

Департамент образования и науки города Москвы

**Государственное автономное образовательное учреждение высшего
образования города Москвы**

«Московский городской педагогический университет»

Институт цифрового образования

Департамент информатики, управления и технологий

Башкатова Анна Денисовна

ОТЧЕТ

по дисциплине «Инструменты для хранения и обработки больших данных»

Тема: «02 Работа в ETL-системе Talend»

Направление подготовки (специальность) 38.03.05 – бизнес-информатика

Направленность (профиль) образовательной программы «Аналитика данных
и эффективное управление»

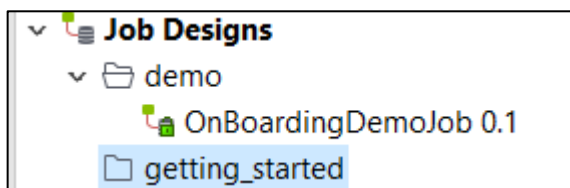
Курс обучения: 3

Форма обучения: очная

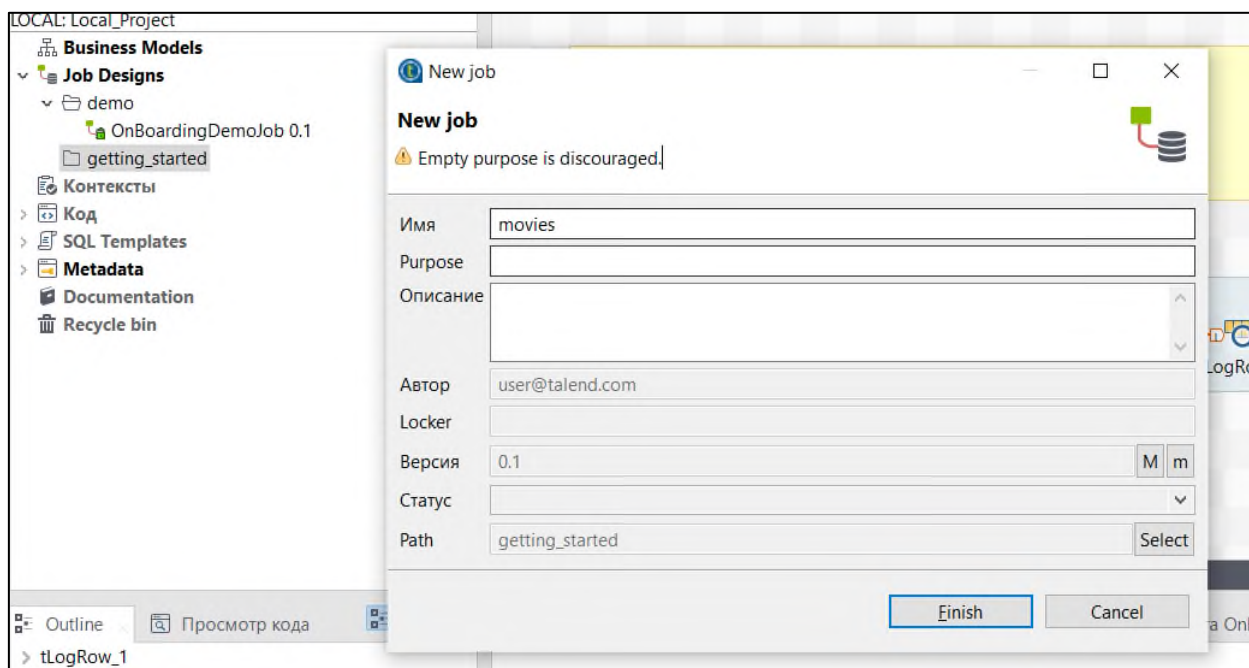
Руководитель: Босенко Т. М.

Москва
2023

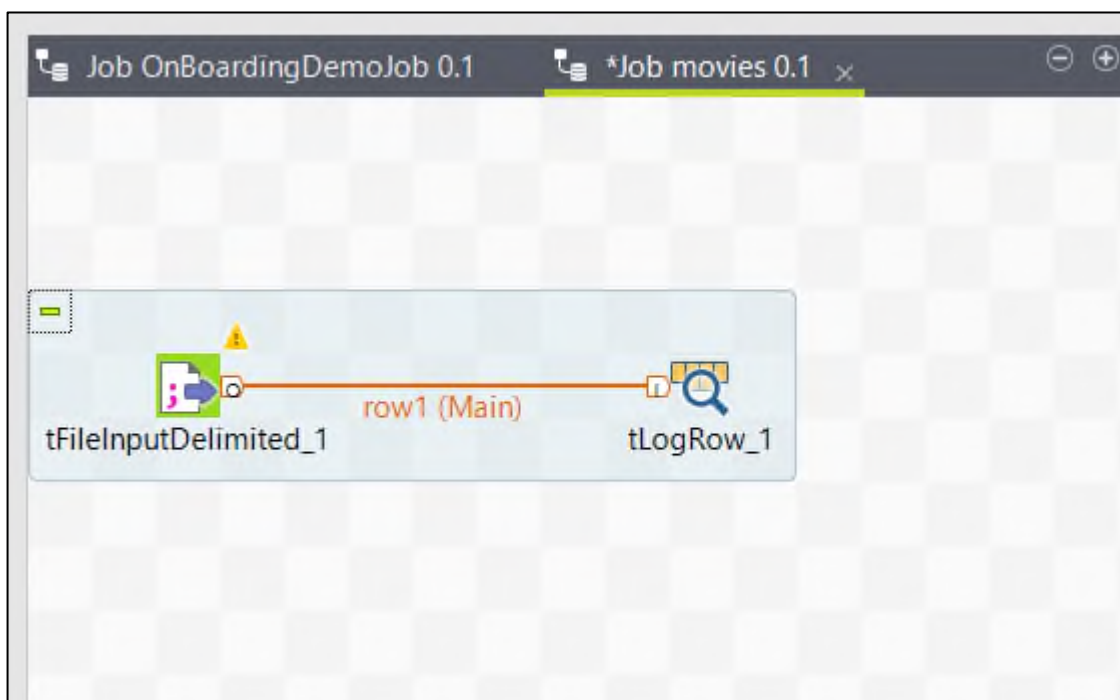
1. Создание проекта.
2. Создали новый проект; создали папку:



3. Создали работу:



4. Добавили элементы и соединили их:



5. Создали файл с разделителями:

Новый файл с разделителями

File - Step 1 of 4

Add a Metadata File on repository
Define the properties

Имя: movies

Purpose: Centralize metadata of movies.csv

Описание: Metada of file movies.csv

Автор: user@talend.com

Locker:

Версия: 0.1

Статус:

Path: Select

< Back Next > Finish Cancel

Новый файл с разделителями

File - Step 2 of 4

Add a Metadata File on repository
Define the path of the file and the format settings

Настройки файла

Сервер: Localhost 127.0.0.1

Файл: C:/Users/Bashk/Downloads/getting_started/movies.csv Обзор...

