

## TP N°8 : Intégration de données avec Talend

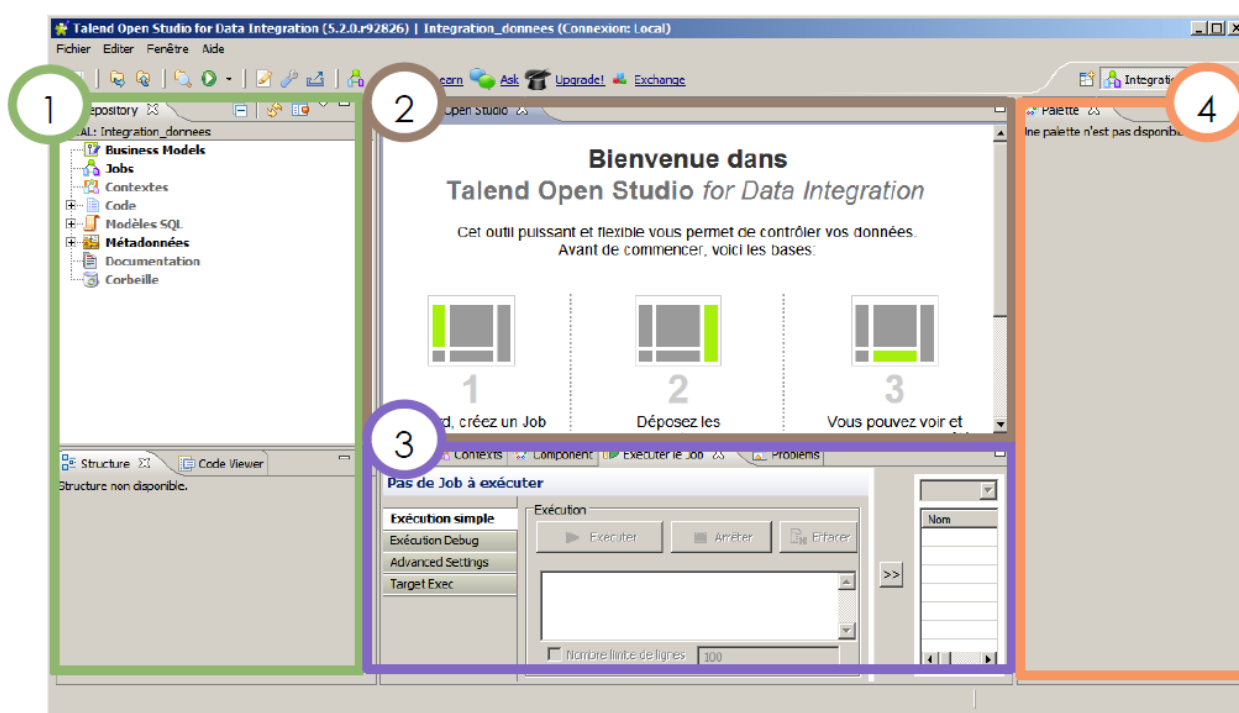
Pour les besoins de notre TP, nous utilisons « Talend Data Integration » pour la transformation des données et leur intégration. Il est possible de télécharger toutes les solutions de Talend Open Studio sur <http://fr.talend.com/products/talend-open-studio>

### 1) Installation et Démarrage

Après avoir installé Talend sur votre machine, le démarrer et créer un nouveau projet intitulé : Intégration\_1.

Remarque : Veiller à ce que votre workspace soit à un emplacement accessible en lecture et en écriture (comme vos documents ou votre bureau) : Éviter de le créer directement dans le répertoire d'installation de Talend.

