonOracleConnection\_1 - CM\_STG\_OLD\_BAT.
- Advanced settings:** Schema: Built-In. Edit schema: [button].
- Dynamic settings:** 테이블명: "AQU\_POPULAR\_PROGRAM\_HIST". Query Type: Built-In. Guess Query: [button]. Guess schema: [button].
- Query:**

```
SELECT
AQU_POPULAR_PROGRAM_HIST.COUNTRY_CD,
AQU_POPULAR_PROGRAM_HIST.WEEK,
AQU_POPULAR_PROGRAM_HIST.DISPLAY_ORDER,
AQU_POPULAR_PROGRAM_HIST.TIMESLOT,
AQU_POPULAR_PROGRAM_HIST.CHANNEL_NM,
AQU_POPULAR_PROGRAM_HIST.PROGRAM_ID,
AQU_POPULAR_PROGRAM_HIST.START_TM,
AQU_POPULAR_PROGRAM_HIST.TITLE,
AQU_POPULAR_PROGRAM_HIST.WATCH_CNT,
AQU_POPULAR_PROGRAM_HIST.REG_DT
FROM AQU_POPULAR_PROGRAM_HIST"
```
- Schema of tAmazonOracleInput\_1:** A table with columns: COUNTRY\_CD, WEEK, DISPLAY\_ORDER, TIMESLOT, CHANNEL\_NM, PROGRAM\_ID, START\_TM, TITLE, WATCH\_CNT, REG\_DT. Data types include VARCHAR, INT, DATE, and NUMBER.

**tAmazonOracleOutput\_1 Configuration:**

- Basic settings:** ☒ Use an existing connection. Connection: AmazonOracleConnection\_1 - CM\_STG\_OLD\_BAT.
- Advanced settings:** Table: "TEST\_INNO". Action on table: truncate table. Action on data: Insert.
- Dynamic settings:** Warning: this component configuration will automatically generate a commit before data insert/update/delete.
- Schema:** Built-In. Edit schema: [button]. Sync columns: [button].
- Die on error:** ☐

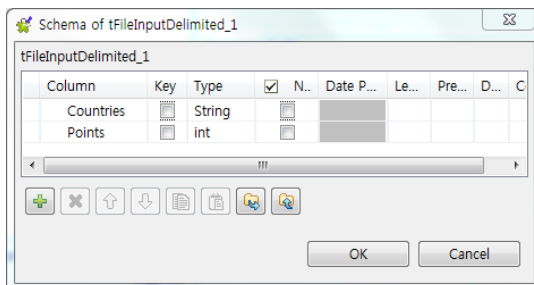
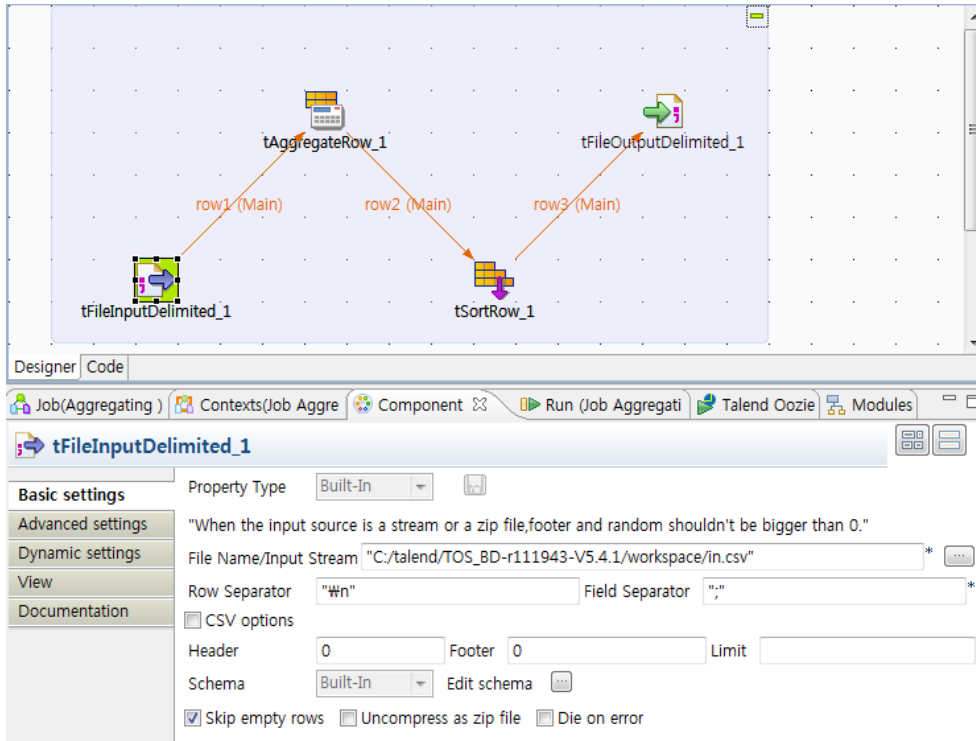
1. tAmazonOracleInput\_1 Component 탭 Basic settings 화면에서 tAmazonOracleConnection\_1을 사용하기 위해서 Use an existing connection을 체크한다.
2. 테이블 정보의 스키마를 편집하기 위해 Edit schema 를 누른 후 팝업창에서 테이블의 정보를 입력한 후 Guest Query 버튼을 누르면 해당 Query가 표시된다.
3. Input data와 Output data를 매핑하기 위하여 Sync columns을 누르면 자동으로 매핑하여 준다.
4. 해당 테이블에 들어있는 데이터를 삭제하고 새로 추가하기 위하여 Action on table은 Truncate table을, Action on data는 Insert를 선택한다.

\* 623,794건 추가 311.75초 소요

2000.94 rows/s

## 5. 컴포넌트 사용 실습

### 5.1 데이터 집계와 정렬하기



국가와 값으로 구성된 CSV파일을 읽어서 국가별로 평균, 최대값, 최소값을 구한 후 정렬하여 CSV파일로 저장한다.

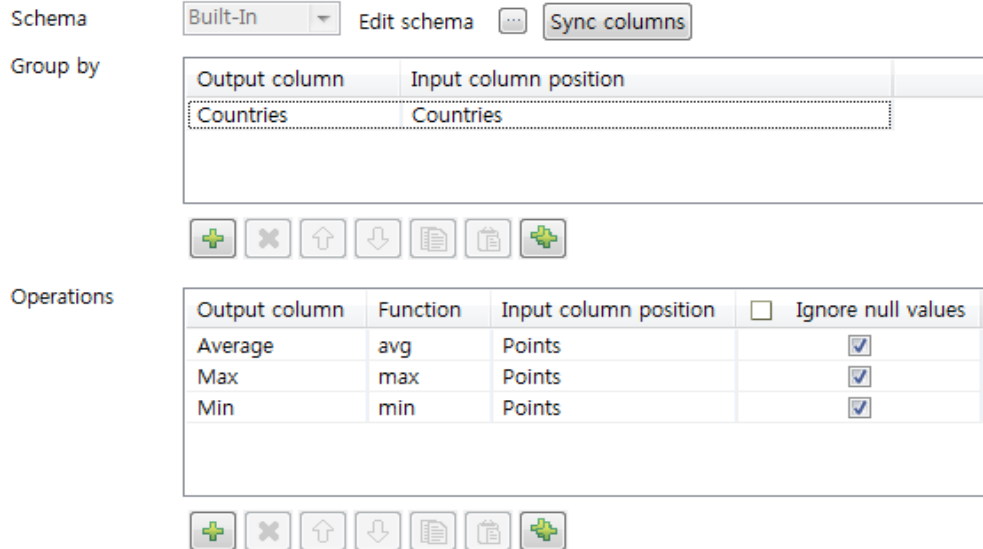
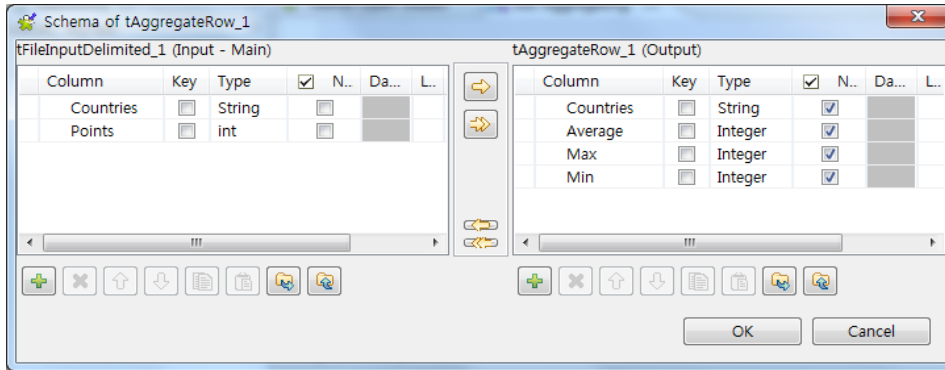
1. 팔레트의 File 폴더에서 tFileInputDelimited 컴포넌트를 design workspace에 drop한다.
2. Basic settings tab 에서 입력할 파일경로를 선택한다
3. Edit schema를 클릭하여 파일구조와 같은 Counties와 Points컬럼을 입력한다.

4. Input data

England:123  
France:5405  
England:123  
France:540  
Germany:937  
Ireland:932  
Italy:2932;  
Spain:3932;  
Korea:4339;  
Germany:935;  
Ireland:932;  
Italy:2934;  
England:123;  
France:540;  
Germany:934;  
Ireland:932;  
England:1232;  
France:5404;  
Germany:937;  
Ireland:9328;  
Italy:29358;  
Spain:393;  
Korea:43397;  
Italy:2938;  
Spain:3935;  
Korea:4339;  
Spain:3939;  
Korea:4339;

