

# Université de Cergy-Pontoise

## RAPPORT

Intégration et Entrepôts de Données  
Master 1 Informatique et Ingénierie des Systèmes Complexes (IISC)  
Parcours Système Intelligents et Distribués(SID)

sur le sujet

Le but de ce projet est de concevoir et réaliser un médiateur simple, exploitant différentes sources de données dans le domaine du cinéma

rédigé par

**Mohamed II Bayo, Benoît FAGOT**

Avril 2019

### Table des matières

<b>1</b>	<b>Description extraction des informations API JSoup :</b>	<b>1</b>
<b>2</b>	<b>La description des transformation sur Talend :</b>	<b>1</b>
2.1	Fusion des fichiers extrait de la page HTML avec l'API JSoup : . . . . .	1
2.2	Configuration des mappages avec tMap : . . . . .	1
2.3	Sauvegarde des résultats de traitement dans la base de données : . . . . .	2
<b>3</b>	<b>Le mapping entre le médiateur et les 3 sources</b>	<b>2</b>
3.1	Source 1 : BD locale . . . . .	2
3.2	Source 2 : DBpedia . . . . .	2
3.3	Source 3 : OpenMovieDatabase . . . . .	3
<b>4</b>	<b>La structure du programme du médiateur et son fonctionnement :</b>	<b>3</b>
<b>5</b>	<b>Annexe</b>	<b>3</b>

### Table des figures

1	Fichier genre extrait avec JSoup . . . . .	1
2	configuration tMap . . . . .	2
3	resultat final . . . . .	2
4	UML . . . . .	3
5	Exemple resultat console : fig 1 : Actor "Brad Pitt", fig 2 : t=Title "Avatar" . . . . .	4

# 1 Description extraction des informations API JSoup :

- On constate un pattern régulier dans le DOM des pages web du site : les données sont toujours rentrées dans la 2ème table de la page, il suffit donc de sélectionner les `<tr>` puis les `<td>` de la table mis à part les deux dernières qui sont des données récapitulatives.
- Pour obtenir les données des années 2000 à 2015, on effectue une boucle sur l'adresse de la page, et on récidive pour chaque genre demandé.
- On extrait les données de la table qui nous intéressent à savoir le distributeur et le titre du film.
- Pour chaque genre, on crée et écrit dans un fichier CSV contenant une colonne Genre éponyme, et on remplit les deux autres colonnes par les données extraites précédemment.
- Chaque fichier csv aura pour entete(genre, titre, distributeur)

Exemple pratique :

- on fait une combinaison de l'année et du genre pour former l'URL correspondant à l'année et le genre du film à extraire  
exemple : `http://www.the-numbers.com/market/" + i + "/genre/" + genres[j]).get()` ici i = année, et genre[j] un élément de la liste des genres. Par exemple pour l'année 2000 et genre Adventure , on aura :  
`http://www.the-numbers.com/market/2000/genre/Adventure`
- ensuite,
  - On sélectionne la balise table -> les tr Exemple : `Elements tables = doc.select("table").get(1).select("tr");`
  - On regarde ensuite pour chaque `<tr>` les `<td>` :  
for (Element headline : tables) {  
Elements tds = headline.select("td");  
}
  - On extrait les `<td>` qui nous intéressent ( `tds.get(1) && tds.get(3)` )

## 2 La description des transformation sur Talend :

### 2.1 Fusion des fichiers extraits de la page HTML avec l'API JSoup :

Après avoir extrait le fichier csv pour chaque genre avec JSoup, on les a fusionné afin de former un seul fichier csv sur talend avec le composant "tUnite".

- Le composant tUnite fusionne des données de diverses sources.
- tLogRow pour afficher le résultat de la fusion des 6 fichiers csv dans la console Run.

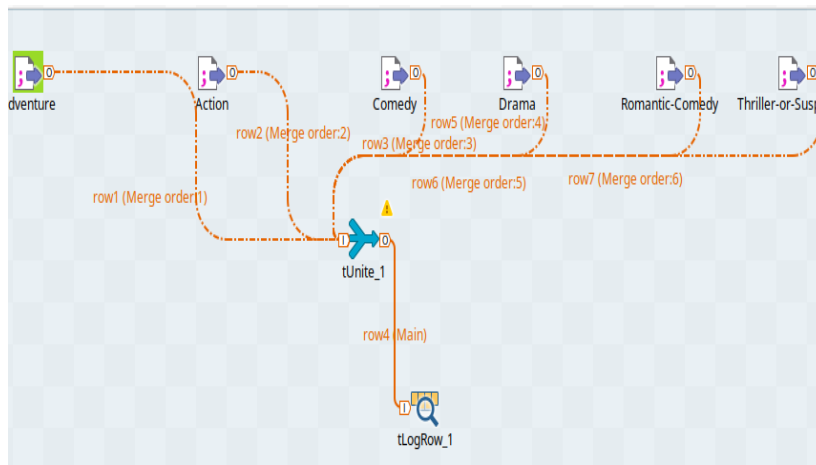


FIGURE 1 – Fichier genre extrait avec JSoup

### 2.2 Configuration des mappages avec tMap :

- L'éditeur de tmap nous affiche trois tables movieBudget, genreFichierFusionJsoup et movieOut correspondant respectivement au schéma du fichier de movieBudget, au schéma du fichier de fusion des six fichiers genre extrait du jsoup et au schéma de sortie pour les informations de film.
- On a sélectionné la colonne movie dans la table movieBudget et déposez dans la colonne titre de la table genreFichierFusionJsoup pour créer une jointure entre les deux ensembles de données en entrée en fonction des titres. Puis sur le bouton paramètre tMap de la table genreFichierFusionJsoup, on a choisi la jointure externe gauche