Après la fermeture de la page de Bienvenue, la fenêtre qui s'affiche aura la forme suivante (selon la version) :

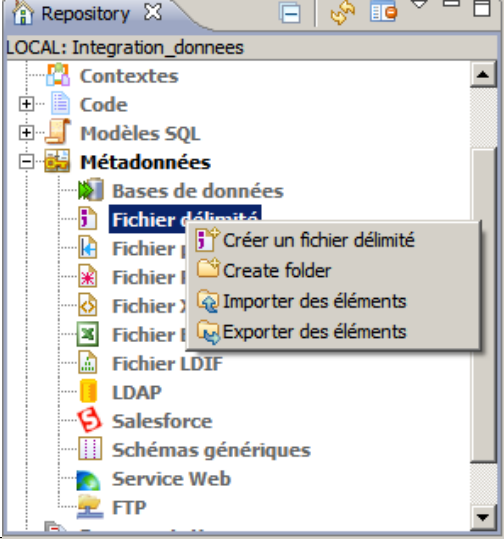
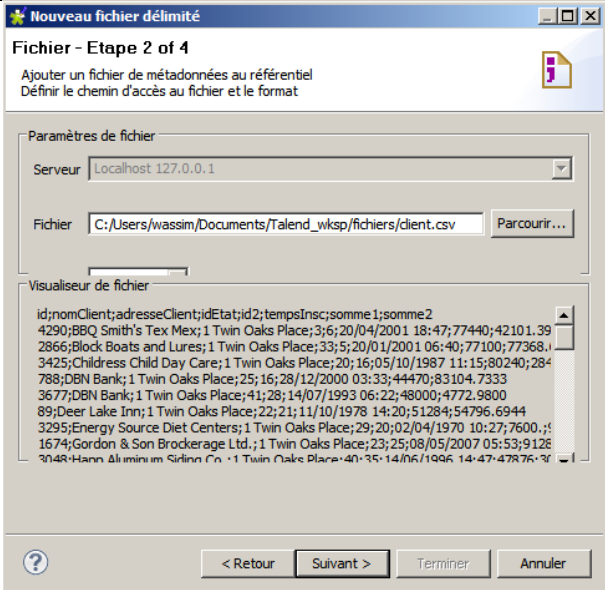


1	Panneau représentant la structure de votre projet.
2	Panneau affichant l'architecture des Jobs et le code
3	Onglets contenant les propriétés des composants, la console d'exécution, les problèmes...
4	Palette des différents composants disponibles

## 2) Préparation des sources de données

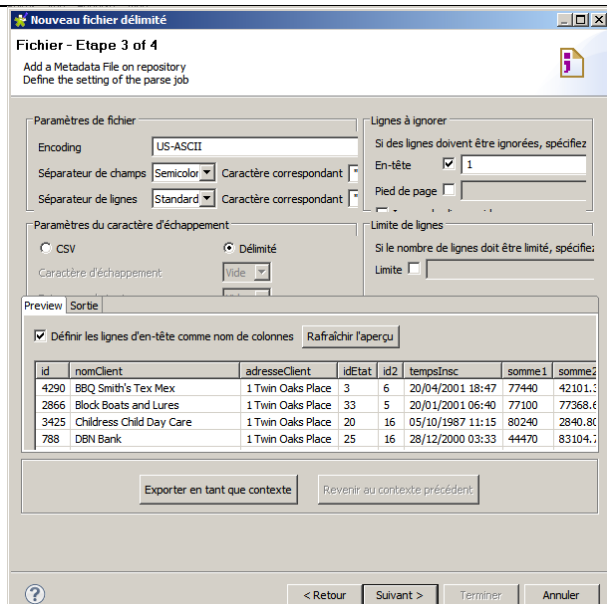
Dans ce TP, nous allons manipuler plusieurs sources de données (fichier CSV, fichier texte et base de données) pour en extraire les données, les transformer et les sauvegarder dans d'autres supports. La première étape à réaliser est de définir ces sources de données dans le Repository pour pouvoir générer leurs schémas et les utiliser dans les activités suivantes.