## 5. 컴포넌트 사용 실습

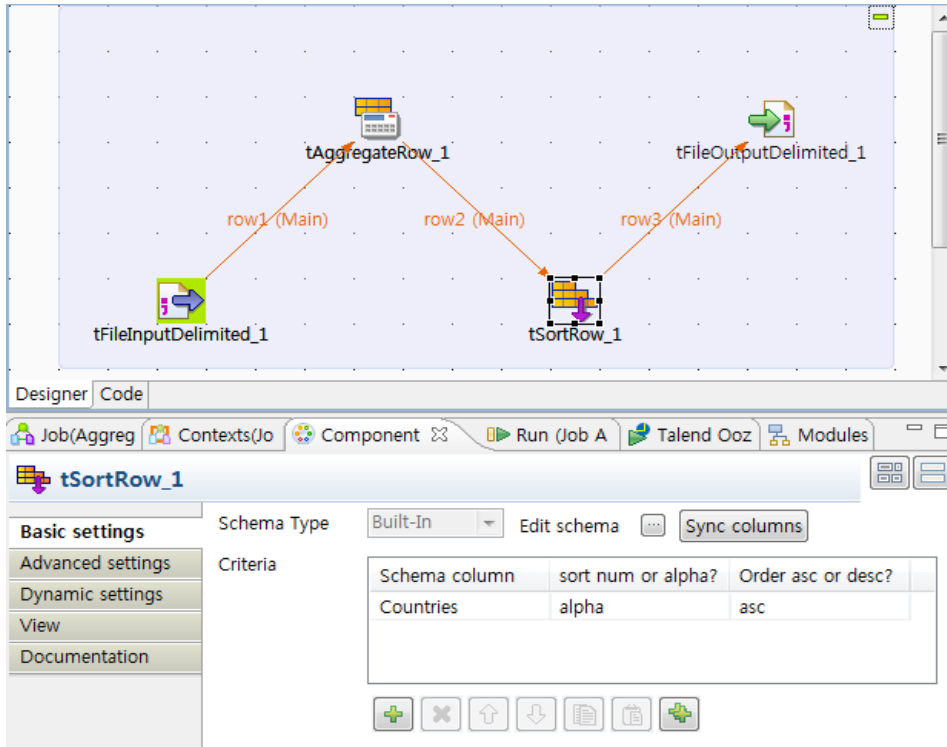
### 5.1 데이터 집계와 정렬하기



1. 팔레트의 Processing 폴더에서 tAggregateRow 컴포넌트를 design workspace에 drop한다.
2. tFileInputDelimited에서 오른쪽 클린한 후 tAggregateRow 에 드롭하면 row1(Main)선이 연결된다
3. 속성값을 정의하기 위하여 tAggregateRow 컴포넌트를 더블클릭한다.
4. Edit schema 버튼을 눌러서 Output 컬럼을 그림과 같이 추가한다
5. 국가별로 그룹을 주기 위하여 Group by의 Output column과 Input column position을 Countries를 선택한다.
6. 아래의 Operations에서 Output컬럼을 선택한 후 국가별로 평균, 최대값, 최소값을 구하기 위하여 이름에 맞는 Function를 선택하고, Input값으로 받을 Points컬럼을 Input column position에서 Points를 선택한다. Null 값은 무시하기 위하여 Ignore null values를 체크한다.

## 5. 컴포넌트 사용 실습

### 5.1 데이터 집계와 정렬하기



Schema of tSortRow\_1

tAggregateRow\_1 (Input - Main)

Column	Key	Type	✓	N.	Dat...	L.
Countries	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		
Average	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		
Max	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		
Min	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		

tSortRow\_1 (Output)

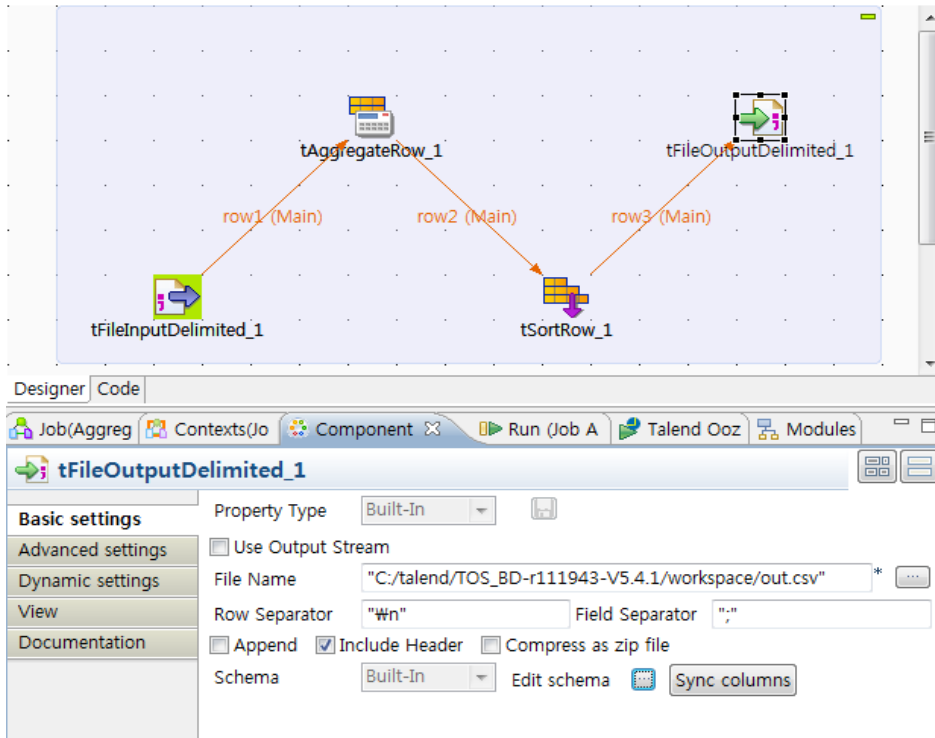
Column	Key	Type	✓	N.	Dat...	L.	P.
Countries	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>			
Average	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>			
Max	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>			
Min	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>			

OK Cancel

1. 팔레트의 Processing 폴더에서 tSortRow 컴포넌트를 design workspace에 drop한다.
2. tAggregateRow 에서 오른쪽 클린한 후 tSortRow 에 드롭하면 row2(Main)선이 연결된다
3. 속성값을 정의하기 위하여 tSortRow 컴포넌트를 더블클릭한다.
4. Input column과 output column을 맞추기 위하여 Sync columns 버튼을 누른 후 Edit schema 버튼을 눌러서 Input, Output 컬럼이 제대로 정의되었는지 확인한다.
5. 국가별로 정렬하기 위하여 Criteria의 Schema column에서 Country를 선택하고 문자, 숫자로 정렬할 지 선택하고, 오름차순, 내림차순을 선택한다.

## 5. 컴포넌트 사용 실습

### 5.1 데이터 집계와 정렬하기



The dialog shows the schema for **tSortRow\_1 (input - Main)** and **tFileOutputDelimited\_1 (Output)**. Both schemas are identical, with columns: Countries, Average, Max, and Min. All columns are of type 'Int...' and have checkboxes for 'N...', 'Date P...', 'Le...', 'Pr...', 'D...', and 'Co...'.

Column	Key	Type	✓	N...	Date P...	Le...	Pr...	D...	Co...
Countries		St...	<input checked="" type="checkbox"/>						
Average		Int...	<input checked="" type="checkbox"/>						
Max		Int...	<input checked="" type="checkbox"/>						
Min		Int...	<input checked="" type="checkbox"/>						

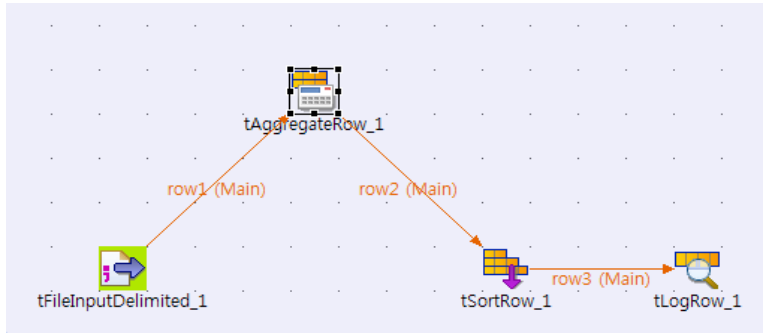
1. 팔레트의 Processing 폴더에서 tFileOutputDelimited 컴포넌트를 design workspace에 drop한다.
2. tSortRow 에서 오른쪽 클릭한 후 tFileOutputDelimited 에 드롭하면 row3(Main)선이 연결된다
3. 속성값을 정의하기 위하여 tFileOutputDelimited 컴포넌트를 더블클릭한다.
4. File Name에서 출력할 폴더와 파일이름을 선택한다.
5. 출력파일에서 헤더를 추가하기 위해서 Include Header를 클릭한다.
6. Input column과 output column을 맞추기 위하여 Sync columns 버튼을 누른 후 Edit schema 버튼을 눌러서 Input, Output 컬럼이 제대로 정의되었는지 확인한다.
7. 실행을 하면 위에서 정의한 파일에 그림과 같은 결과파일이 생성된다.

8. Input data

Countries:	Average:	Max:	Min:
England:	400:	1232:	123
France:	2972:	5405:	540
Germany:	935:	937:	934
Ireland:	3031:	9328:	932
Italy:	9540:	29358:	2932
Korea:	14103:	43397:	4339
Spain:	3049:	3939:	393

## 5. 컴포넌트 사용 실습

### 5.1 데이터 집계와 정렬하기



Schema Type: Built-In Edit schema Sync columns

Mode

☐ Basic

☒ Table (print values in cells of a table)

☐ Vertical (each row is a key/value list)

Execution

Run Kill Clear

```
Starting job Aggregating at 00:01 13/03/2014.
[statistics] connecting to socket on port 3514
[statistics] connected
```

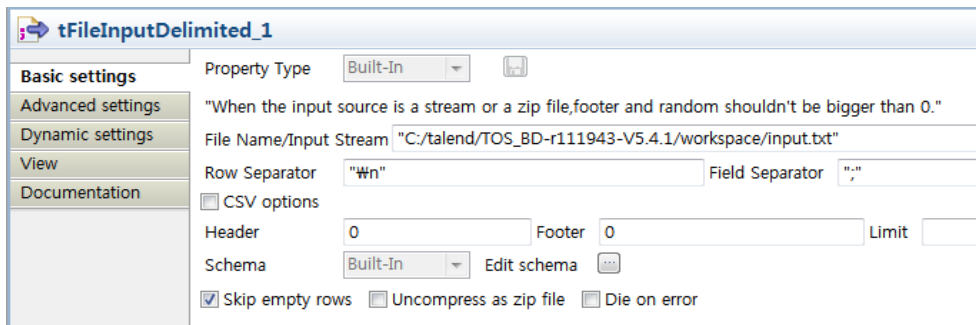
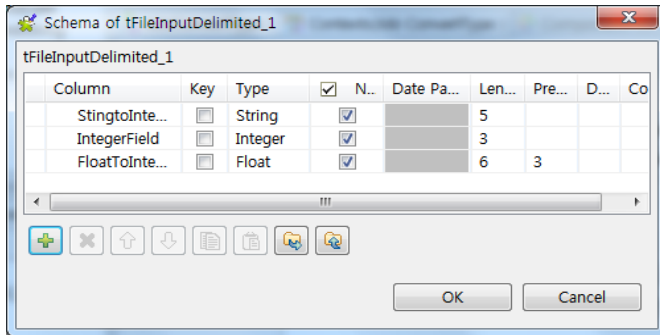
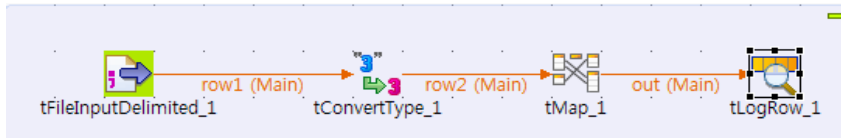
Countries	Average	Max	Min
England	400	1232	123
France	2972	5405	540
Germany	935	937	934
Ireland	3031	9328	932
Italy	9540	29358	2932
Korea	14103	43397	4339
Spain	3049	3939	393

```
[statistics] disconnected
Job Aggregating ended at 00:01 13/03/2014.
[exit code=0]
```

1. 결과를 화면으로 보고 싶으면 팔레트의 Logs&Errors폴더에서 tLogRow 컴포넌트를 design workspace에 drop한다.
2. tSortRow 에서 오른쪽 클릭한 후 tLogRow 에 드롭하면 row3(Main) 선이 연결된다
3. 속성값을 정의하기 위하여 tLogRow 컴포넌트를 더블클릭한다.
4. Input column과 output column을 맞추기 위하여 Sync columns 버튼을 누른 후 Edit schema 버튼을 눌러서 Input, Output 컬럼이 제대로 정의되었는지 확인한다.
5. 출력모드에서 테이블 양식으로 출력하기 위해서 Table을 선택한다.
6. 실행을 하면 그림과 같이 결과화면이 출력된다.

