

J. R. Almeida, L. Coelho and J. L. Oliveira

BIcenter: A collaborative Web ETL solution



Contents

List of Figures	v
List of Tables	vii
Preface	ix
1 Introduction	1
1.1 Main requirements	1
2 System Implementation	3
2.1 Requirements	3
2.2 Architecture	4
2.3 ETL SDK	6
3 Software User Manual	9
3.1 Initial views	9
3.2 Institutional features	10
3.3 ETL Task Editor	13
4 Guidelines for Developers	17
4.1 Branching Convention	17
4.2 Adding an REST API endpoint	18
4.3 Adding web view	20
4.4 Adding new PDI step	22



List of Figures

2.1	System workflow to build, manage and execute ETL tasks.	4
2.2	Components diagram of system architecture.	5
2.3	ETL task execution flowchart.	7
3.1	Login page.	9
3.2	Home page after login with all institutions assigned to the user.	10
3.3	Modal to create new institution.	10
3.4	Creation of connection to a data source (first step).	11
3.5	Modal to insert information regarding the new connection to a data source (second step).	11
3.6	Creation of connection to a remote server (first step).	12
3.7	Modal to insert information regarding the new connection to a remote server (second step).	12
3.8	Creation of a new ETL task (first step).	13
3.9	Overview of the ETL Task Editor.	13
3.10	Simple example of an ETL pipeline represented on this editor	14
3.11	Overview of part of the ETL components currently deployed in the system.	14
3.12	Main features of the ETL Task Editor.	15



List of Tables



Preface



1

Introduction

BIcenter is a web-based platform that allows the building and management of ETL pipelines, by non-IT users, in a multi-institution environment. Each institution manages and maintains ETL tasks and provides the resources for the execution of the associated tasks. Thus, each institution owns its private data sources, servers for ETL task execution and a task scheduler that allows periodic execution. In order to provide access and management control of ETL tasks and institutions, there are four distinct types of users:

- Administrator: Entity that moderates the platform. This actor has permissions to create and delete institutions.
- Resource Manager: Entity that manages private data sources and execution servers. This actor has permissions to create and delete private data sources and execution servers, within specific institutions.
- Task Manager: Entity that builds and executes ETL tasks. This actor can create and configure ETL tasks, within specific institutions.
- Data Analyst: This actor has permissions to inspect task execution history, namely the resulting data, execution logs and performance metrics.

1.1 Main requirements

Information Security Since ETL tasks parse and handle sensitive data that belongs to a particular institution, the system must be designed and implemented taking in account these security issues, namely user authentication, access control, data protection and isolation.

System Reliability Considering the periodic execution of ETL tasks, it is important to ensure that each execution is correctly initialized, started, motorized and concluded. When some fatal error occurs during an ETL

task execution, the system must be able to handle the error and successfully conclude the execution.

Solution Scalability Since a complete ETL tool typically encompasses a wide variety of components, it is crucial to build an agile approach to the development and integration of new ETL components.

2

System Implementation

This chapter contains the description of the main aspects of BIcenter implementation.

2.1 Requirements

The main goal of this application was to allow the building and management of ETL pipelines, by non-IT users, in a multi-institution environment. Therefore, some functional requirements were considered to implement a proper and reliable solution.

Each institution manages and maintains ETL tasks and provides the resources for the execution of the associated tasks. Thus, each institution owns its private data sources, servers for ETL task execution and a task scheduler that allows periodic execution. In order to provide access and management control of ETL tasks and institutions, there are four distinct types of users: 1) Administrator, entity that moderates the platform, and it has permissions to create and delete institutions; 2) Resource Manager, entity that manages private data sources and execution servers, and it has permissions to create and delete private data sources and execution servers, within specific institutions; 3) Task Manager, entity that builds and executes ETL tasks, and it can create and configure ETL tasks, within specific institutions; and 4) Data Analyst, actor with permissions to inspect task execution history, namely the resulting data, execution logs and performance metrics.