Pour faire cela, suivre les étapes suivantes :

<p>Dans le panneau (1) représentant le Repository, développer la section Métadonnées.</p> <p>Pour définir des sources sous forme de champs séparés par des délimiteurs (comme des fichiers csv ou texte), choisir : Créer un fichier délimité.</p> <p>Entrer le nom du fichier dans la fenêtre qui apparaît : client (dans notre cas, nous allons ajouter le fichier <i>client.csv</i>)</p>	
<p>Choisir ensuite le fichier que vous désirez ajouter.</p> <p>Naviguer pour cela vers le fichier client.csv qui vous a été fourni. Le visualiseur de fichier vous permet d'avoir une idée sur le contenu de ce fichier.</p> <p>Cliquer sur suivant.</p>	

Dans la fenêtre suivante, cliquer sur la case *Définir les lignes d'en-tête comme nom de colonne*. Cliquer ensuite sur *Rafraîchir* l'aperçu. L'aperçu du fichier extrait sera mis à jour, de manière à ce que la première ligne du fichier représente les noms des champs.

Cliquer sur *suivant*.

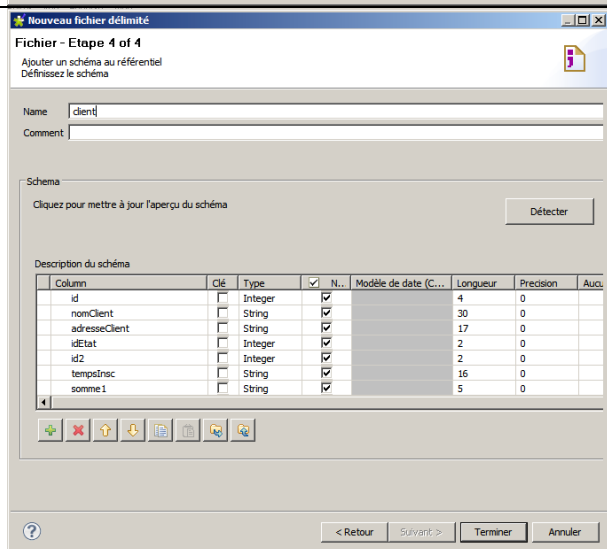


Modifier le nom du schéma du fichier délimité (*client*), et observer la composition des champs dans le panneau Description du schéma. Vous pourrez ainsi modifier les données du schéma à votre guise.

Dans notre cas, **ne pas oublier de cocher la case Clé pour le champ *id***.

Vous pourrez également modifier les longueurs des champs (les valeurs par défaut ont été calculées par Talend selon les données déjà présentes dans le fichier).

Cliquer sur *terminer*.