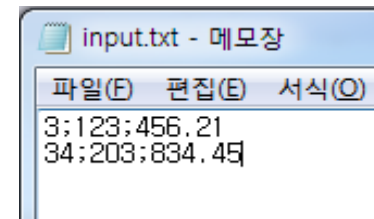
## 5. 컴포넌트 사용 실습

### 5.2 데이터 타입 변환



3개의 컬럼으로 구성된 CSV파일을 읽어서 첫번째 컬럼은 숫자타입으로, 두번째 컬럼과 세번째 컬럼의 합계를 구하여 출력파일로 저장한다.

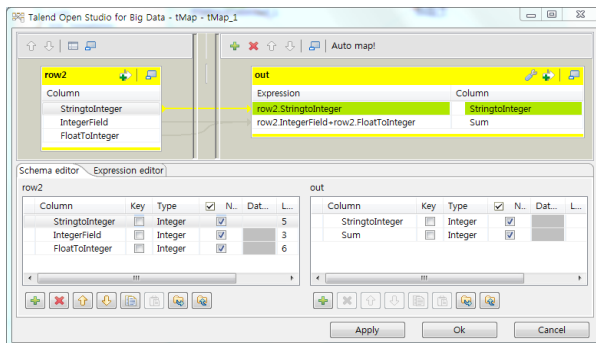
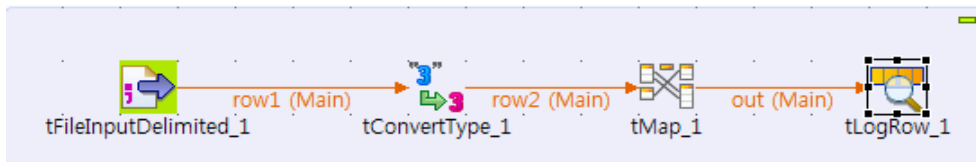
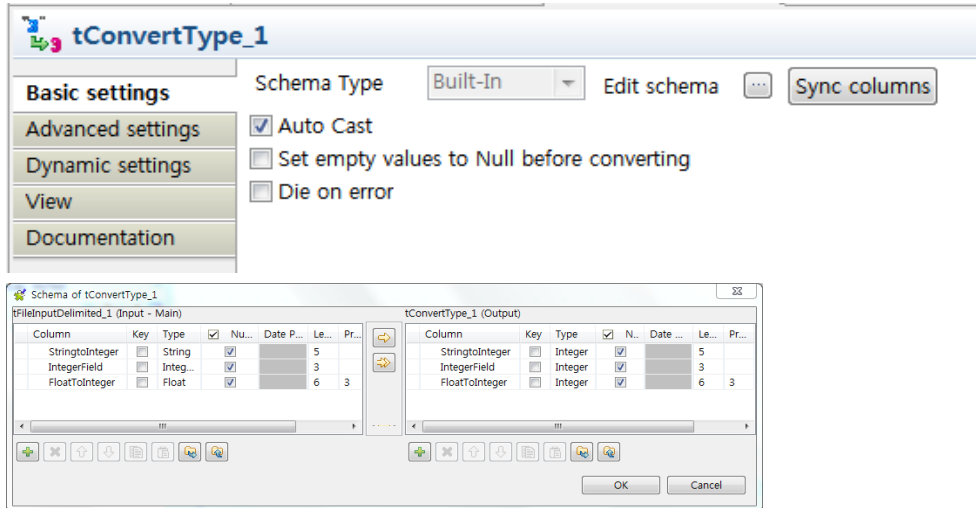
1. 팔레트의 File 폴더에서 tFileInputDelimited 컴포넌트를 design workspace에 drop한다.
2. Basic settings tab 에서 입력할 파일경로를 선택한다
3. Edit schema를 클릭하여 파일구조와 같이 컬럼과 타입을 입력한다.
4. 팔레트의 File 폴더에서 tConvertType\_1컴포넌트를 design workspace에 drop한다.
5. Basic settings tab 에서 입력할 파일경로를 선택한 tFileInputDelimited에서 오른쪽 클릭한 후 tConvertType\_1 에 드롭하면 row1(Main)선이 연결된다
6. Input data





## 5. 컴포넌트 사용 실습

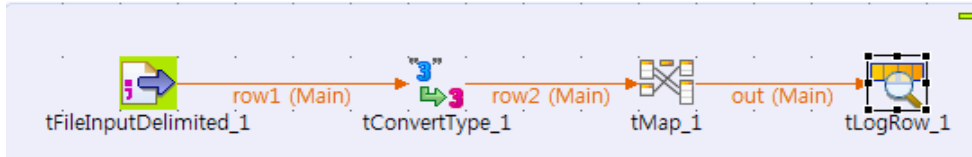
### 5.2 데이터 타입 변환



1. 속성값을 정의하기 위하여 tConvertType\_1 컴포넌트를 더블 클릭한다.
2. Input과 Output 컬럼을 맞추기 위하여 Sync columns를 클릭 후 Edit schema를 누른다
3. 그림과 같이 컬럼이 매칭되는지 확인하고 Output 컬럼의 Type은 Integer로 변경한다.
4. 팔레트의 Processing 폴더에서 tMap컴포넌트를 design workspace에 drop한다.
5. tConvertType\_1 에서 오른쪽 클릭한 후 tMap 에 드롭하면 row2(Main)선이 연결된다.
6. tMap을 더블 클릭한 후 팝업창이 뜨면 그림과 같이 설정한다. Output의 Sum Column은 입력 컬럼 IntegerField와 FloatToInteger의 합계를 구하기 위하여 Expression에 row2.IntegerField + row2.FloatToInteger 를 입력한다.

## 5. 컴포넌트 사용 실습

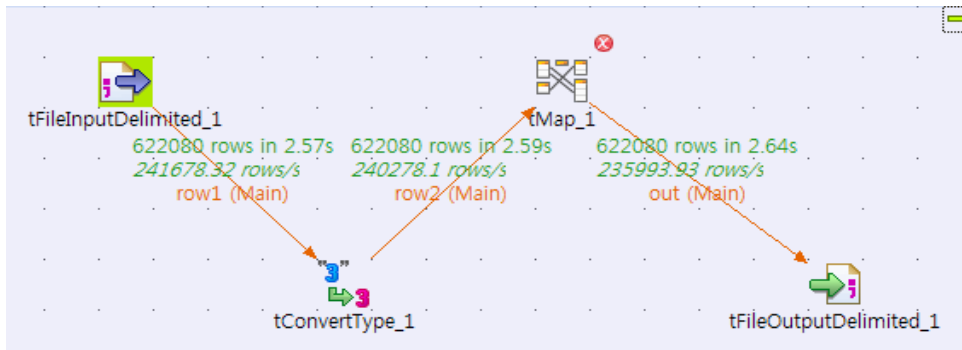
### 5.2 데이터 타입 변환



Schema Type: Built-In Edit schema Sync columns

Mode

- ☐ Basic
- ☒ Table (print values in cells of a table)
- ☐ Vertical (each row is a key/value list)



1. 결과를 화면으로 출력하기 위하여 팔레트의 Logs&Errors폴더에서 tLogRow 컴포넌트를 design workspace에 drop한다.
2. tMap에서 오른쪽 클린한 후 tLogRow 에 드롭하면 out(Main)선이 연결된다
3. 속성값을 정의하기 위하여 tLogRow 컴포넌트를 더블클릭한다.
4. Input column과 output column을 맞추기 위하여 Sync columns 버튼을 누른 후 Edit schema 버튼을 눌러서 Input, Output 컬럼이 제대로 정의되었는지 확인한다.
5. 출력모드에서 테이블 양식으로 출력하기 위해서 Table을 선택한다.
6. 실행을 하면 그림과 같이 결과화면이 출력된다.

Execution

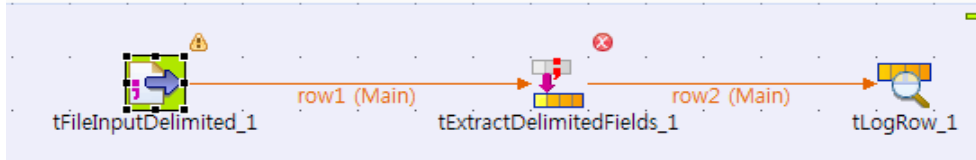
Run Kill Clear

Starting job ConvertType at 01:18 13/03/2014.  
[statistics] connecting to socket on port 3942  
[statistics] connected

tLogRow_1	
StringtoInteger	Sum
2343	579
3465	1037
565	1495
654	15260
6987	38596
87897	7189
241	72036
312	10531
345	1979

## 5. 컴포넌트 사용 실습

### 5.3 한 필드가 콤마로 구분된 파일 추출하기



Property Type: Built-In

File Name/Input Stream: C:/talend/TOS\_BD-r111943-V5.4.1/workspace/test5\*

Row Separator: \n Field Separator: ,

Header: 1 Footer: 0 Limit:

Schema: Built-In Edit schema

☒ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

Column	Key	Type	N.	Date Patt...	Len...	Prec...	De...	Co...
name		String	<input checked="" type="checkbox"/>					

OK Cancel

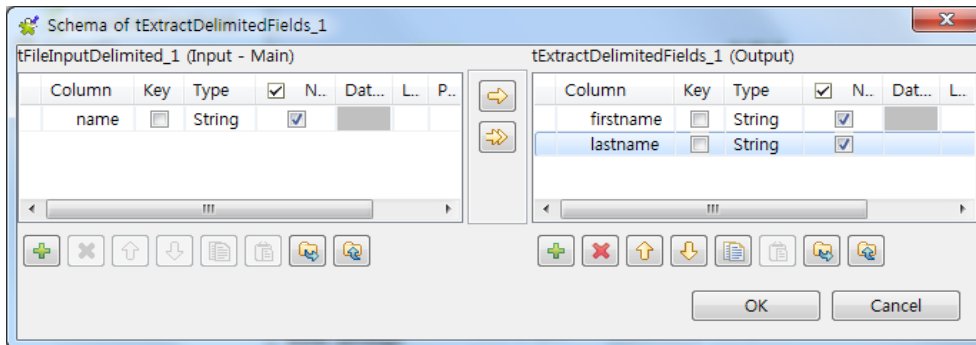
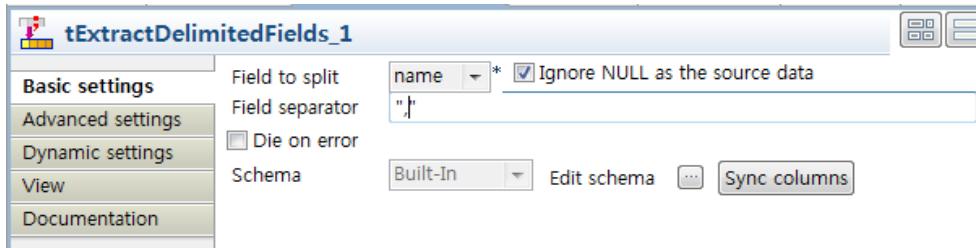
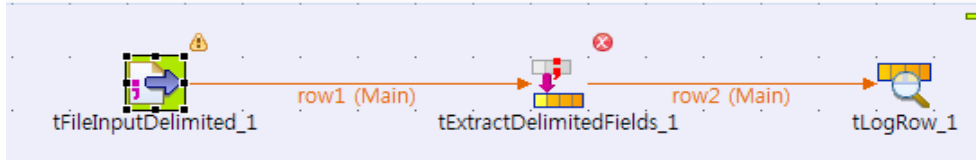
콤마로 구분되어 있는 파일을 읽어서 두 컬럼을 추출한다.