The functional model was derived from the identification of all key components from a top-down analysis of the system requirements. Figure 2.1 illustrates step by step all actions and interactions between the main actors in order to build and execute ETL tasks and to analyze the output results. Firstly, the administrator must create an institution. A resource

manager with access to the institution can create and configure all desired private data sources and task execution servers. Thereafter, a task manager can build and configure the ETL task and then schedule it for periodic or non-periodic execution. At the configured times the task will be sent to the execution server and executed. At the end of the execution, a notification will be sent to the data analysts, so they can check the execution results.

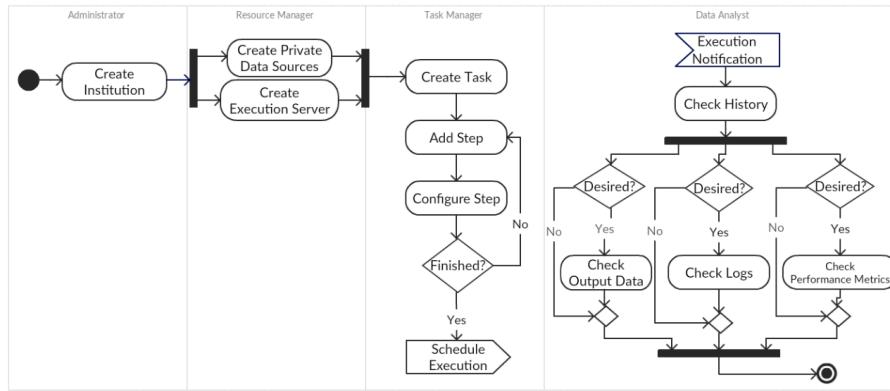


FIGURE 2.1: System workflow to build, manage and execute ETL tasks.

2.2 Architecture

The architecture considers four different tiers: 1) Application tier controls all application functionalities and maintains the system business logic; 2) Data tier is responsible for the maintenance of the private data sources; 3) Processing tier has the duty of executing and monitoring ETL task executions; and 4) Client tier is responsible for the solution presentation and page rendering. Figure 2.2 illustrates this architecture.

The system logic is built upon five different controllers: 1) RBAC controller provides user identity and evaluates access requests to resources; 2) Institution controller allows the creation and destruction of institutions, and resource allocation in institutions; 4) Task controller enables the creation, configuration and destruction of ETL tasks; and 5) Execu-

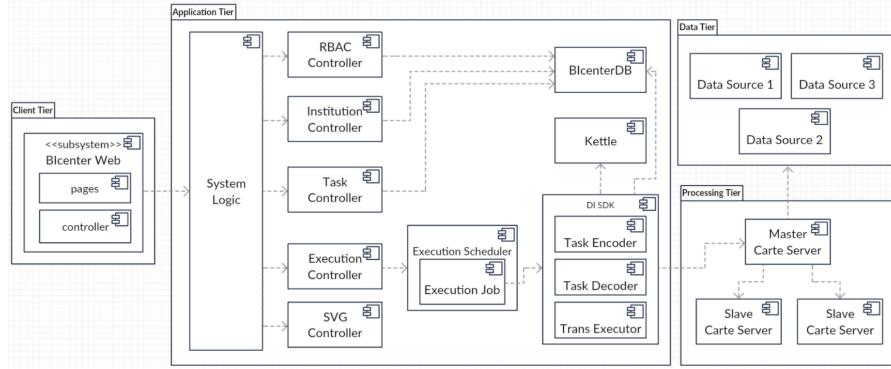


FIGURE 2.2: Components diagram of system architecture.

tion controller uses Kettle and Carte servers, so it can build, remotely execute and monitor ETL tasks.