Формат: WINDOWS

File Viewer

movieID;title;releaseYear;url;directorID
315;Apt Pupil;1998;http://us.imdb.com/Title?Apt+Pupil+(1998);26
1294;Ayn Rand: A Sense of Life;1998;http://us.imdb.com/Title?Ayn+Rand%3A+A+Sense+of+Life+(1997);123
1679;B. Monkey;1998;http://us.imdb.com/M/title-exact?B%2E+Monkey+(1998);124
1649;Big One, The;1998;http://us.imdb.com/Title?Big+One,+The+(1997);122
362;Blues Brothers 2000;1998;http://us.imdb.com/M/title-exact?Blues+Brothers+2000+(1998);86
1645;Butcher Boy, The;1998;http://us.imdb.com/M/title-exact?imdb-title-118804;134
1650;Butcher Boy, The;1998;http://us.imdb.com/M/title-exact?imdb-title-118804;134
1234;Chairman of the Board;1998;http://us.imdb.com/Title?Chairman+of+the+Board+(1998);6
1654;Chairman of the Board;1998;http://us.imdb.com/Title?Chairman+of+the+Board+(1998);6
918;City of Angels;1998;http://us.imdb.com/Title?City+of+Angels+(1998);22
909;Dangerous Beauty;1998;http://us.imdb.com/M/title-exact?imdb-title-118892;113

< Back Next > Finish Cancel

Новый файл с разделителями

File - Step 3 of 4

Add a Metadata File on repository
Define the setting of the parse job

Настройки файла

КодировкаUS-ASCII

Разделитель полейCorresponding Character";"

Разделитель строкCorresponding Character"\n"

Escape Char Settings

CSV

Разделенный

Escape CharПустой

Text EnclosureПустой

☐ Split row before field

Rows To Skip

Определите следующие параметры, если какие-либо строки должны быть пропущены

Header☒ 1

Footer☐

☐ Skip empty row

Ограничение строк

Если количество строк должно быть ограничено, определите это количество

Ограничение

Предпросмотр

Выход

☒ Установить строку заголовка в качестве имен столбцов

Refresh Preview

Колонка 0	Колонка 1	Колонка 2	Колонка 3
movieID	title	releaseYear	url
315	Apt Pupil	1998	http://us.imdb.com/Title?Apt+Pupil+(1998)
1294	Ayn Rand: A Sense of Life	1998	http://us.imdb.com/Title?Ayn+Rand%3A+A+Sense+of+Life+(1997)
1679	B. Monkey	1998	http://us.imdb.com/M/title-exact?B%2E+Monkey+(1998)

Экспортировать как контекст

Восстановить контекст

< Back

Next >

Finish

Cancel

Новый файл с разделителями

File - Step 4 of 4

Add a Schema on repository
Define the Schema

Имяmetadata

Комментарий

Схема

Click to update schema preview

Guess

Описание схемы

Колонка	К...	Тип	<input checked="" type="checkbox"/> N..	Date Pattern (Ctrl...	Длина	Precision	Default	Коммент...
movieID	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		4	0		
title	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		29	0		
releaseYear	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		4	0		
url	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		66	0		
directorID	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		3	0		

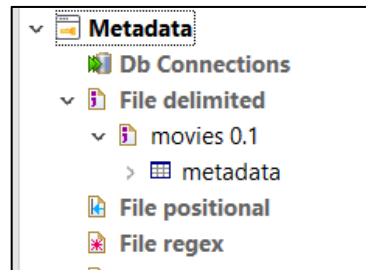
< Back

Next >

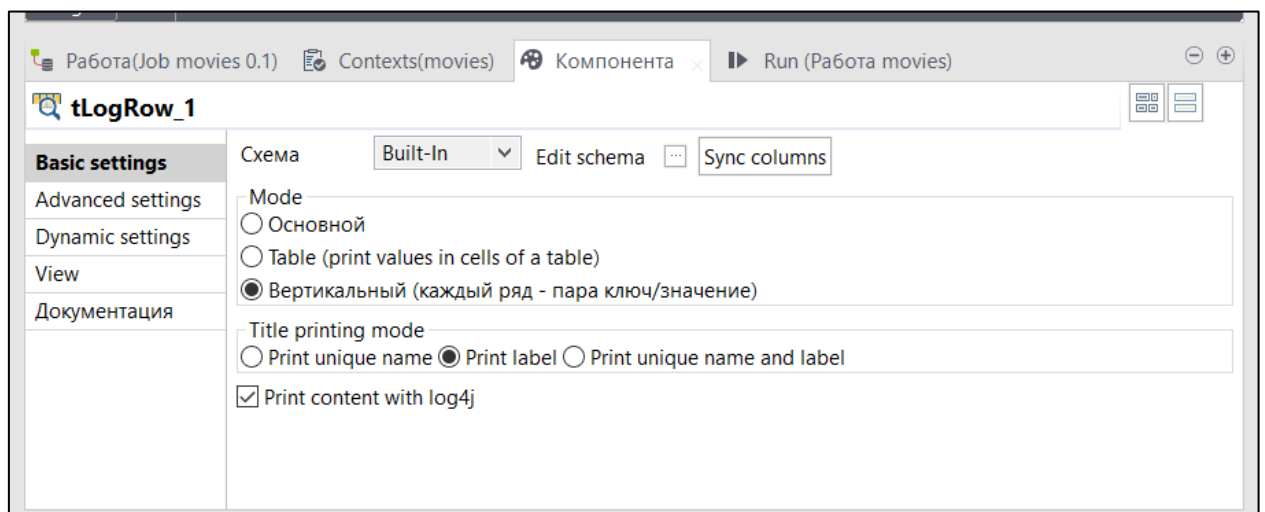
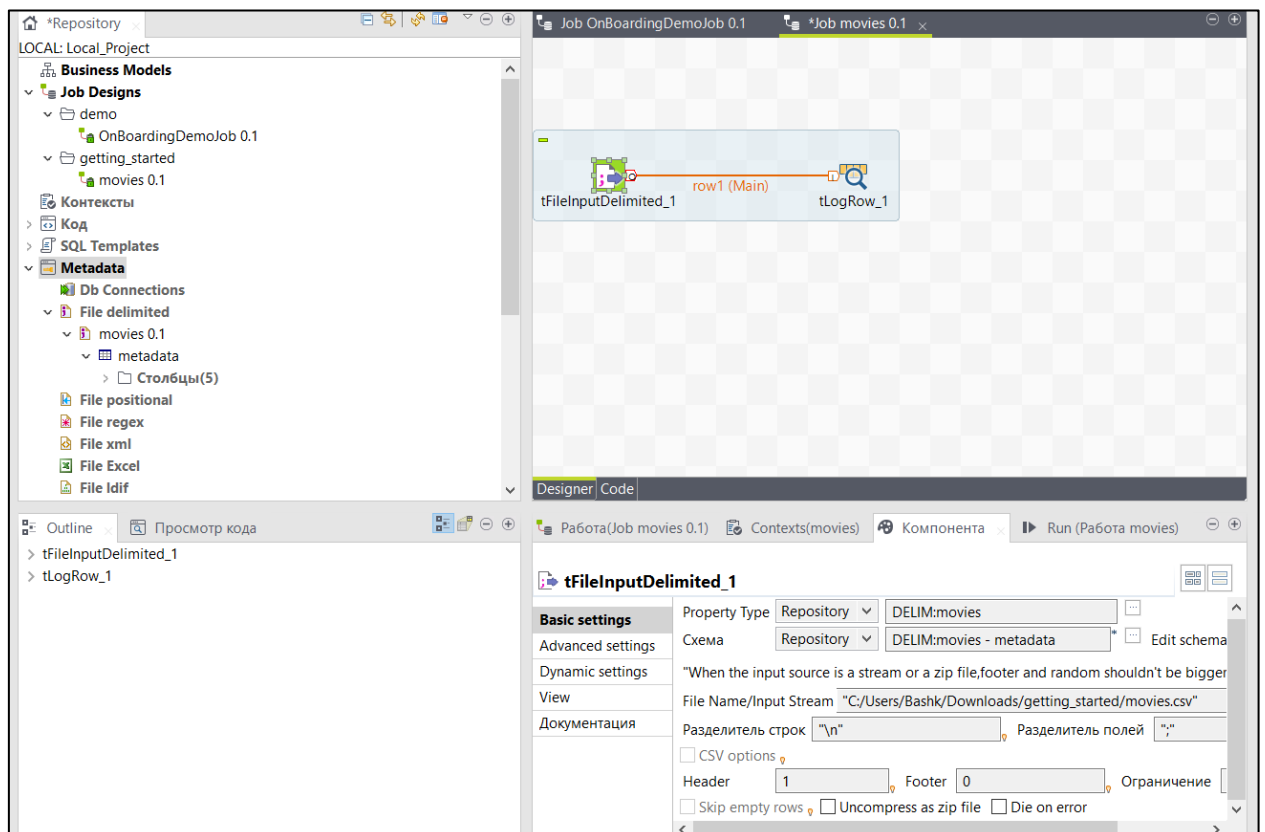
Finish