1. 팔레트의 File폴더에서 tFileInputDelimited 컴포넌트를 design workspace에 drop한다.
2. Basic settings tab 에서 입력할 파일경로를 선택한다. Input 파일에 헤더가 1라인 있으므로 Header에 1을 입력한다.
3. Edit schema를 클릭하여 파일구조와 같이 컬럼과 타입을 입력한다.
4. Input data

```
test5.txt - 메모장
파일(F) 편집(E) 서식(O) 보기(V)
firstname, lastname; number
Janet, Aderson; 19988
Martin, Chairman; 9889
Lily, Massy; 9988
```

## 5. 컴포넌트 사용 실습

### 5.3 한 필드가 콤마로 구분된 파일 추출하기



1. 팔레트의 Processing/Fields폴더에서 **tExtractDelimitedFields** 컴포넌트를 design workspace에 drop한다.
2. tFileInputDelimited에서 오른쪽 클린한 후 **tExtractDelimitedFields**에 드롭하면 row1(Main)선이 연결된다
3. 속성값을 정의하기 위하여 **tExtractDelimitedFields** 컴포넌트를 더블 클릭한다.
4. 분리할 필드명을 name을 선택하고, 필드 구분자는 "," 를 입력한다.
5. Edit schema를 클릭하여 Output 에 firstname, lastname을 추가한다.
6. 실행을 하면 그림과 같이 firstname과 lastname이 구분되어 출력된다.

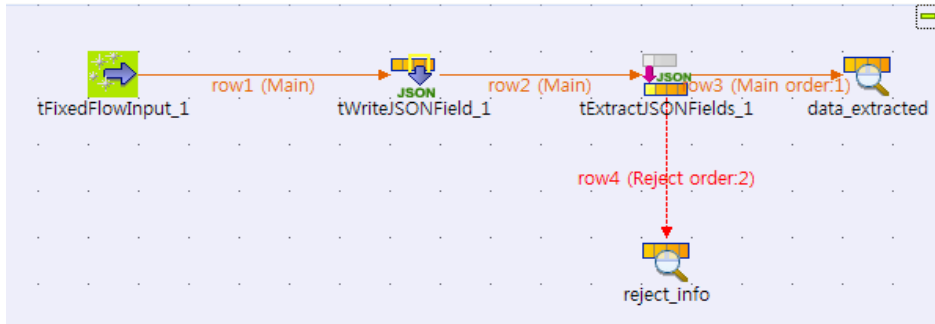
```
Starting job ExtractDelimitedFields at 01:56
13/03/2014.
[statistics] connecting to socket on port 3646
[statistics] connected

+-----+
| tLogRow_1 |
+-----+
| firstname|lastname|
+-----+
| Janet    |Aderson  |
| Martin   |Chairman |
| Lily     |Massy    |
+-----+

[statistics] disconnected
Job ExtractDelimitedFields ended at 01:56
13/03/2014. [exit code=0]
```

## 5. 컴포넌트 사용 실습

### 5.4 Error messages 출력하기



**tFixedFlowInput\_1**

Basic settings

Schema Type: Built-In Edit schema

Advanced settings

Number of rows: 1

Dynamic settings

View

Documentation

Mode

☐ Use Single Table

☐ Use Inline Table

☒ Use Inline Content(delimited file)

Row Separator: "\n" \* Field Separator: ","

Content: Andrew;Wallace;Doc  
John;Smith;R&D  
Christian;Dior;Sales

Schema of tFixedFlowInput\_1

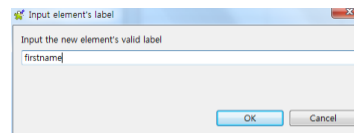
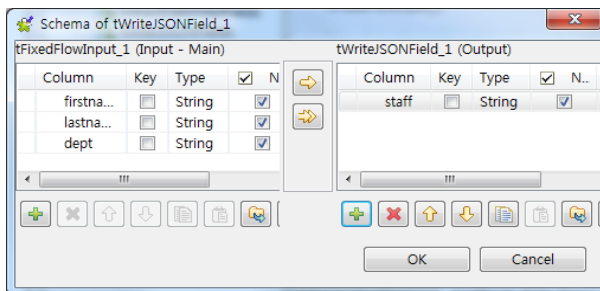
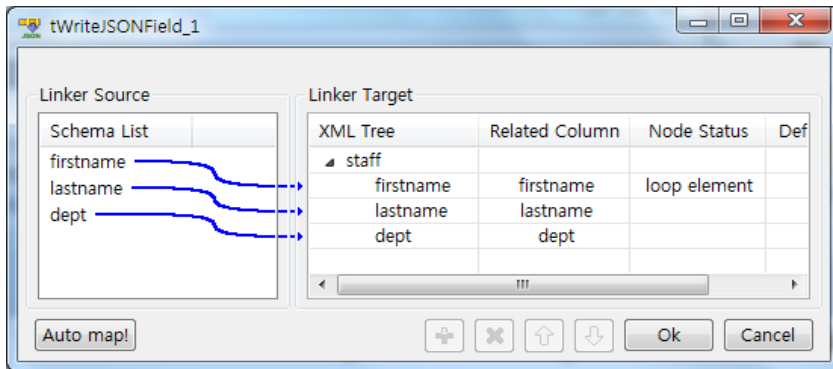
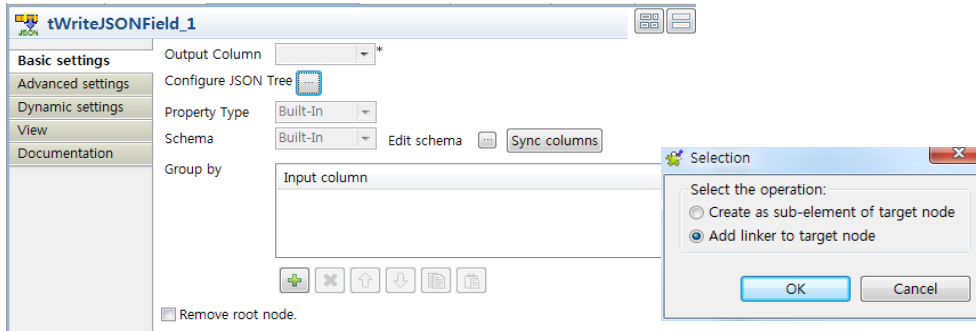
Column	Key	Type	✓	N..	Date...	L...	P...
firstname		String		✓			
lastname		String		✓			
dept		String		✓			

OK Cancel

1. 팔레트의 Misc폴더에서 tFixedFlowInput 컴포넌트를, Misc폴더에서 tWriteJSONField를 Processing/Fields폴더에서 tExtractJSONFields 컴포넌트를, Logs&Errors폴더에서 tLogRow 2개를 design workspace에 drop한다.
2. 그림과 같이 각 컴포넌트를 Main으로 연결하고, tExtractJSONFields and reject\_info 는 Reject로 연결한다.
3. tFixedFlowInput\_1의 Basic settings tab 에서 Use Inline Content를 선택하고, content 내용을 그림과 같이 입력한다.
4. Edit schema를 클릭하여 파일구조와 같이 컬럼과 타입을 입력한다.

## 5. 컴포넌트 사용 실습

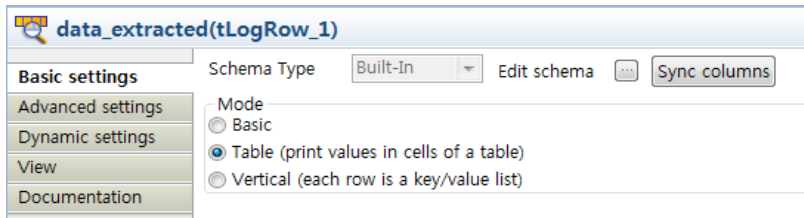
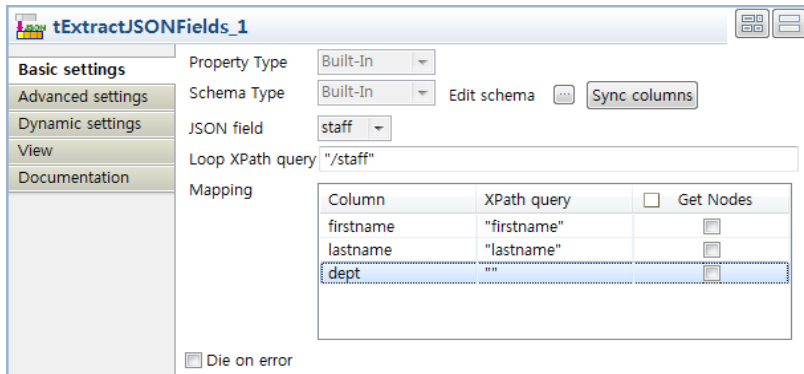
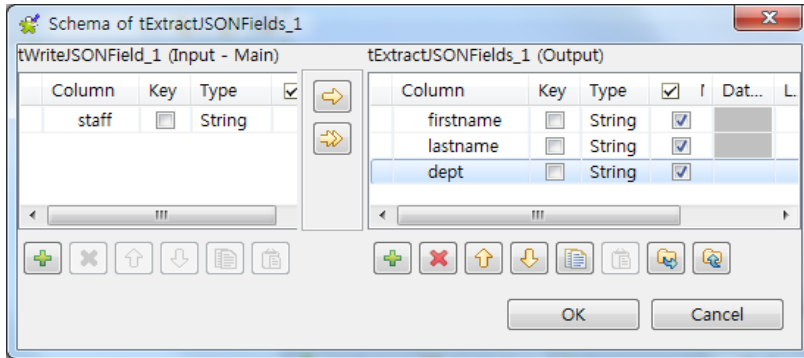
### 5.4 Error messages 출력하기



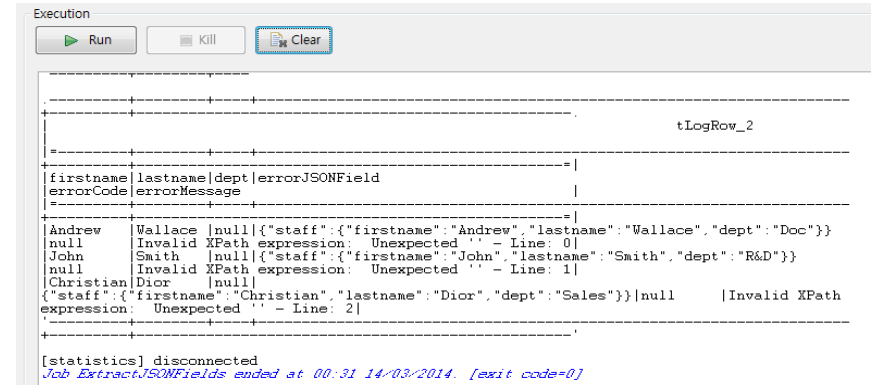
1. Basic settings tab 에서 Configure JSON Tree 를 클릭하여 XML tree editor를 연다
2. tFixedFlowInput 스키마 창의 Linker target 패널에서 XML Tree에 rootTag로 되어 있는 부분을 staff로 수정한다. 이 부분이 JSON field의 루트 노드이다
3. Staff에서 오른쪽 클릭하고 context menu 에서 Add Sub-element 를 선택한다.
4. 팝업창에서 sub-node 이름으로 firstname를 입력한다. 같은 방법으로 lastname과 dept.를 입력한다.
5. firstname 에서 오른쪽 클릭 후 나온 context menu 에서 Set As Loop Element 를 선택한다.
6. Linker source에서 firstname 을 Linker target 의 firstname 에 드롭한다.
7. Selection pop-up dialog box가 나타나면 Add linker to target node 를 선택한다. 같은 방법으로 lastname과 dept도 연결한다.
8. Edit schema을 눌러서 오른쪽 패널 tWriteJSONField\_1에서 JSON 데이터를 가져올 필드로 staff라고 입력한다.