Data Integration (DI) Software Development Kit (SDK) seeks to bridge the gap between Kettle and the stored information in BIcenter database (DB). It provides methods to build Pentaho's ETL processes according to the stored information, and also to execute them. Nevertheless, task execution is a periodic process. Therefore, Execution Scheduler manages the task execution scheduling. When the appropriated time arrives, a Execution Job is triggered. That job will communicate with DI SDK so that the given task is indeed executed. SVG controller maps between images and components, with the intent to build the visual pipeline.

BIcenter Web client adopted Model-View-Controller (MVC) architectural pattern. This pattern divides the web application in three different components:

- Model: represents the knowledge and data in an application. It has the responsibility to respond to information requests, proceed to information changes according to given instructions requests, and to notify observers in event-driven systems when information changes. Typically, the application data is stored in a database.
- View: it represents the user interface. The View updates the UI upon changes in the Model, by rendering the data into the suitable UI form.
- Controller: it handles events that occurs in the View, such as user interactions, and updates the Model accordingly.

2.3 ETL SDK

The ETL Task Editor was idealized to allow users to build a visual representation of the desired ETL pipeline. This editor must be able to process a similar structure to the ones used by the desktop ETL tools and generate the correspondent visual representation of the ETL task. MxGraph was used to build the ETL pipeline editor, and this Java/JavaScript diagramming library enables the building of interactive graphs. A graph consists in a set of cells. A cell can either be a vertex or an hop. Thus, a graph is formed by a group of vertices connected by edges. The vertices correspond to the ETL steps and the edges correspond to the ETL hops. Hence, in order to exhibit the desired task within the mxEditor, it is necessary to translate the task to the corresponding mxGraph object. GraphDecoder is responsible for creating a mxGraph and defining the appropriated model, according to the given Task. To insert vertices and edges in the graph model, a transaction must be created (`beginUpdate` and `endUpdate`). This is required for the model to remain in a consistent state. A default parent is automatically created and represents the first child of the root cell in the model. Subsequent elements must be added to the default parent.

The ETL SDK was built on top of PDI SDK. Kettle contains a rich set of data integration functionality that is exposed in a set of data integration tools. However, we can also use Kettle as a library in our own software and solutions. Pentaho Data Integration can be used as a Java API composed by three main components: 1) Core, that contains the core classes for Kettle; 2) Database, that contains the database-related classes; and 3) Engine, that contains the Kettle's runtime classes.

Figure 2.3 explains step by step how an ETL task is executed. The initialization of the Kettle environment loads all available plugins, initializes the logging environment and set up and reads the system variables. After the environment initialization, the transformation metadata is loaded and a transformation engine object is instantiated. Then, the execution is prepared and the transformation threads are started. Finally, because the whole transformation runs multi-threaded, it waits until all processing is completed.

In Kettle, an ETL process can be represented by six classes [31]:

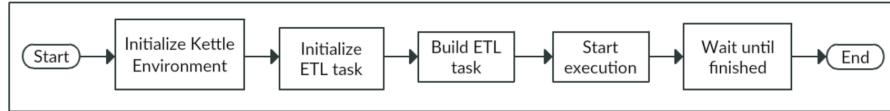


FIGURE 2.3: ETL task execution flowchart.

- TransMeta: is the class that defines the information about the ETL process and offers methods to save and load it from XML, as well as methods to alter an ETL process by adding/removing databases, steps, hops, etc.
- Trans: represents the information and operations associated with the concept of an ETL process. It can loads, instantiates, initializes, runs, and monitors the execution of the ETL process.
- DatabaseMeta: defines the database specific parameters for a certain database type.
- StepMeta: is the class that defines the information about a ETL process's Step.
- TransHopMeta: defines a link between two steps in an ETL process.
- BaseStep: represents the information and operations associated with the concept of an ETL process Step. It offers initialization, row processing and step clean-up methods.

Initially a TransMeta must be properly configured with all proper DatabaseMetas, StepMetas and TransHopMetas. TransDecoder is responsible to create a TransMeta and define all underlying information, according to the given Task. GenericStep is a class that has a generic algorithm to encode or decode a given Step to or from a Kettle's Step-Meta. In order to provide interface segregation, TransDecoder must use the StepDecoder interface that only provides the GenericStep decoding method.