Cancel

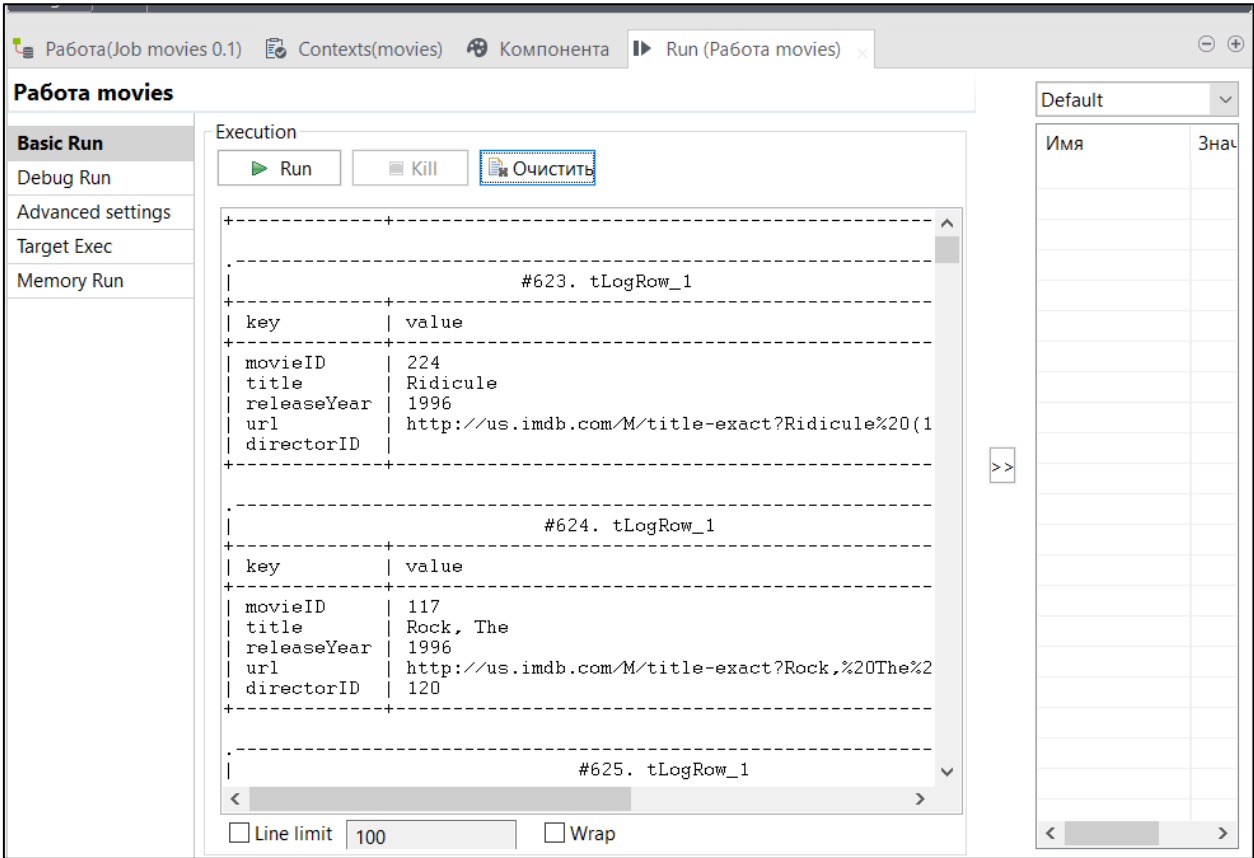
6. Схема с загруженными метаданными:



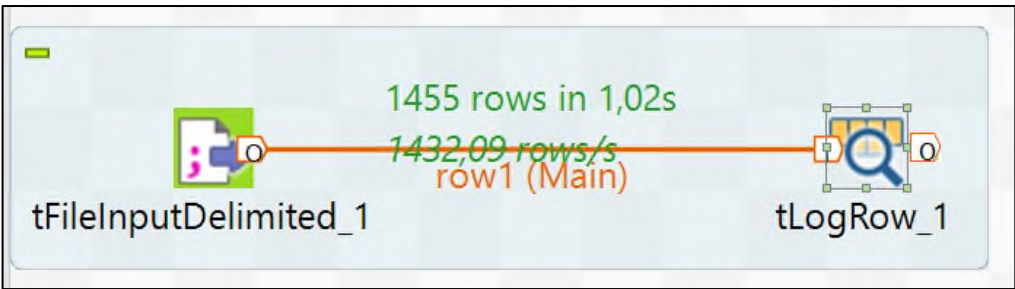
7. Настройка и выполнение работы:



8. Консоль Запуска отображает информацию о фильмах, считанную из исходного файла:



9. Фильтрация информации о фильмах:



10. Подготовка метаданных:

Новый файл с разделителями

File - Step 1 of 4

Add a Metadata File on repository
Define the properties

Имя: directors

Purpose: Centralize the metadada of direcrors info

Описание: Metadata of the directors dataset

Автор: user@talend.com

Locker:

Версия: 0.1 M m

Статус:

Path: Select

< Back Next > Finish Cancel

Новый файл с разделителями

File - Step 2 of 4

Add a Metadata File on repository
Define the path of the file and the format settings

Настройки файла

Сервер: Localhost 127.0.0.1

Файл: C:/Users/Bashk/Downloads/getting_started/directors.txt O6sop...