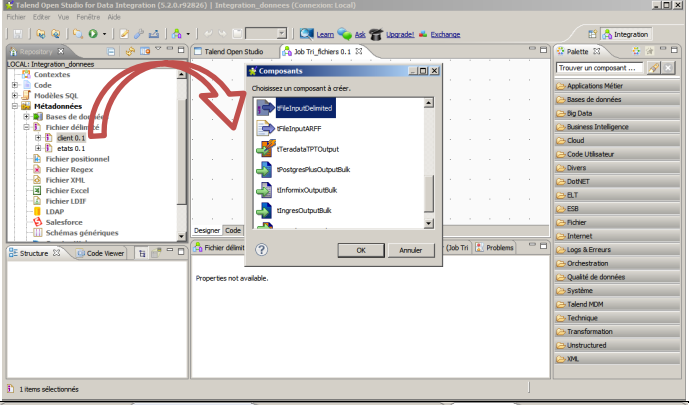
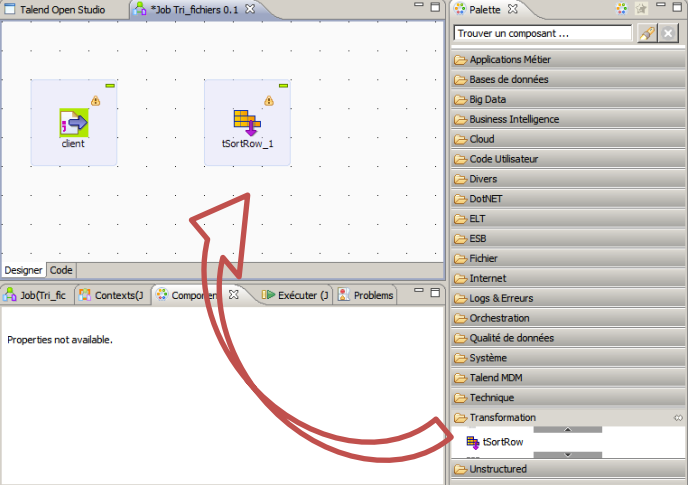
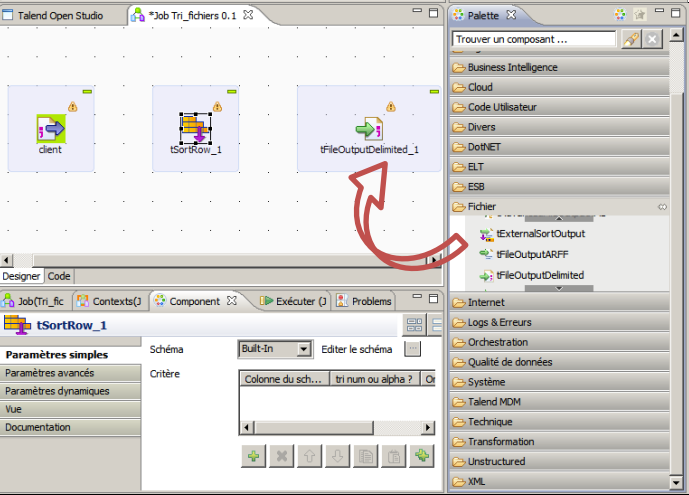
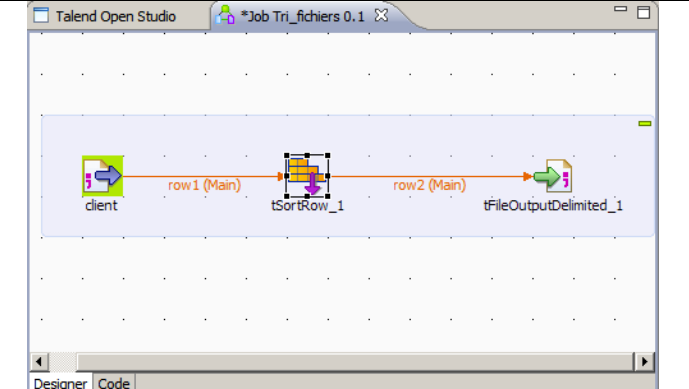
Vous avez ainsi ajouté un fichier source, dont le schéma pourra être utilisé dans toute l'application.

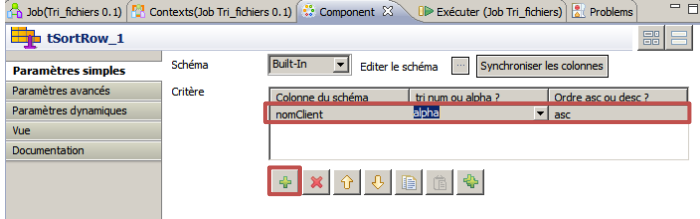
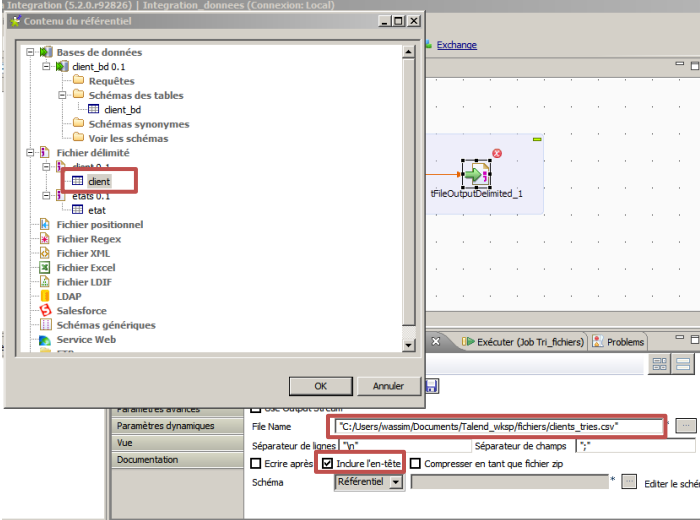
## Activité 1

- 1) Générer de la même manière le schéma du fichier état.txt qui vous est fourni.
- 2) Dans la base de données Oracle XE, créer une table appelée client. La structure de cette table n'a pas d'importance, elle sera écrasée plus tard.
- 3) Ajouter la base de données comme source dans la partie Métadonnées, et ajouter la table client aux schémas des tables.