## 5. 컴포넌트 사용 실습

### 5.4 Error messages 출력하기

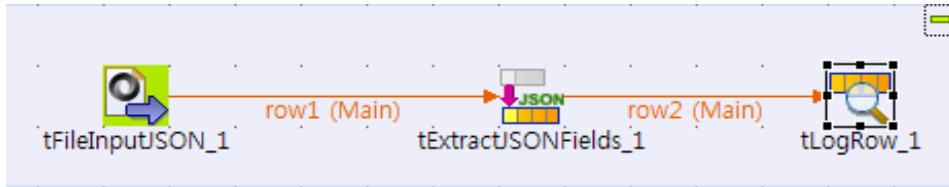


1. tExtractJSONFields\_1의 Basic settings tab 에서 Edit schema 창에서 오른쪽 판넬에 *firstname, lastname, dept*를 입력한다. 이 부분은 *JSON 필드 staff*의 노드가 된다
2. Loop XPath query field에서 JSON data의 루트 노드로 *"/staff"*를 입력한다.
3. Mapping 부분의 XPath query 컬럼에 JSON data 의 노드이름을 입력한다. Dept 부분의 ""은 에러 처리부분으로 빠질것이다.
4. data\_extracted 의 Basic settings에서 Mode 를 결과를 보기 좋게 볼 수 있도록 Table로 선택한다.
5. 실행을 하게 되면 그림과 같이 extraction에 실패한 reject된 부분의 원인과 상세 에러 메시지가 출력된다.



## 5. 컴포넌트 사용 실습

### 5.5 소셜 네트워크에서 데이터 추출하기



Column	JSONPath query
friends	\$.user.friends[*]

```
{
  "user": {
    "id": "9999912388",
    "name": "Kelly Clarkson",
    "friends": [
      {
        "name": "Tom Cruise",
        "id": "55555555555555",
        "likes": {
          "data": [
            {
              "category": "Movie",
              "name": "The Shawshank Redemption",
              "id": "103636093053996",
              "created_time": "2012-11-20T15:52:07+0000"
            },
            {
              "category": "Community",
              "name": "Positive contribution",
              "id": "4713689265413",
              "created_time": "2012-12-16T21:13:26+0000"
            }
          ]
        }
      }
    ]
  },
  "name": "Tom Hanks",
  "id": "8888888888888888",
  "likes": {
    "data": [
      {
        "category": "Journalist",
        "name": "Janelle Wang",
        "id": "135009829148951",
        "created_time": "2013-01-01T08:22:17+0000"
      },
      {
        "category": "Tv show",
        "name": "Now With Alex Wagner",
        "id": "305948749433410",
        "created_time": "2012-11-20T06:14:10+0000"
      }
    ]
  }
}
```

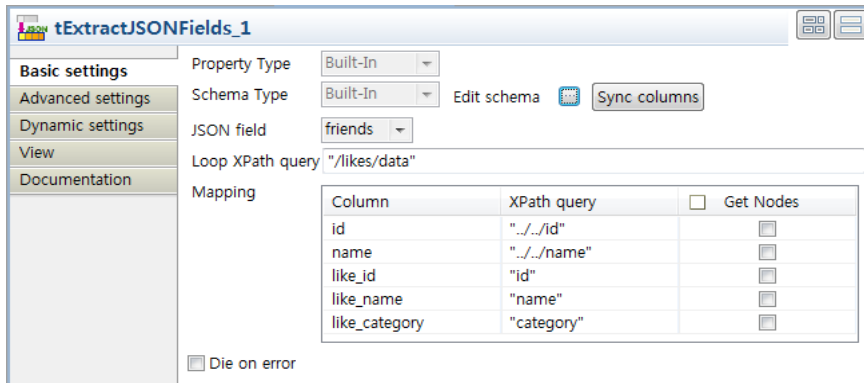
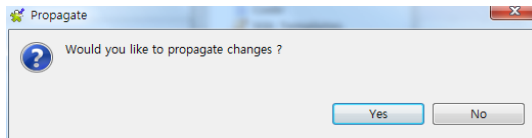
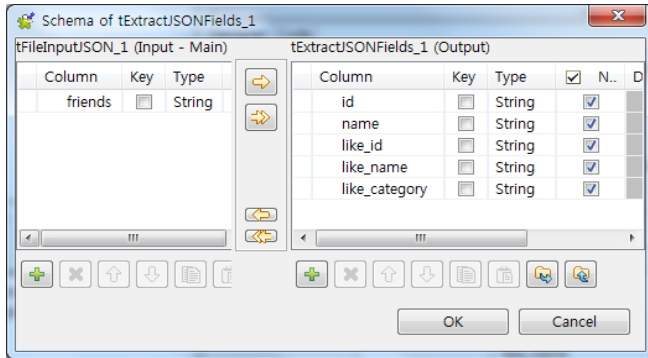
Column	Key	Type	N...	Dat...	L...	P...
friends		String	<input checked="" type="checkbox"/>			

1. 팔레트의 File/Input폴더에서 tFileInputJSON 컴포넌트, Processing/Fields폴더에서 tExtractJSONFields 컴포넌트, Logs&Errors폴더에서 tLogRow를 각각 design workspace에 drop한다.
2. 그림과 같이 각 컴포넌트를 Main으로 연결한다.
3. tFileInputJSON\_1 의 Basic settings tab 에서 입력할 파일경로를 선택한다. 선택한 JSON file *facebook.json* 의 내용을 그림과 같다
4. Edit schema를 클릭하여 컬럼에 friends를 입력한다.
5. Mapping table의 friends 컬럼 오른쪽 JSONPath query에 *\$.user.friends[\*]*를 입력한다.
6. Read by XPath check box는 선택하지 않는다.

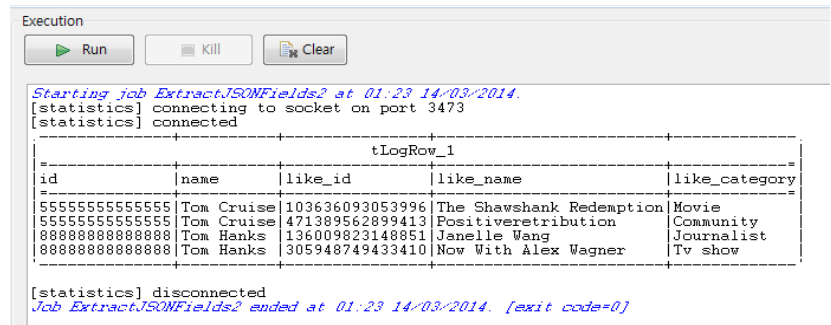


## 5. 컴포넌트 사용 실습

### 5.5 소셜 네트워크에서 데이터 추출하기

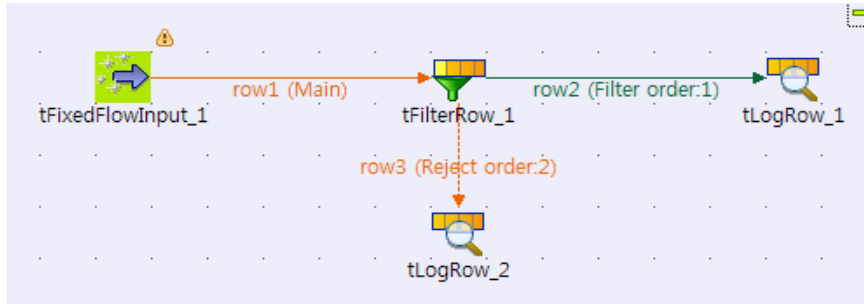


1. tExtractJSONFields 의 Basic settings에서 Edit schema창을 열어 JSON field *friends*의 노드인 *id*, *name*, *like\_id*, *like\_name*, *like\_category* 을 입력한다.
2. Propagate 팝업창이 나오면 Yes를 클릭한다.
3. Loop XPath query field에 *"/likes/data"*를 입력한다.
4. Mapping 부분의 XPath query column에 그림과 같이 입력한다. *".././id"* 는 *"/friends/id"* 노드, *".././name"* 는 *"/friends/name"* 노드를 가져온다.
5. tLogRow\_1의 Basic settings에서 Mode 를 결과를 보기 좋게 볼 수 있도록 Table로 선택한 후 실행하면 그림과 같은 결과가 출력된다.



## 5. 컴포넌트 사용 실습

### 5.6 이름으로 필터링하기



Schema Type Built-In Edit schema

Number of rows

Mode

☐ Use Single Table

☐ Use Inline Table

☒ Use Inline Content(delimited file)

Row Separator  \* Field Separator  \*

Content

```

romain;m:french;16
roman;m:russian,polish,czech;55
romano;n:italian41
romeo;m:italian;29
romolo;m:italian;0
romulusmlroamn mythology;42
  
```

tFixedFlowInput_1							
Column	Key	Type	<input checked="" type="checkbox"/>	N..	Date Patt...	Len...	P
firstname	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
gender	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
language	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
frequency	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				

컴마로 구분되어 있는 파일을 읽어서 두 컬럼을 추출한다.

1. 팔레트의 Misc폴더에서 tFixedFlowInput컴포넌트, Processing폴더에서 tFilterRow 컴포넌트, Logs\$Errors폴더에서 tLogRow 컴포넌트 2개를 가져와서 design workspace에 drop한다.
2. tFixedFlowInput\_1에서 tFilterRow\_1로는 Row > Main을 이용하여 연결하고, tFilterRow\_1에서 tLogRow\_1로는 Row > Filter를 이용하여 연결한다. tFilterRow\_1에서 tLogRow\_2로는 Row > Reject를 이용하여 연결한다. tLogRow\_2는 필터에서 리젝트된 내용을 출력한다.
3. tFixedFlowInput\_1의 Basic settings View에서 Edit schema를 클릭하여 파일구조와 같이 컬럼과 타입을 입력한다.
4. Mode 부분은 Use Inline Content(delimited file)을 선택하고 샘플 데이터로 그림과 같이 입력한다.