Формат: WINDOWS

File Viewer

1, Gregg Araki
2, P.J. Hogan
3, Alan Rudolph
4, Alex Proyas
5, Alex Sichel
6, Alex Zamm
7, Alfonso Cuarón
8, Alfred Hitchcock
9, Allison Anders
10, Andrew Davis
11, Andrew Niccol
12, Antoine Fuqua

< >

< Back Next > Finish Cancel

Новый файл с разделителями

File - Step 3 of 4

Add a Metadata File on repository
Define the setting of the parse job

Настройки файла

Кодировка

UTF-8

Разделитель полей

Corresponding Character

" "

Разделитель строк

Corresponding Character

"\n"

Escape Char Settings

CSV

Разделенный

Escape Char

Пустой

Text Enclosure

Пустой

Split row before field

Rows To Skip

Определите следующие параметры, если какие-либо строки должны быть пропущены

Header

Footer

Skip empty row

Ограничение строк

Если количество строк должно быть ограничено, определите это количество

Ограничение

Предпросмотр

Выход

Установить строку заголовка в качестве имен столбцов

Refresh Preview

Колонка 0	
1, Gregg Araki	
2, P.J. Hogan	
3, Alan Rudolph	
4, Alex Proyas	

Экспортировать как контекст

Восстановить контекст

< Back

Next >

Finish

Cancel

Новый файл с разделителями

File - Step 4 of 4

Add a Schema on repository
Define the Schema

Имя

metadata

Комментарий

Схема

Click to update schema preview

Guess

Описание схемы

Колонка	К...	Тип	<input checked="" type="checkbox"/> N.	Date Pattern (Ctrl...	Длина	Precision	Default	Коммент...
Column0	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0		
Column1	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		20	0		

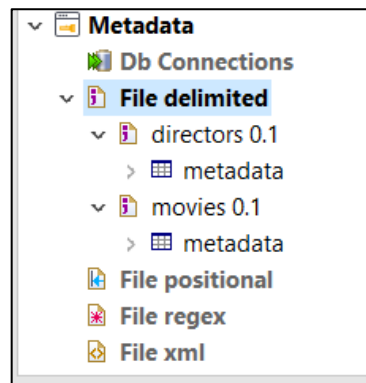
< Back

Next >

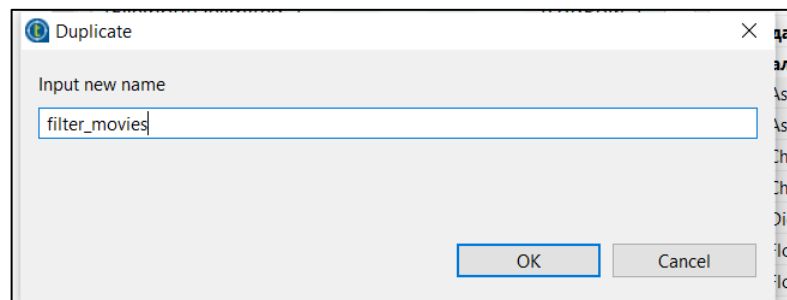
Finish

Cancel

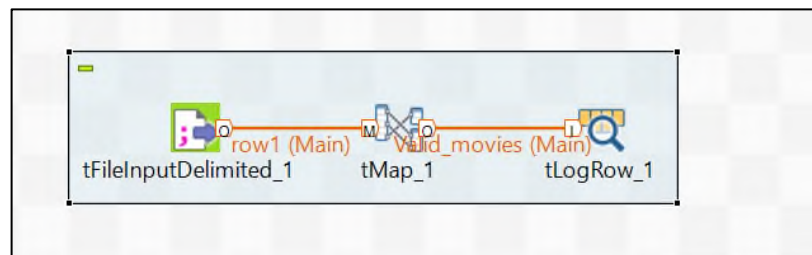
11.Представления репозитория:



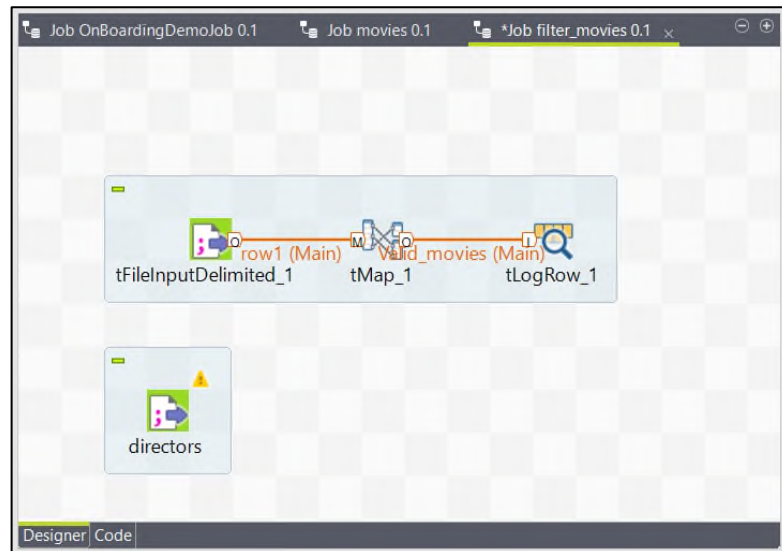
12. Дублирование задания:



13. Добавление компонента сопоставления:



14. Добавление компонента поиска:



Designer | Code

Работа(Job filter_movies 0.1) Contexts(filter_movies) Компонента x Run (Работа filter_movies)

directors(tFileInputDelimited_2)

Basic settings

Property Type: Repository DELIM:directors

Схема: Repository DELIM:directors - metadata Edit schema

Dynamic settings

View

Документация

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

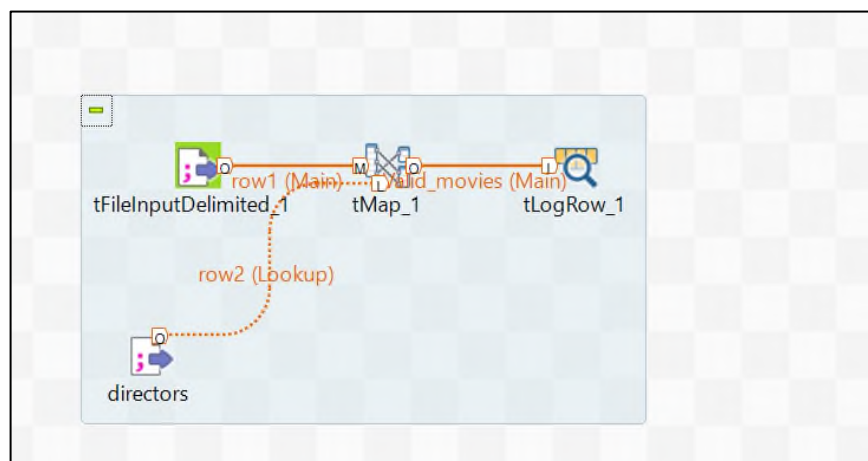
File Name/Input Stream: "C:/Users/Bashk/Downloads/getting_started/directors.txt"