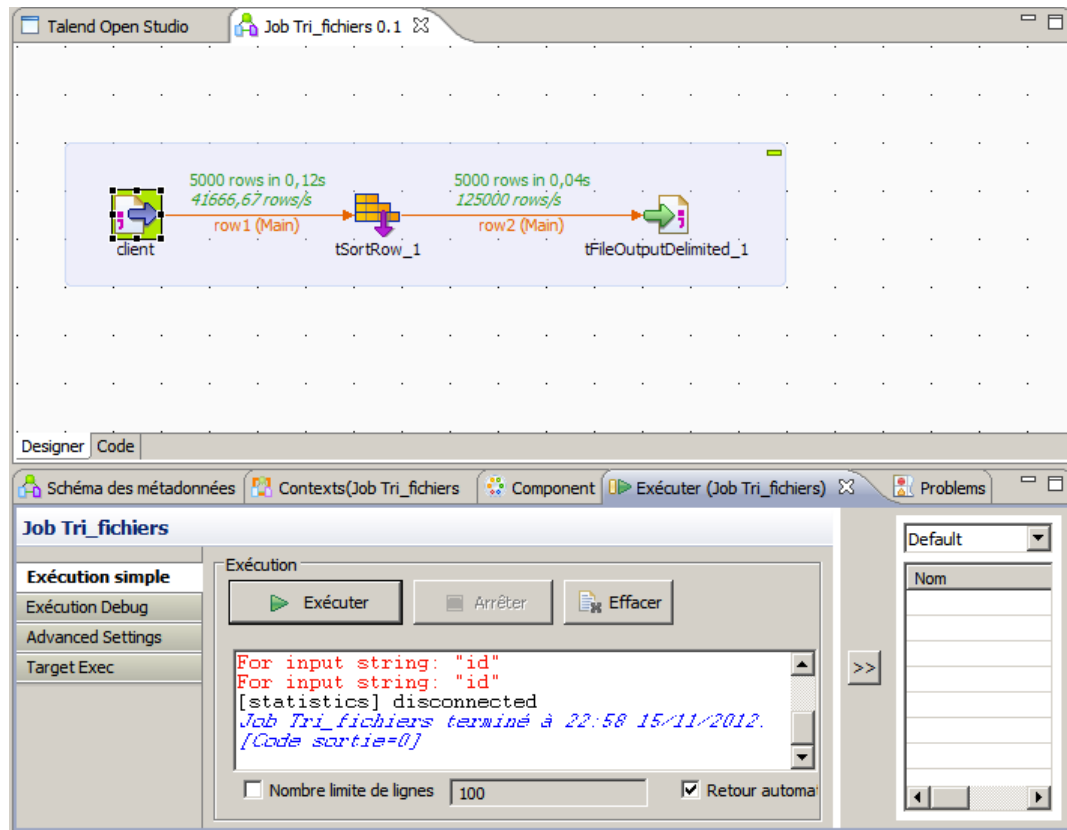
## 3) Tri de documents

Dans cette partie, on veut trier le contenu du fichier *client.csv* de manière automatique, en utilisant les composants Talend. Pour cela, suivre les étapes suivantes :

<p>Créer un nouveau Job que vous appellerez <i>Tri_fichiers</i></p>	
<p>Glisser le fichier délimité <i>client 0.1</i>, que vous avez créé précédemment, dans le panneau (2). Indiquer dans la fenêtre qui apparaît que c'est un <i>tFileInputDelimited</i>. Cliquer sur OK.</p>	
<p>Dans le panneau (4), représentant la palette, choisir le composant <i>tSortRow</i> dans la catégorie <i>Transformation</i>. Ce composant permet, comme son nom l'indique, de faire le tri d'un ensemble de données, selon une colonne particulière. Faire glisser ce composant dans la fenêtre principale.</p>	
<p>Pour représenter le fichier de sortie, faire glisser le composant <i>tFileOutputDelimited</i> dans la fenêtre principale. Il se trouve sous la catégorie <i>Fichier -&gt; Ecriture</i>.</p>	
<p>Relier les trois éléments pour représenter la chaîne d'exécution. Pour cela, faire un clic droit sur le composant <i>client</i>, maintenir enfoncé, et glisser vers le composant de <i>tri</i>. Faire de même entre le composant de <i>tri</i> et le fichier de sortie.</p>	