## 5. 컴포넌트 사용 실습

### 5.6 이름으로 필터링하기

Schema Built-In Edit schema Sync columns

Logical operator used to combine conditions And\*

Conditions

Input column	Function	Operator	Value
firstname	Length	<=	6

+ × ↑ ↓ 📄 📋

☒ Use advanced mode

Advanced

```
// code sample : use input_row to define
the condition.
// input_row.columnName1.equals("foo") |||
(input_row.columnName2.equals("bar"))
// replace the following expression by
your own filter condition
input_row.language.equals("italian")
```

```
Starting job FilterRow at 01:40 17/03/2014.
[statistics] connecting to socket on port 3485
[statistics] connected
```

tLogRow_1			
firstname	gender	language	frequency
romeo	m	italian	29
romolo	m	italian	0

```

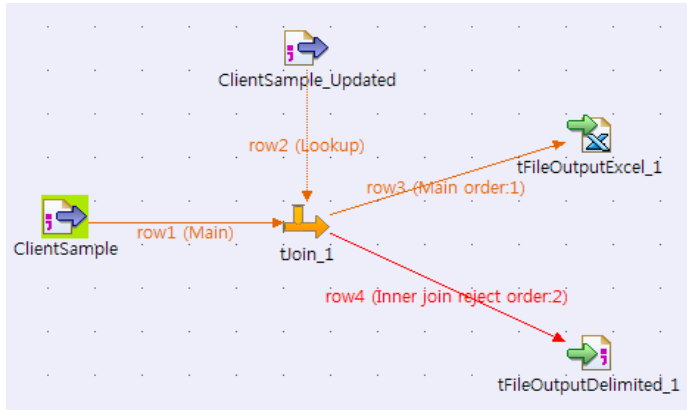
[statistics] disconnected
Job FilterRow ended at 01:40 17/03/2014. [exit code=0]
```

tLogRow_2				
firstname	gender	language	frequency	errorMessage
romain	m	french	16	advanced condition failed
roman	m	russian, polish, czech	55	advanced condition failed
romano	n	italian41	null	advanced condition failed
romulus	m	roamn mythology	42	firstname.length() <= 6 failed advanced condition failed

1. tFilterRow\_1의 Basic settings View에서 Conditions 테이블의 + 버튼을 눌러 1라인을 추가한 후 InputColumn에서 firstname, Function에서는 Length, Operator에서는 <= 을 선택하고, Value 는 6을 입력한다.
2. Language가 italian을 찾기 위해서 Use advanced mode 를 체크하고, Advanced 부분에 input\_row.language.equals("italian")를 입력한다.
3. 실행을 하면 그림과 같이 첫번째 테이블에는 firstname이 6자 이하이고, language가 italian 인 결과가 나타나고, 두번째 테이블에는 리젝트된 내용과 이유가 표시된다.

## 5. 컴포넌트 사용 실습

### 5.7 두개의 파일 조인하여 엑셀과 리젝트 파일로 출력하기



Property Type: Built-In

"When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."

File Name/Input Stream: "C:\talend\TOS\_BD-r111943-V5.4.1\workspace\ClientSample.csv" \*

Row Separator: "\n" Field Separator: ";" \*

☐ CSV options

Header: 1 Footer: 0 Limit:

Schema: Built-In Edit schema ...

☒ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

Schema of ClientSample

Column	Key	Type	✓	N...	Ds
idClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
firstnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
lastnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
ageClient	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		

ClientSample.csv - 메모장

```


파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)
idClient;firstnameClient;lastnameClient;ageClient
1:Dwight;Madison;21
2:Franklin;Jackson;42
3:Ronald;Buchanan;35
4:Bill;Cleveland;49
5:William;Harrison;15
6:William;Fillmore;64
7:Harry;Adams;27
8:Harry;McKintey;39
9:Herbert;Reagan;46
10:Lyndon;Jefferson;25
11:Bill;Jackson;31
12:John;Hayes;22
13:Ulvsses;Reagan;56
28:Jerbert;Eisenhower;42
36:Chester;Grant;42
  
```

컴마로 구분되어 있는 파일을 읽어서 두 컬럼을 추출한다.

1. 팔레트의 File/Input폴더에서 tFileInputDelimited 컴포넌트 2개, Processing폴더에서 tJoin컴포넌트, File/Output 폴더에서 tFileOutputExcel, tFileOutputDelimited 컴포넌트를 가져와서 design workspace에 drop한다.
2. tFileInputDelimited\_1은 ClientSample, tFileInputDelimited\_2는 ClientSample\_Updated 로 이름을 변경한다
3. 여기서 입력으로 사용하는 ClientSample File은 네개의 컬럼으로 되어 있고, *firstnameClient* 과 *lastnameClient* 이 포함되어 있고 이 두 컬럼들이 정확하게 매치되도록 할 것이다
4. ClientSample에서 tJoin\_1로 Main link를 이용하여 연결한다. ClientSample\_Update에서 tJoin\_1로 Main link를 이용하여 연결하면 Lookup link로 연결된다
5. tJoin\_1에서 tFileOutputExcel\_1로 Main link, tFileOutputDelimited\_1로 Inner join reject link로 연결한다.
6. ClientSample의 의 Basic settings View에서 File Name에서 InputFile을 선택하고, 헤더라인이 1이므로 Header는 1을 입력한다.
7. Edit schema를 클릭하여 세번째 그림과 같이 입력필드를 정의한다.
8. ClientSample.csv파일의 내용은 네번째 그림과 같다

## 5. 컴포넌트 사용 실습

### 5.7 두개의 파일 조인하여 엑셀과 리젝트 파일로 출력하기

Property Type Built-In 


"When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."

File Name/Input Stream "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/join\_input.csv" \*

Row Separator "\n" Field Separator "," \*

☐ CSV options

Header 1 Footer 0 Limit

Schema Built-In Edit schema 

☒ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

Schema of ClientSample\_Updated

Column	Key	Type	<input checked="" type="checkbox"/>	N..	Da
firstnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
lastnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		


join\_input.csv - 메모장

파일(F) 편집(E) 서식(O) 보기(V) ?

```
|firstnameClient;lastnameClient  
Jerbert;Eisenhower  
Chester;Grant
```

Schema of tJoin\_1







Column	Key	Type	<input checked="" type="checkbox"/>	N..	Da
idClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
firstnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
lastnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
ageClient	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		

Schema Built-In Edit schema 

☐ Include lookup columns in output

Key definition

Input key attribute	Lookup key attribute
firstnameClient	row2.firstnameClient
lastnameClient	row2.lastnameClient

☒ Inner join ( with reject output )

1. ClientSample\_Update의 Basic settings View에서 File Name에서 Join\_input.csv를 선택하고, 헤더라인이 1이므로 Header는 1을 입력한다.
2. Edit schema를 클릭하여 두번째 그림과 같이 입력필드를 정의한다.
3. Join\_input.csv 파일의 내용은 세번째 그림과 같다
4. tJoin\_1의 Basic settings View에서 Edit schema를 클릭하여 네번째 그림과 같이 입력한다.
5. Basic setting View의 Key definition에서 Input key attribute와 Lookup key attribute를 그림과 같이 지정한다. Input key attribute와 Lookup key attribute가 일치하는 내용은 Excel 파일로 출력이 되고, 그 이외의 데이터는 csv output 파일로 생성할 것이다.
6. 일치하지 않는 내용을 파일로 생성하기 위하여 Inner join (with reject output)를 체크한다.

## 5. 컴포넌트 사용 실습

### 5.7 두개의 파일 조인하여 엑셀과 리젝트 파일로 출력하기

Property Type: Built-In

☐ Write excel2007 file format(xlsx)

☐ Use Output Stream

File Name: "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/join\_out.xls" \*

Sheet name: "Sheet1"

☒ Include header

☐ Append the exist file

☐ Is absolute Y pos.

Font: None

☐ Define all columns auto size

Define column auto size

Column	Auto size
idClient	<input type="checkbox"/>
firstnameClient	<input type="checkbox"/>
lastnameClient	<input type="checkbox"/>

Schema: Built-In Edit schema Sync columns

Schema of tFileOutputExcel\_1

Column	Key	Type	✓	N.
idClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	
firstnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	
lastnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	

tFileOutputExcel\_1 (Output)

Column	Key	Type	✓	N..
idClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	
firstnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	
lastnameClient	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>	

Property Type: Built-In

☐ Use Output Stream

File Name: "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/join\_rejectedout.csv" \*

Row Separator: "\n" Field Separator: ","

☐ Append ☒ Include Header ☐ Compress as zip file

Schema: Built-In Edit schema Sync columns

1. tFileOutputExcel 의 Basic settings View에서 File Name에서 엑셀 파일로 출력할 Join\_out.xls를 입력하고, 헤더라인을 같이 출력하기 위하여 Include header 를 체크한다.
2. Edit schema를 클릭하여 두번째 그림과 같이 입력필드를 정의한다.
3. tFileOutputDelimited 의 Basic settings View에서 File Name에서 리젝트된 내용을 CSV파일로 출력할 Join\_rejectedout.csv를 입력하고, 헤더라인을 같이 출력하기 위하여 Include header 를 체크한다.
4. 실행 후 나온 결과 출력은 아래 그림과 같다.

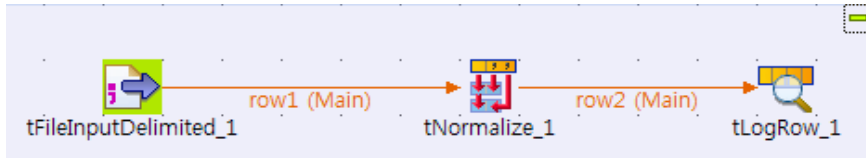
join\_out.xls [호환 모드]

	A	B	C
1	idClient	firstnameClient	lastnameClient
2	28	Jerbert	Eisenhower
3	36	Chester	Grant
4			
5			

idClient:firstnameClient:lastnameClient  
 1:Dwight:Madison  
 2:Franklin:Jackson  
 3:Ronald:Buchanan  
 4:Bill:Cleveland  
 5:William:Harrison  
 6:William:Fillmore  
 7:Harry:Adams  
 8:Harry:McKinley  
 9:Herbert:Reagan  
 10:Lyndon:Jefferson  
 11:Bill:Jackson  
 12:John:Hayes  
 13:Ulysses:Reagan

## 5. 컴포넌트 사용 실습

### 5.8 데이터 표준화하기



Property Type Built-In

"When the input source is a stream or a zip file, footer and random shouldn't be bigger than the header"

File Name/Input Stream 'C:/talend/TOS\_BD-r111943-V5.4.1/workspace/labels\_raw.txt' \*

Row Separator "\n" Field Separator ";"