Разделитель строк: "\n" Разделитель полей: ","

☐ CSV options

Header: 0 Footer: 0 Ограничение:

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



Designer | Code

Работа(Job filter_movies 0.1) Contexts(filter_movies) Компонента x Run (Работа filter_movies)

directors(tFileInputDelimited_2)

Basic settings

☐ Advanced separator(for number)

☐ Extract lines at random Кодировка: UTF-8

Advanced settings

☒ Trim all column ☐ Проверить каждый ряд на соответствие схеме ☐ Check date

Dynamic settings

View

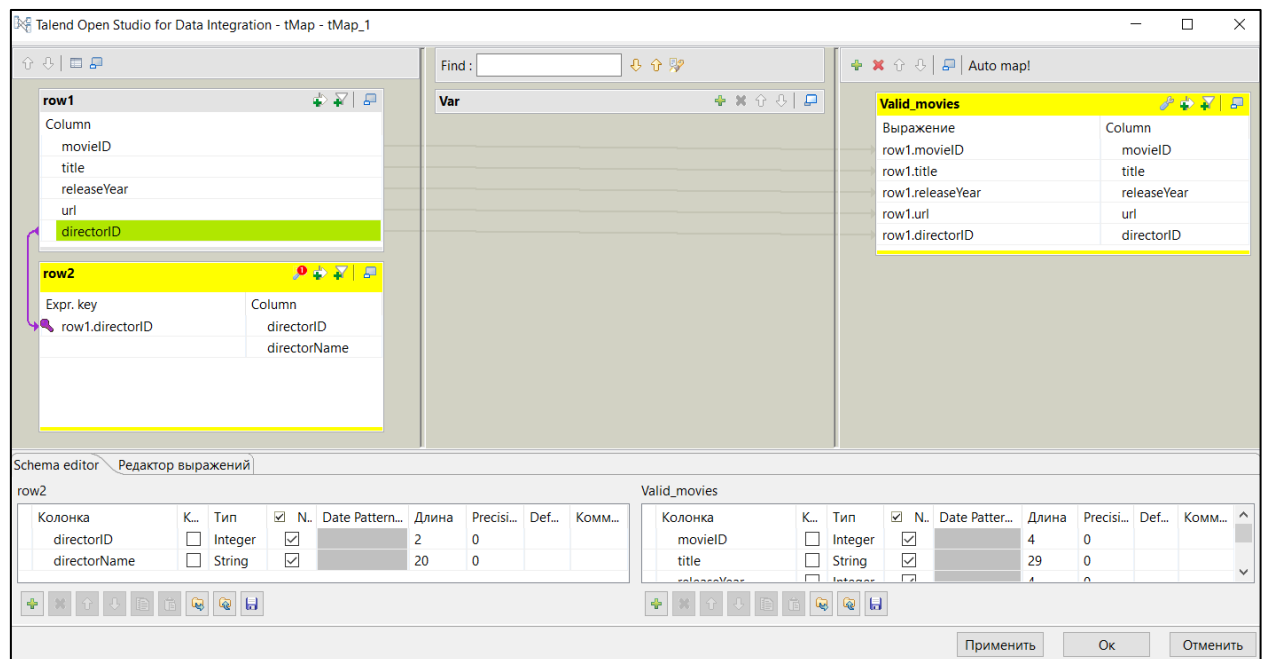
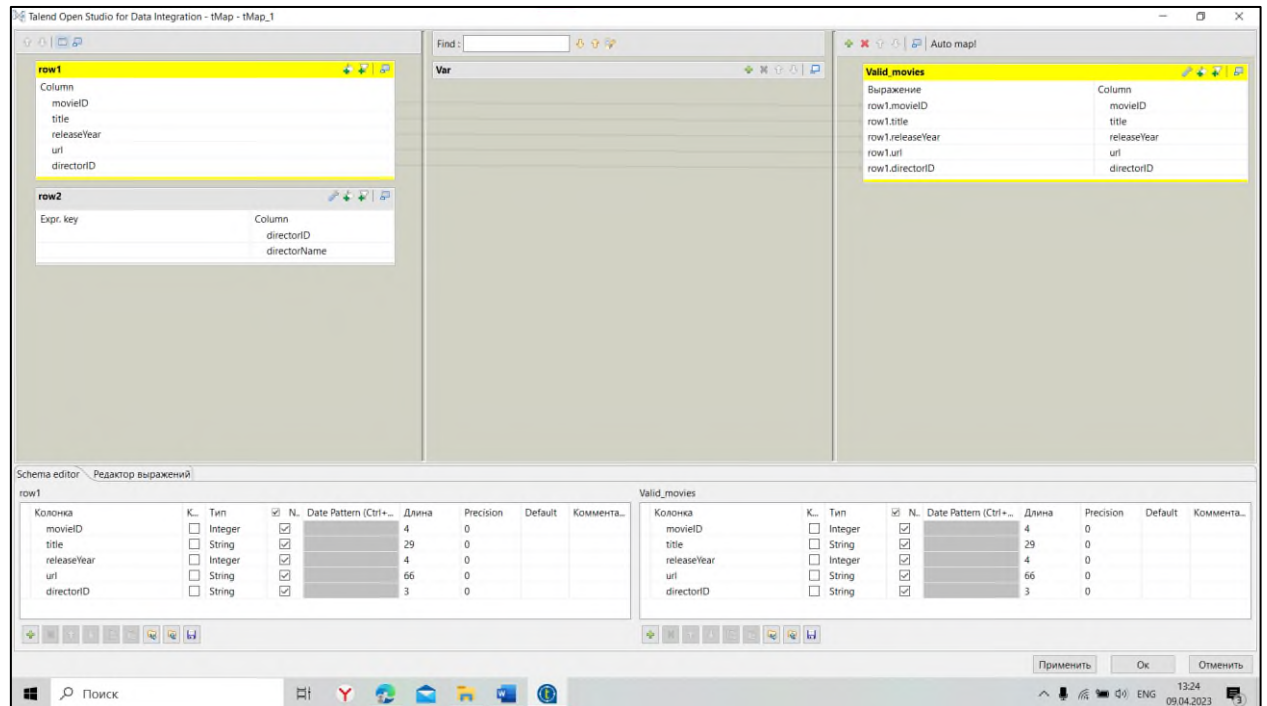
Документация