<p>Nous allons maintenant configurer les trois composants. Nous allons d'abord définir le nom du client comme critère de tri, par ordre alphabétique, du fichier source.</p> <p>Cliquer sur le composant de tri. Sous l'onglet <i>Composant</i> du panneau (3), cliquer sur (+). Modifier la valeur des champs insérés, pour faire le tri selon le nom de client, par ordre alphabétique ascendant.</p>	
<p>Cliquer ensuite sur le composant de sortie.</p> <ul style="list-style-type: none"> <li>- Choisir l'emplacement où on désire sauvegarder le fichier de sortie</li> <li>- Cocher la case : Inclure l'en-tête pour que l'entête des colonnes s'affiche dans le fichier de sortie</li> <li>- Devant la case Schéma, changer le type de schéma vers Référentiel, puis cliquer sur [...] à côté de Editer le schéma. Cela permettra de définir la structure des champs du fichier de sortie.</li> <li>- Dans la fenêtre affichée, choisir le schéma <i>client</i> du fichier délimité que vous avez créé.</li> </ul>	

Une fois ces étapes terminées, enregistrer le projet. Pour exécuter le processus, Cliquer sur l'onglet *Exécuter* du panneau (3), puis cliquer sur *Exécuter*. Ou alors taper F6. A la fin de l'exécution, la trace suivante est affichée sur la fenêtre principale :



Vérifier que le fichier trié a bien été créé dans le répertoire que vous avez spécifié plus tôt.

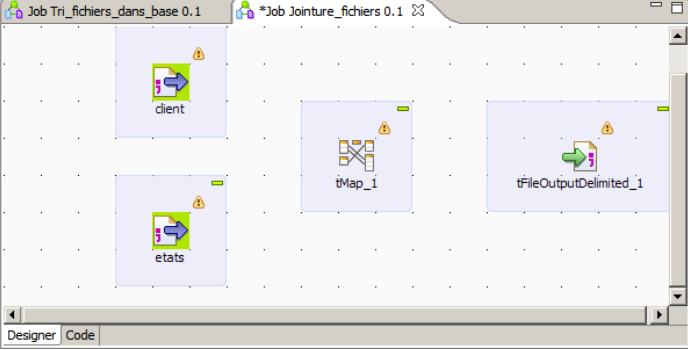
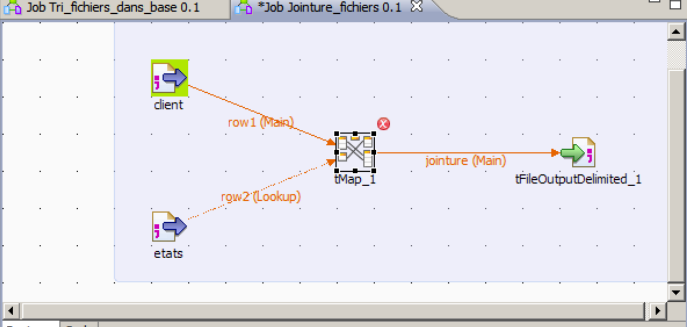
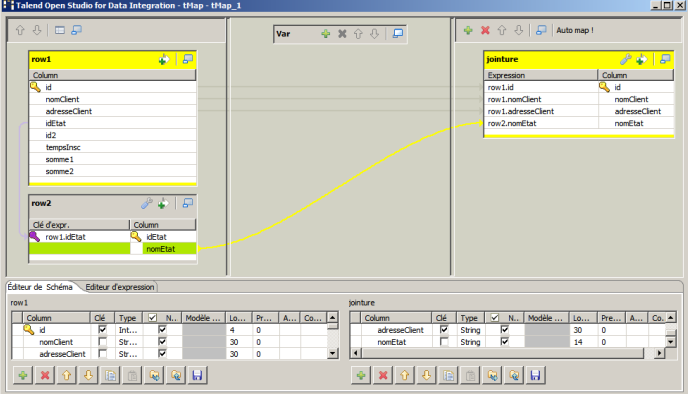
## Activité 2