☐ CSV options

Header 0 Footer 0 Limit

Schema Built-In Edit schema

☒ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

Schema of tFileInputDelimited\_1

Column	Key	Type	✓	N..	Date
Tags	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		

tNormalize\_1

Basic settings

Advanced settings

Dynamic settings

View

☒ Get rid of duplicate rows from output

☐ Use CSV parameters

☒ Discard the trailing empty strings ☒ Trim resulting values

☐ tStatCatcher Statistics

1. 팔레트의 File/Input폴더에서 tFileInputDelimited 컴포넌트, Processing/Field폴더에서 tNormalize 컴포넌트, Logs&Errors 폴더에서 tLogRow 컴포넌트를 가져와서 design workspace에 drop한다.
2. tFileInputDelimited 에서 tNormalize 로 Main link를 이용하여 연결하고, 다시 tNormalize 에서 tLogRow 로 Main link를 이용하여 연결한다
3. tFileInputDelimited의 Basic settings View의 File Name에서 입력파일로 사용할 labels\_raw.txt를 선택한다.
4. Edit schema 창을 열어서 Column에 Tags를 추가한다.
5. tNormalize의 Basic settings 에서 입력 컴포넌트와 스키마를 맞추기 위해서 Sync columns 을 클릭한다.
6. Advanced settings view에서 Get rid of duplicate rows from output, Discard the trailing empty strings, Trim resulting values 를 체크한다.
7. 입력파일 내용은 아래와 같다

labels\_raw.txt - 메모장

```

파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)
l|dap,
db2, jdbc driver,
grid computing, talend architecture ,
content, environment,,
tmap,,
eclipse,
database, java, postgresql,
tmap,
database, java, sybase,
deployment,,
repository,
database, informix, java
    
```

## 5. 컴포넌트 사용 실습

### 5.8 데이터 표준화하기

Schema Type Built-In Edit schema Sync columns

Mode

☐ Basic

☒ Table (print values in cells of a table)

☐ Vertical (each row is a key/value list)

Execution

Run Kill Clear

```
[statistics] connected
tLogRow_1
-----=
Tags
-----=
ldap
db2
jdbc driver
grid computing
talend architecture
content
environment
tmap
eclipse
database
java
postgresql
sybase
deployment
repository
informix
-----=
[statistics] disconnected
```

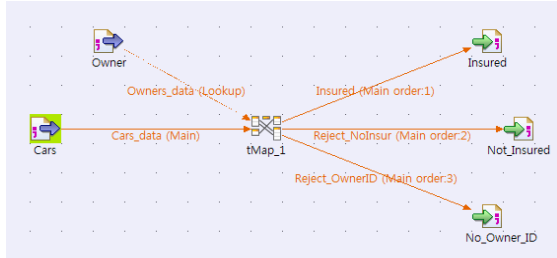
☐ Line limit 100 ☒ Wrap

1. tLogRow 에서 Print values in the cells of table 을 선택한다.
2. 실행을 하면 그림과 같이 입력데이터가 표준화 되어서 출력이 된다.



## 5. 컴포넌트 사용 실습

### 5.9 데이터 매핑과 Reject하기



**Cars(tFileInputDelimited\_1)**

Property Type: Built-in

Advanced settings: "When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."

Dynamic settings: File Name/Input Stream: "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/cars.csv"

View: Row Separator: "\n", Field Separator: ","

Documentation: ☐ CSV options

Header: 1, Footer: 0, Limit:

Schema: Built-in, Edit schema

☒ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

**Cars**

Column	Key	Type	✓	N...	D
ID_Owner	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
Registration	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
Make	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
Color	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
ID_Reseller	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		

**Owner**

Column	Key	Type	✓	N...	D
ID_Owner	<input checked="" type="checkbox"/>	String	<input checked="" type="checkbox"/>		
Name	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
ID_Insurance	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		
Children_Nr	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		

**Cars\_data**

Column: ID\_Owner, Registration, Make, Color, ID\_Reseller

**Owners\_data**

Property: Value

Lookup Model: Load once

Match Model: Unique match

Join Model: Inner Join

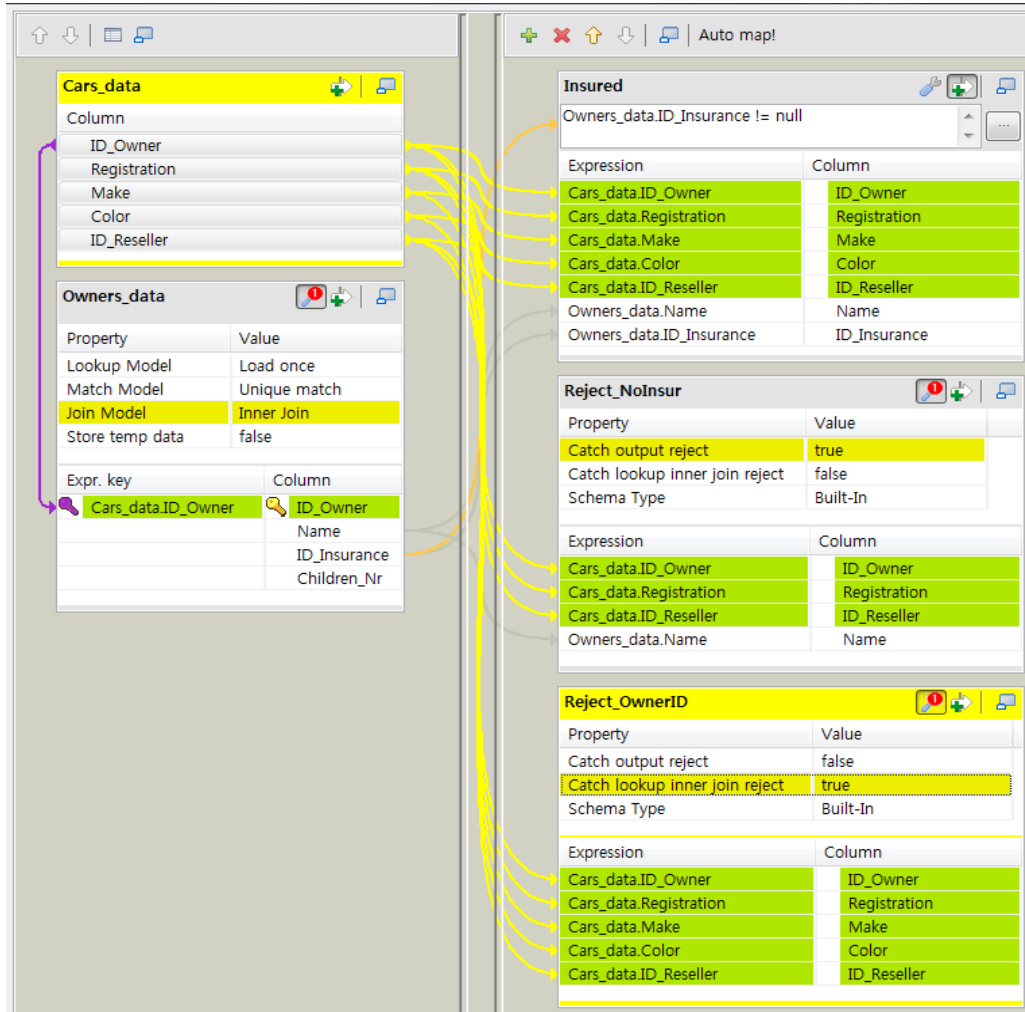
Store temp data: false

Expr. key: Cars\_data.ID\_Owner, Column: ID\_Owner, Name, ID\_Insuran..., Children\_Nr

1. 팔레트의 File/Input폴더에서 tFileInputDelimited 컴포넌트 2개, Processing폴더에서 tMap 컴포넌트, File/Output 폴더에서 tFileOutputDelimited 컴포넌트 3개를 가져와서 design workspace에 drop한다.
2. tFileInputDelimited\_1은 Cars, tFileInputDelimited\_2는 Owner로 이름을 변경한다.
3. Cars에서 tMap 로 Main link를 이용하여 연결하고 이름을 Cars\_data, Owner에서 tMap으로 Lookup link를 이용하여 연결하고 이름은 Owners\_data로 변경한다.
4. Cars 컴포넌트의 Basic settings View의 File Name에서 입력파일로 사용할 cars.csv를 선택하고, Owner도 같은 방식으로 File Name에 owner.csv를 선택한다.
5. Cars와 Owner의 Edit schema 창을 열어서 그림과 같이 컬럼을 추가한다.
6. tMap컴포넌트의 맵 에디터 창을 열어서 Cars\_data 테이블의 ID\_Owner 컬럼을 Owner\_data 테이블에 drop하여 두 테이블간의 조인을 생성한다.
7. Owner\_data 테이블의 tMap settings 버튼을 클릭한 후 Join Model은 Inner Join을 선택한다.

## 5. 컴포넌트 사용 실습

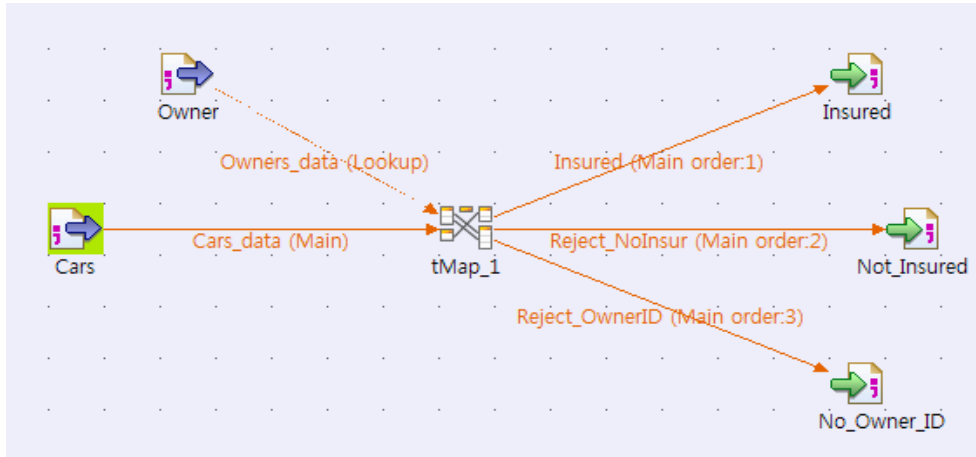
### 5.9 데이터 매핑과 Reject하기



1. Map Editor의 output 영역에 + 키를 눌러서 Insured, Reject\_NoInsur, Reject\_OwnerID 라는 3개의 output 테이블을 생성한다.
2. Cars\_data의 모든 컬럼을 Insured 테이블에 drag&drop한다.
3. Insured 테이블의 우측 상단의 +키를 클릭하여 필터 row를 추가한 후 Owners\_data 테이블의 ID\_Insurance 컬럼을 필터 조건에 drop하고 "!= null"을 입력한다. insurance ID가 있는 모든 레코드는 이 테이블에 모여지게 된다.
4. Cars\_data테이블의 ID\_Owner, Registration, ID\_Reseller 컬럼과 Owners\_data 테이블의 Name, ID\_Insurance 컬럼을 Reject\_NoInsur 테이블에 drag&drop한다.
5. Cars\_data의 모든 컬럼을 Reject\_OwnerID 테이블에 drag&drop한다.
6. Reject\_NoInsur 테이블의 우측 상단의 +키를 클릭하여 Catch output reject의 value를 true로 변경한다. insurance ID가 없는 레코드는 이 테이블에 모여지게 된다
7. Reject\_OwnerID 테이블의 우측 상단의 +키를 클릭하여 Catch lookup inner join reject의 value를 true로 변경한다. Cars\_data에서 owner ID가 매치되지 않거나 빠지는 레코드들이 모여지게 된다.