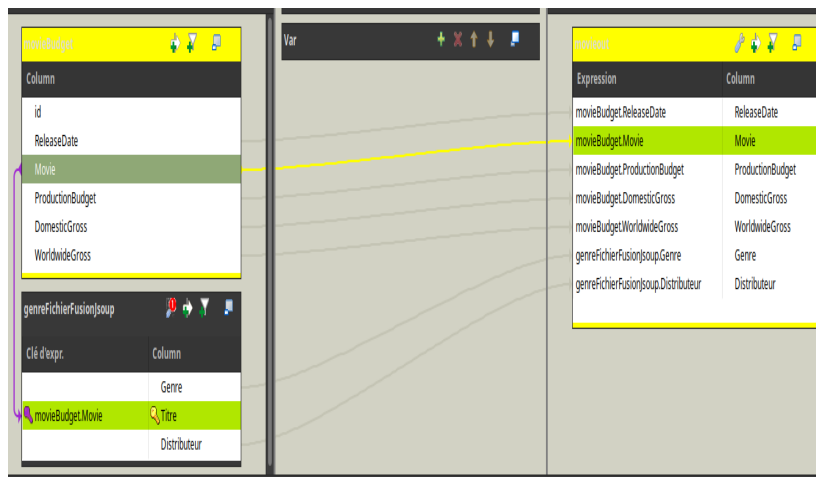


FIGURE 2 – configuration tMap

### 2.3 Sauvegarde des résultats de traitement dans la base de données :

Avec le composant tDBOutput on a créé une connexion à notre base de donnée "Mysql", puis on a cliqué sur tMap pour le lier à tDBOutput. La base de donnée recevra les données de sortie de tMap de la table "movieOut"

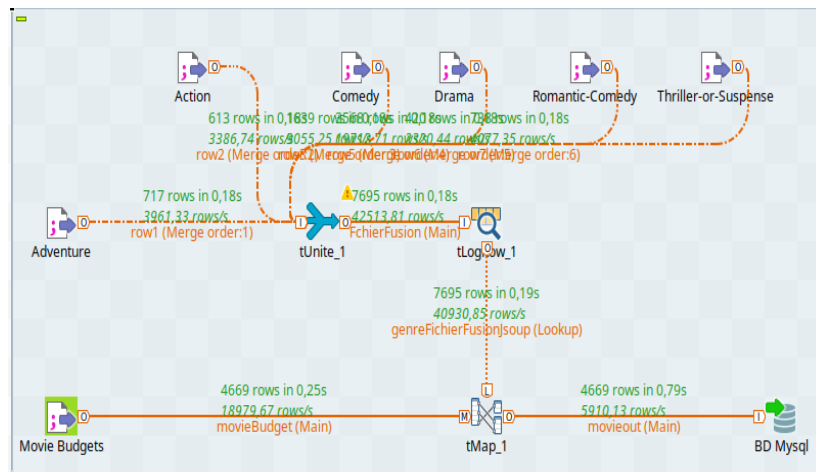


FIGURE 3 – resultat final

## 3 Le mapping entre le médiateur et les 3 sources

### 3.1 Source 1 : BD locale

Avec JDBC on fait une requête SQL classique pour obtenir l'année de sortie, les différents revenus du film, le budget, le genre, le distributeur.

```
SELECT * FROM movie WHERE Movie LIKE title
```

### 3.2 Source 2 : DBpedia

Avec l'API Jena, nous effectuons nos requêtes SPARQL puis stockons les résultats en mémoire. Nous effectuons 4 requêtes différentes pour obtenir le réalisateur d'un film, la liste des acteurs d'un film, la liste des producteurs d'un film, la liste des films où un acteur a joué. Ce sont les requêtes vu en TP. Il est primordial de renseigner l'ontologie suivante : PREFIX dbo : <http://dbpedia.org/ontology/>

```
SELECT ?nr WHERE { ?f a dbo :Film ; foaf :name "title"@en ; dbo :starring ?ac . ?ac foaf :name ?nr . }
SELECT ?nr WHERE { ?f a dbo :Film ; foaf :name "title"@en ; dbo :director ?d . ?d foaf :name ?nr . }
SELECT ?nr WHERE { ?f a dbo :Film ; foaf :name "title"@en ; dbo :producer ?d . ?d foaf :name ?nr . }
SELECT ?t WHERE { ?f a dbo :Film ; foaf :name ?t ; dbo :starring ?a . ?a foaf :name "actor"@en }
```