- 1) Dupliquer le job *Tri\_fichiers* et le nommer *Tri\_fichier\_dans\_base*
- 2) Copier les données générées dans le fichier délimité de sortie dans la base de données *client\_bd* que vous avez créé dans l'activité précédente (au lieu d'un fichier CSV). A la création, la table cible sera écrasée et remplacée par la table contenant les données triées.

## 4) Jointure de fichiers

Le fichier *etat.txt* permet d'associer l'identifiant d'un état américain avec le nom de cet état. On se propose de faire la jointure des fichiers *client.csv* et *etat.txt* pour remplacer l'identifiant de l'état dans les données du client par son nom.

Pour faire cela, créer un nouveau Job *Jointure\_fichiers* et suivre les étapes suivantes :

<p>Glisser les deux fichiers délimités <i>client</i> et <i>etats</i> dans le panneau principal.</p> <p>Glisser le composant <i>tMap</i>, de la catégorie Transformation dans le panneau principal. Ce composant permet de transformer et diriger les données à partir d'une ou plusieurs sources vers une ou plusieurs destinations.</p> <p>Enfin, faire glisser un fichier délimité de sortie</p>	
<p>Relier les différents composants.</p> <p>Relier le fichier d'entrée <i>client</i> d'abord à la <i>tMap</i>, puis le fichier <i>etats</i>. Relier enfin le la <i>tMap</i> vers le fichier de sortie.</p> <p>Appeler la sortie <i>jointure</i>.</p>	
<p>Double cliquer sur la <i>tMap</i> pour la configurer. Une fenêtre s'ouvre.</p> <p>Commencer par relier le champ <i>idEtat</i> de la première table <i>row1</i>, au champ <i>idEtat</i> de la table <i>row2</i>.</p> <p>Faire glisser ensuite les champs <i>id</i>, <i>nomClient</i>, et <i>adresseClient</i> de <i>row1</i>, puis <i>nomEtat</i> de <i>row2</i> vers la table de destination <i>jointure</i>.</p>	

Configurer ensuite le fichier de sortie en précisant son chemin, et en incluant l'en-tête.

Exécuter le Job, et vérifier le fichier de sortie.

### Activité 3

1) Créer un nouveau Job *Jointure\_Tri\_fichiers\_de\_base*.

2) Ce job doit permettre de :

- Faire la jointure entre la table *client* créée dans l'activité 2 et le fichier *etat.txt* pour obtenir les champs *id*, *nomClient*, *adresseClient* et *nomEtat*.
- Trier ces données jointes par nom d'état, avant de les stocker dans un fichier texte *clients-etat.txt* dont les champs sont délimités par le caractère « | ».

## 5) Sélection des données

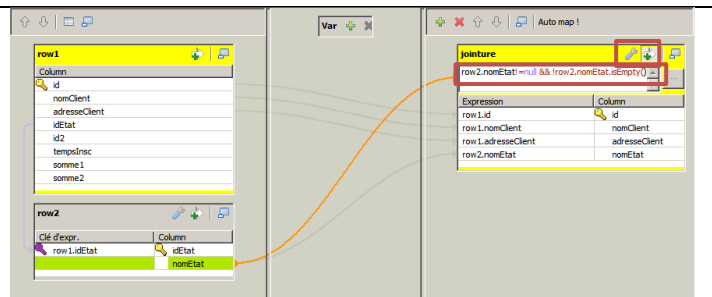
Il est possible de filtrer les données, en rejetant par exemple les entrées erronées. On peut remarquer dans les données du fichier *client.csv* que certaines entrées ne comportent pas de nom d'état. On désire filtrer ces données, et n'enregistrer dans le fichier de sortie que les données comportant un nom d'état. Les autres données pourront être affichées dans la console.