## 5. 컴포넌트 사용 실습

### 5.9 데이터 매핑과 Reject하기



**Insured(tFileOutputDelimited\_1)**

**Basic settings**

Property Type: Built-In

☐ Use Output Stream

File Name: "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/Insured\_all.csv"

Row Separator: "\n" Field Separator: ","

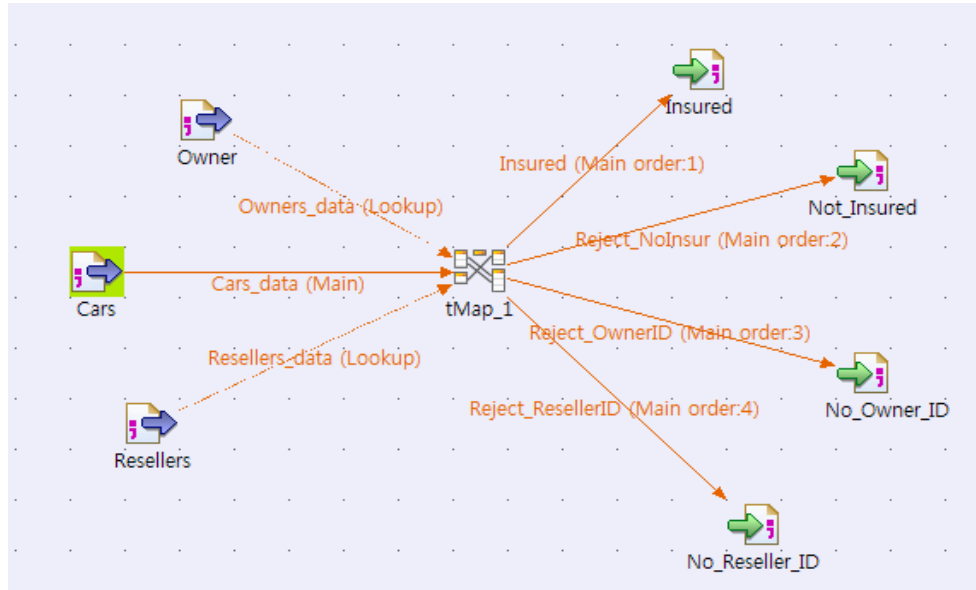
☐ Append ☒ Include Header ☐ Compress as zip file

Schema: Built-In Edit schema Sync columns

1. 앞에서 가져온 3개의 tFileOutputDelimited 컴포넌트 이름을 Insured, Not\_Insured, No\_Owner\_ID로 변경한다.
2. tMap\_1에서 오른쪽 클릭하여 Row/Insured는 Insured output 컴포넌트를 클릭하면 Insured link가 연결된다. 같은 방식으로 Reject\_NoInsur, Reject\_OwnerID를 연결한다.
3. Insured 컴포넌트의 Basic settings View의 File Name에서 출력파일로 사용할 insured\_all.csv를 입력하고 Include Header를 체크한다. 같은 방식으로 Not\_Insured, No\_Owner\_ID를 설정한다.

## 5. 컴포넌트 사용 실습

### 5.10 Inner join rejection



Property Type: Built-In

"When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."

File Name/Input Stream: "C:/talend/TOS\_BD-r111943-V5.4.1/workspace/Resellers.csv" \*

Row Separator: "#n" Field Separator: "," \*

☐ CSV options

Header: 1 Footer: 0 Limit:

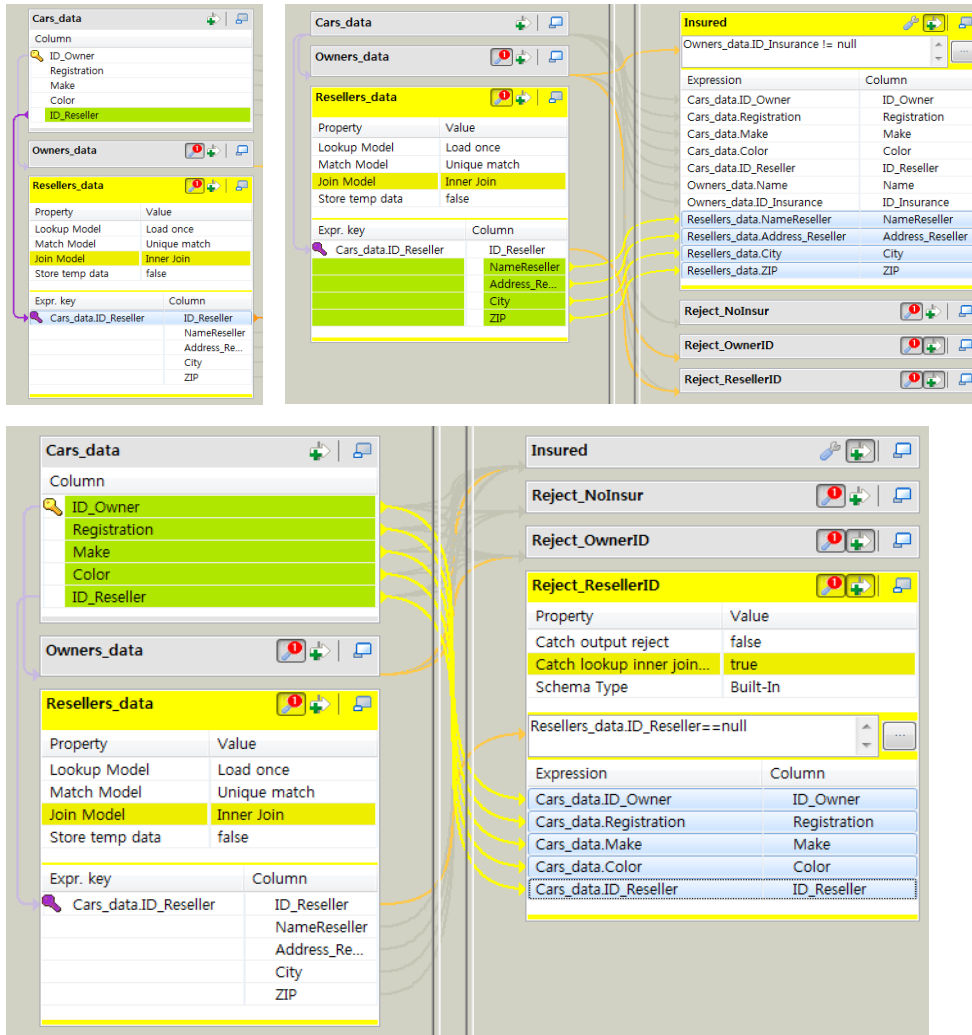
Schema: Built-In Edit schema

☐ Skip empty rows ☐ Uncompress as zip file ☐ Die on error

1. 앞의 시나리오에 이어서 작업한다.
2. 팔레트의 File/Input폴더에서 tFileInputDelimited 컴포넌트, File/Output 폴더에서 tFileOutputDelimited 컴포넌트를 가져와서 design workspace에 drop한다.
3. tFileInputDelimited 이름은 Resellers, tFileOutputDelimited 이름은 No\_Reseller\_ID로 변경한다.
4. Resellers에서 tMap 로 Main link를 이용하여 연결하고 이름을 Resellers\_data로 변경하고, tMap\_1에서 오른쪽 클릭하여 Row/No\_Reseller\_ID를 클릭한 후 No\_Reseller\_ID 컴포넌트를 클릭하면 output 이름을 입력하는 창이 표시되고 Reject\_ResellerID라고 입력한 후 OK를 누르면 Reject\_ResellerID link가 연결된다.
5. Resellers 컴포넌트의 Basic settings View의 File Name에서 출력파일로 사용할 resellers.csv를 입력하고 Include Header를 체크하고 Header는 1을 입력한다.

## 5. 컴포넌트 사용 실습

### 5.10 Inner join rejection



1. Cars\_data 와 Resellers\_data 를 조인하기 위하여 Cars\_data 테이블의 ID\_Reseller 컬럼을 Resellers\_data 테이블에 드롭하면 Resellers\_data 테이블에 ID\_Reseller 컬럼이 추가된다
2. 화면 하단의 Schema Editor를 이용하여 Name\_Reseller, Address\_Reseller, City, ZIP 컬럼을 추가한다.
3. Resellers\_data 테이블의 우측 상단의 +키를 클릭하여 Join Model 속성을 Inner Join으로 변경한다
4. Resellers\_data 테이블의 모든 컬럼을 메인 output 테이블인 Insured 에 드롭한다.
5. Cars\_data의 모든 컬럼을 선택하여 Reject\_ResellerID테이블에 드롭하면 컬럼이 추가된다.
6. Reject\_ResellerID 테이블의 우측 상단의 +키를 클릭하여 Catch lookup inner join reject 속성을 true로 변경한다

## 5. 컴포넌트 사용 실습

### 5.10 Inner join rejection

**Owners\_data**

Property	Value
Lookup Model	Load once
Match Model	Unique match
Join Model	Inner Join
Store temp data	false

**Expr. key**

Expr. key	Column
Cars_data.ID_Owner	ID_Owner
	Name
	ID_Insurance
	Children_Nr

**Resellers\_data**

Property	Value
Lookup Model	Load once
Match Model	Unique match
Join Model	Inner Join
Store temp data	false

**Expr. key**

Expr. key	Column
Cars_data.ID_Reseller	ID_Reseller

**Reject\_OwnerID**

Property	Value
Catch output reject	false
Catch lookup inner join...	true
Schema Type	Built-In

Expression: Owners\_data.ID\_Owner==null

**Reject\_ResellerID**

Property	Value
Catch output reject	false
Catch lookup inner join...	true
Schema Type	Built-In

Expression: Resellers\_data.ID\_Reseller==null

1. 첫번째 Inner Join output 테이블인 Reject\_OwnerID 테이블의 +화살표 버튼을 눌러서 필터를 추가할 수 있는 라인을 추가하고, reject할 조건 Owners\_data.ID\_Owner==null 을 입력한다.
2. 두번째 Inner Join output 테이블인 Reject\_ResellerID 의 필터 조건은 Resellers\_data.ID\_Reseller==null 을 입력한다.
3. No\_Reseller\_ID 컴포넌트의 Basic settings View의 File Name에서 출력파일로 사용할 No\_Reseller\_ID.csv를 입력하고 Include Header를 체크한다.
4. 잡을 저장한 후 실행을 하면 출력 파일 지정한 폴더에 4개의 출력파일이 생성된다.

**No\_Reseller\_ID(tFileOutputDelimited\_4)**

**Basic settings**

Property Type: Built-In

**Advanced settings**

Use Output Stream: ☒

File Name: C:/talend/TOS\_BD-r111943-V5.4.1/workspace/No\_Reseller\_ID.csv

Row Separator: \n

Field Separator: ;

Append: ☐ Include Header: ☒ Compress as zip file: ☐

Schema: Built-In