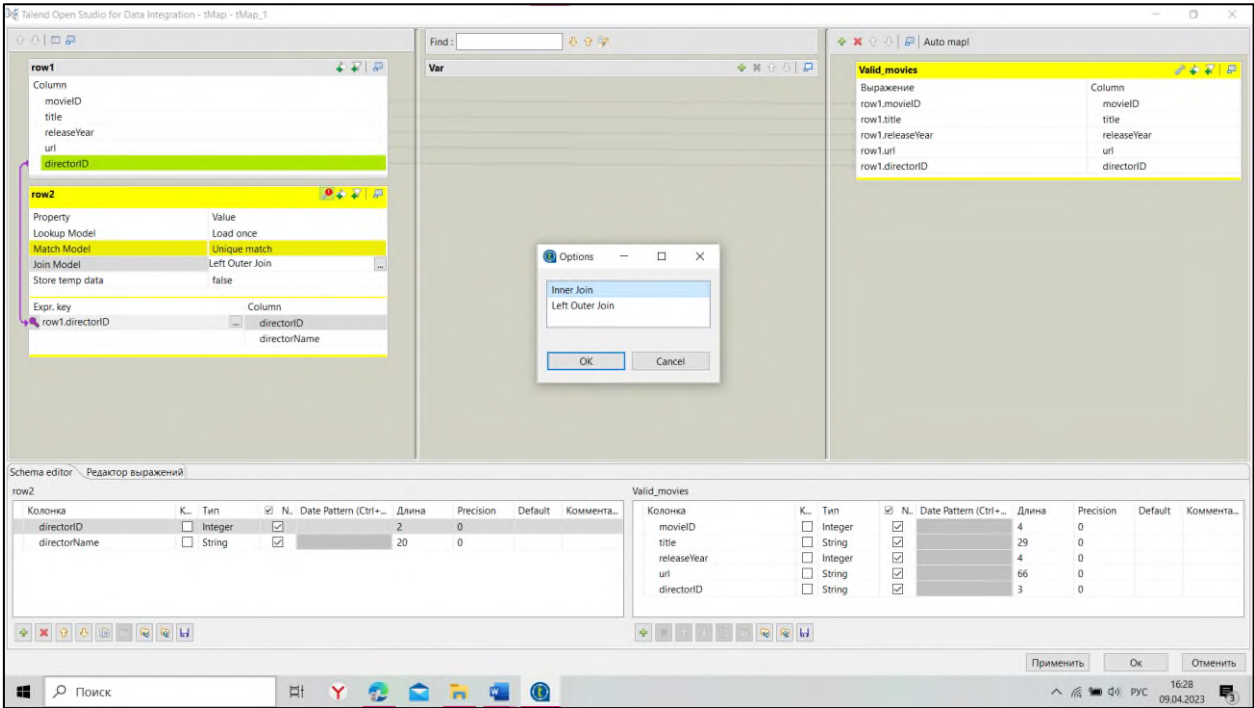
☐ Split row before field

☐ Permit hexadecimal (0xNNN) or octal (0NNNN) for numeric types - it will act the opposite for Byte

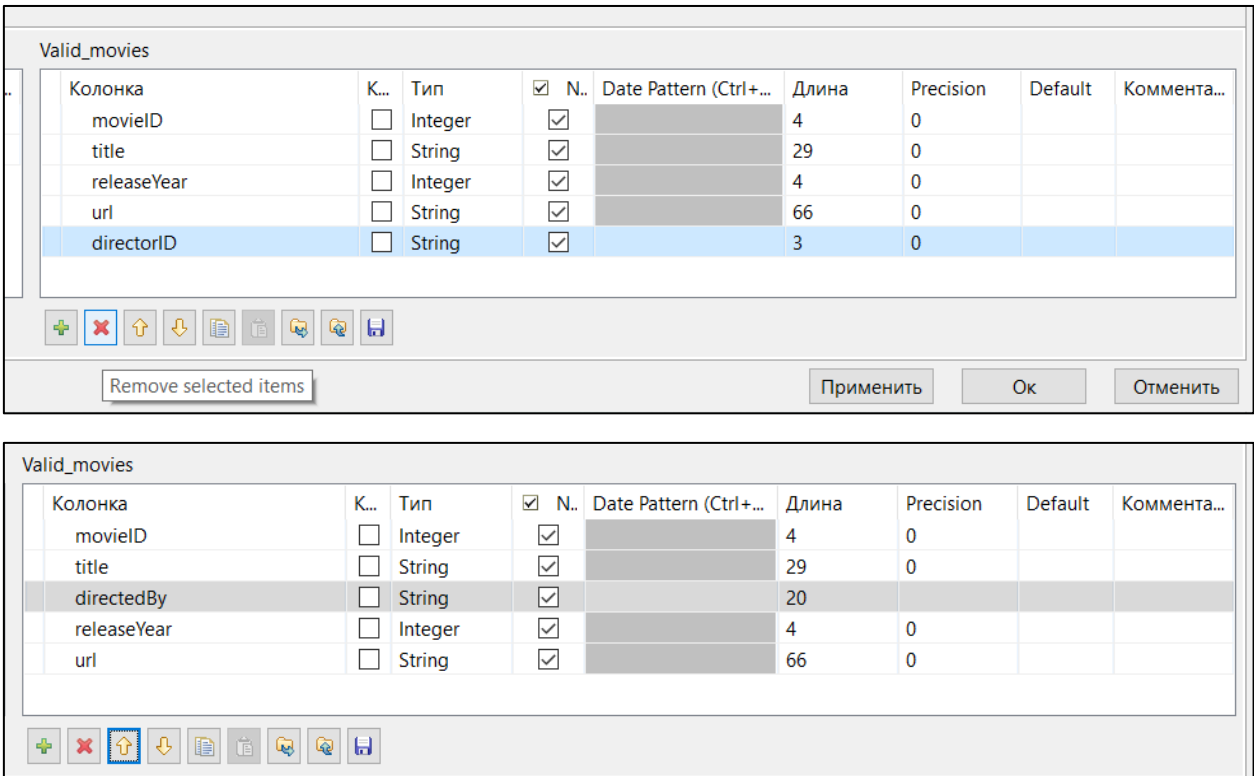
☐ tStatCatcher Statistics

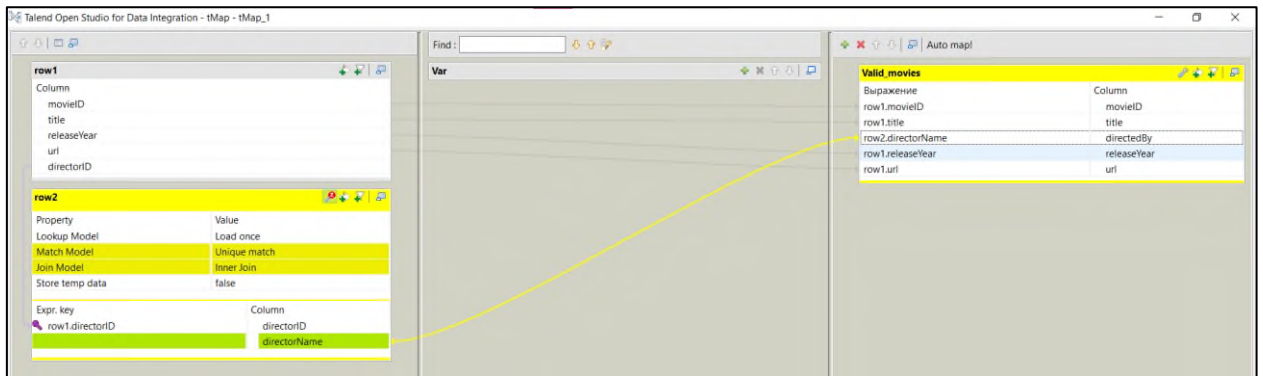
15. Настройка сопоставлений и выполнение:



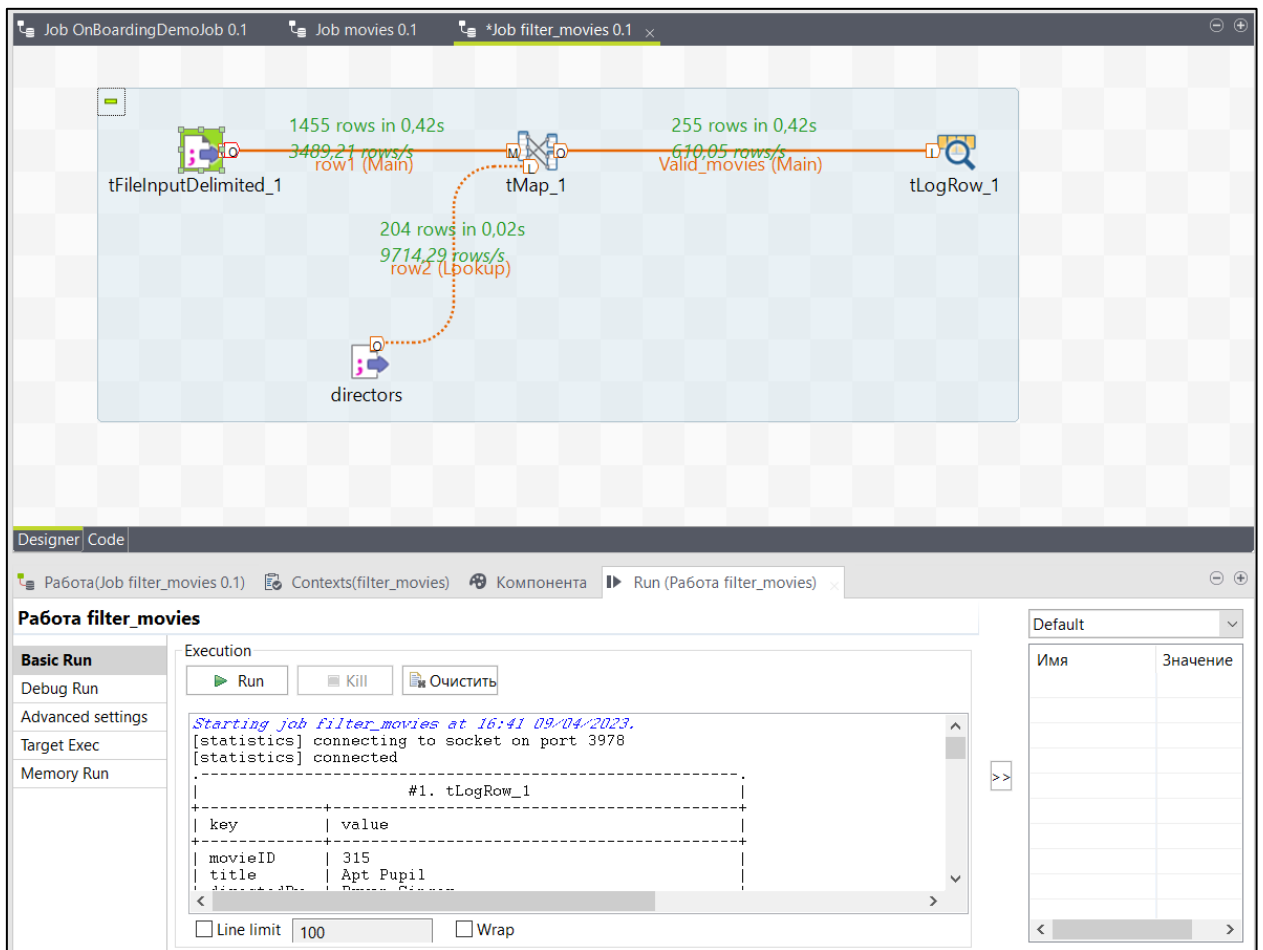


16.Удаление и добавление строк

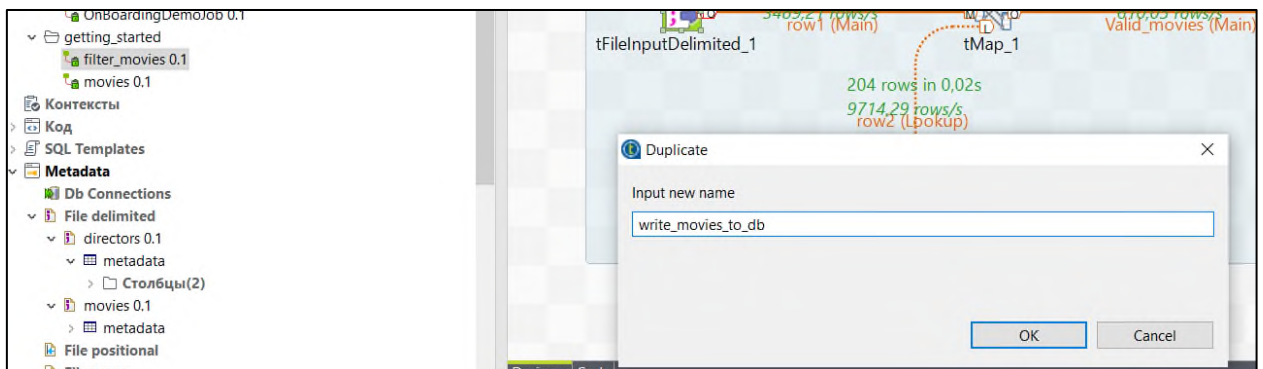


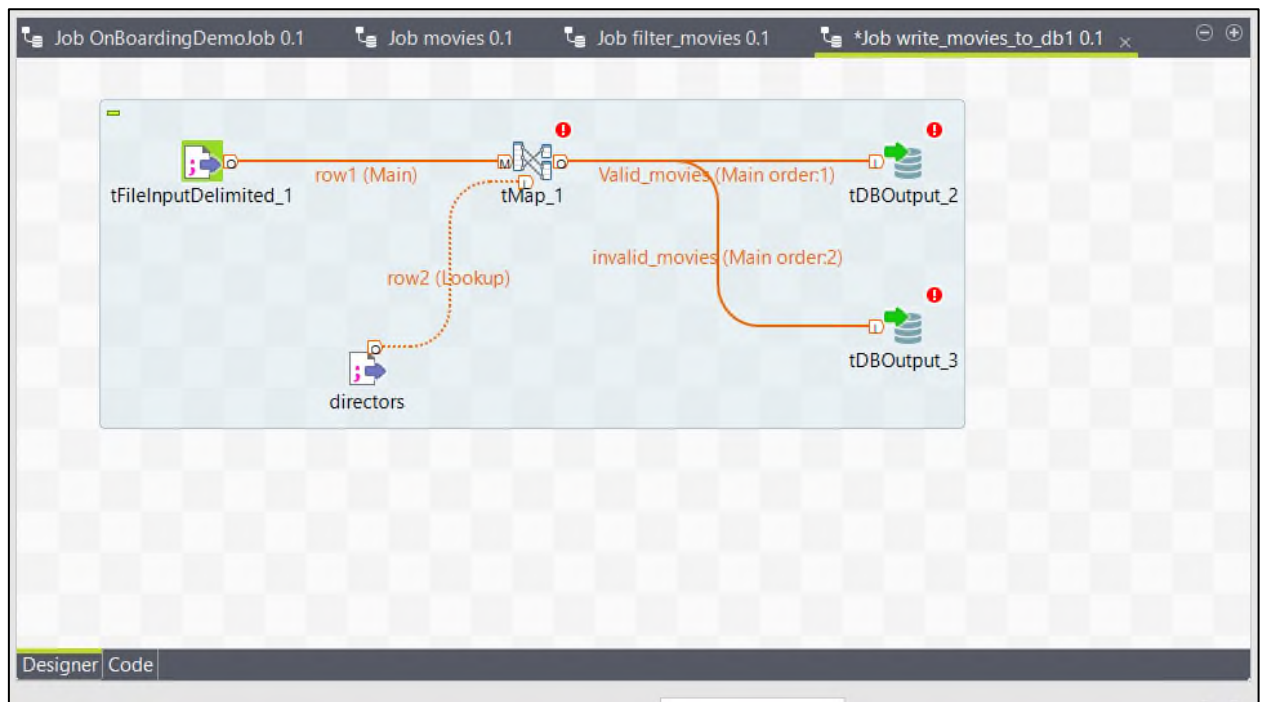


17.Выполнение:



18.Сбор информации об отклоненных фильмах и сохранение результатов обработки в базе данных





The screenshot shows the Talend Studio interface for configuring the **tMap_1** component. The **Schema editor** is open, displaying the **Invalid_movies** output configuration. The **Invalid_movies** output is configured with the following properties:

Property	Value
Catch output reject	false
Catch lookup inner join reject	true
Schema Type	Column
row1.movieID	movieID
row1.title	title

The **Schema editor** also displays the **Invalid_movies** output configuration, showing the **Invalid_movies** output is configured with the following properties:

Property	Value
Catch output reject	false
Catch lookup inner join reject	true
Schema Type	Column
row1.movieID	movieID
row1.title	title

The screenshot shows the **Options** dialog box for the **Invalid_movies** output. The dialog box is titled **Options** and contains the following configuration:

Property	Value
row1.movieID	true
row1.title	false

The **Options** dialog box is also showing the **Invalid_movies** output configuration, which includes the **Invalid_movies** output and the **Invalid_movies** output configuration.

19. Настройка выходных данных базы данных MySQL

tDBOutput_1(MySQL)

Basic settings Database: MySQL Apply

Advanced settings Property Type: Built-In

Dynamic settings Версия БД: Mysql 8

View ☐ Использовать существующее соединение

Документация Хост: "localhost"

Порт: "3306"

Database: "bashkatova"

Имя пользователя: "root"

Пароль: "*****"

Таблица: "Valid_movies"

Action on table: Drop table if exists and create

Действие над данными: Вставить

Схема: Built-In Edit schema Sync columns

tDBOutput_2(MySQL)

Basic settings Database: MySQL Apply

Advanced settings Property Type: Built-In

Dynamic settings Версия БД: Mysql 8

View ☐ Использовать существующее соединение

Документация Хост: "localhost"

Порт: "3306"

Database: "bashkatova"