Dupliquer le Job *Jointure\_fichiers* et le renommer *Selection\_fichiers*.

Double-cliquer sur le composant *tMap* pour en modifier les propriétés.

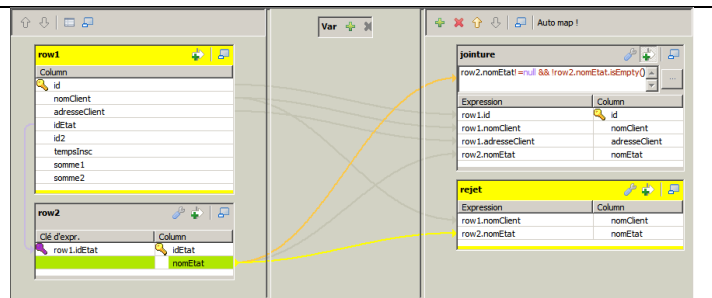
Activer le filtre des données, en cliquant sur la flèche de la table *jointure*. Entrer ensuite le critère de sélection des données (en Java) suivant :

```
row2.nomEtat !=null
&& !row2.nomEtat.isEmpty()
```



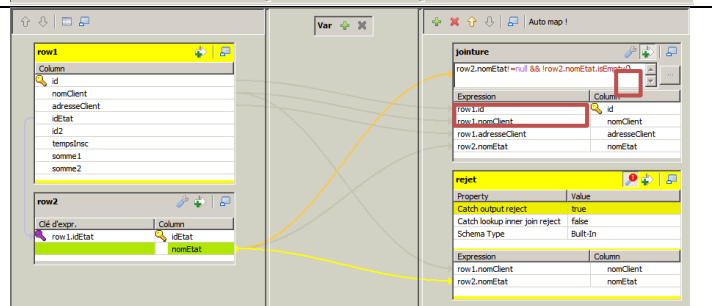
Créer une deuxième sortie appelée *rejets* en cliquant sur (+) au dessus de *jointure*.

Faire glisser les champs *nomClient* et *nomEtat* dans la table *rejets*.



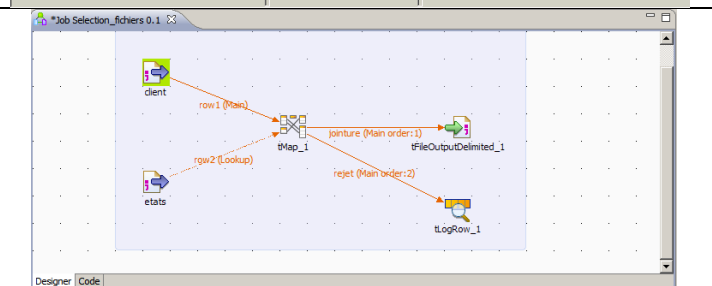
Indiquer que cette table contient les données rejetées par la sortie principale en cliquant sur le "tournevis" sur la table *rejets*, et en mettant le champ *catch output reject* à *true*.

Cliquer sur OK.



Faire glisser le composant *tLogRow* de la catégorie Logs et Erreurs dans la fenêtre principale.

Clic-droit sur la *tMap*, choisir *Ligne->rejet* et cliquer sur le composant de Log pour relier ces deux composants, et envoyer la sortie *rejet* vers le log.



Exécuter le Job et observer le résultat.



#### Activité 4

A partir des fichiers *client.csv* et *etat.txt*, réaliser les opérations suivantes :

1) Stocker dans une nouvelle table de la base de données les données jointes de ces deux fichiers, en respectant les règles suivantes :

- Les champs de la table seront : *id*, *nomClient*, *adresseClient*, *nomEtat*, *somme1*, *somme2*, *total* (où *total* est calculé en faisant la somme entre *somme1* et *somme2*)
- Stocker uniquement les clients dont l'état est « Alabama ».

2) Stocker le reste des enregistrements dans un fichier *reste.csv* dont la structure contient uniquement le *nomClient* et l'*état*.