Имя пользователя: "root"

Пароль: "*****"

Таблица: "invalid_movies"

Action on table: Drop table if exists and create

Действие над данными: Вставить

Схема: Built-In Edit schema Sync columns

Data source

Designer Code

Работа(Job Copy_of_movies 0.1) Contexts(Copy_of_movies) Компонента Run (Работа Copy_of_movies)

Работа Copy_of_movies

Basic Run Execution

Run Kill Очистить

Starting job Copy_of_movies at 13:13 11/04/2023.

[statistics] connecting to socket on port 4029

[statistics] connected

[statistics] disconnected

Job Copy_of_movies ended at 13:13 11/04/2023. [exit code = 0]

Line limit 100 Wrap

phpMyAdmin

Сервер: 127.0.0.1:3306 База данных: bashkatova

Структура SQL Поиск Запрос по шаблону Экспорт Импорт Операции Привилегии

Фильтры

Содержит слово:

Таблица	Действие	Строки	Тип	Сравнение	Размер	Фрагментировано
invalid_movies		1 290	InnoDB	utf8_general_ci	16.0 КиБ	-
valid_movies		255	InnoDB	utf8_general_ci	64.0 КиБ	-
2 таблицы	Всего	1 455	InnoDB	utf8_general_ci	80.0 КиБ	0 Байт

Отметить все С отмеченными:

phpMyAdmin

Надавене

Избранное

Создать БД

bashkatova

Новая

invalid_movies

valid_movies

information_schema

mysql

performance_schema

sys

Сервер: 127.0.0.1:3306

База данных: bashkatova

Таблица: invalid_movies

Обзор

Структура

SQL

Поиск

Вставить

Экспорт

Импорт

Привилегии

Операции

Триггеры

Данное выделение не содержит уникального столбца. Изменение сетки, выставление галочки, редактирование, копирование и удаление невозможно.

Отображение строк 0 - 49 (1200 всего, Запрос занял 0.0006 сек.)

SELECT * FROM `invalid_movies`

Профилирование Построение редактирование Изменить Анализ SQL запроса Создать PHP-код Обновить

1 > >>

Показать все

Количество строк: 50

Фильтровать строки: Поиск в таблице

Параметры

movieID	title
1591	Duoluo lianshi
905	Great Expectations
349	Hard Rain
1612	Leading Man, The
913	Love and Death on Long Island
1582	Magic Hour, The
1678	Mat' i syn
1313	Palmetto
1651	Spanish Prisoner, The
1613	Tokyo Fist
1300	Ta There Was You
1643	Angel Baby
1382	Bonheur, Le
936	Brassed Off
295	Breakdown
Консоль	Cats Don't Dance