After having a correctly configured TransMeta, it can be executed with the associated Trans object. TransExecutor is a singleton responsible for the initialization, execution and monitoring of the Trans object. Moreover, it writes and stores the execution logs, step measures and status, and the resulting data in real-time. To accomplish this, listeners are coupled to the Trans and Steps objects in order to obtain the results and metrics throughout the transformation's execution.



3

Software User Manual

This chapter presents the main interfaces of BIcenter, aiming to help users to navigate through the system.

3.1 Initial views

The login page presented in Figure 3.1 allows users to be authenticated using their credentials. This can be connected with the institution LDAP, which provides a central place to store usernames and passwords.

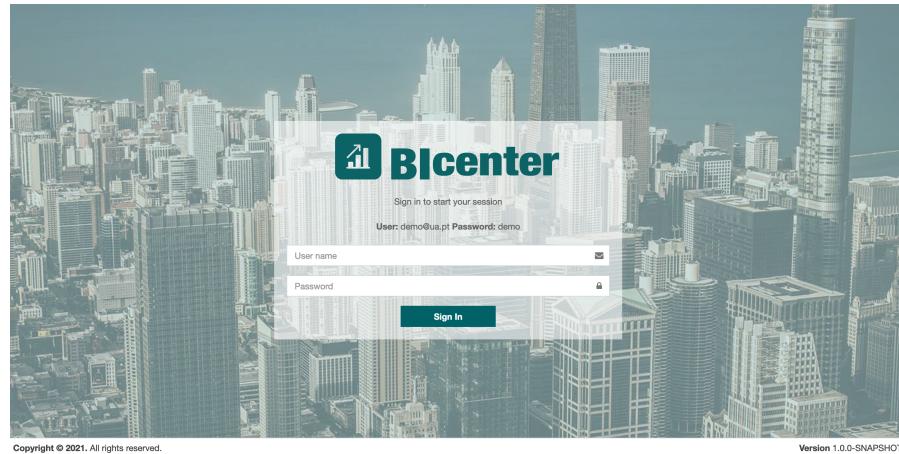


FIGURE 3.1: Login page.

After the authentication process, the user is redirected to the application's home page represented in Figure 3.2. This page contains all the institutions that the user belongs to. In this case, this user has permission to create a new instance of an institution in the application using the modal illustrated in Figure 3.3.

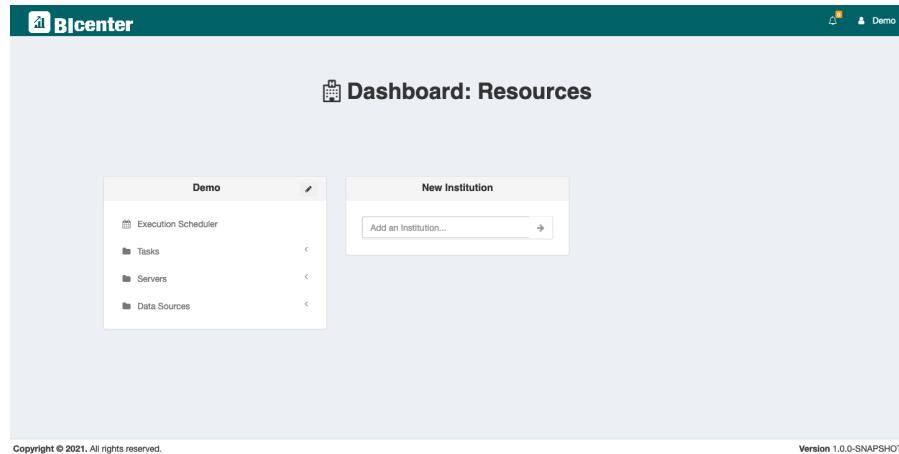


FIGURE 3.2: Home page after login with all institutions assigned to the user.

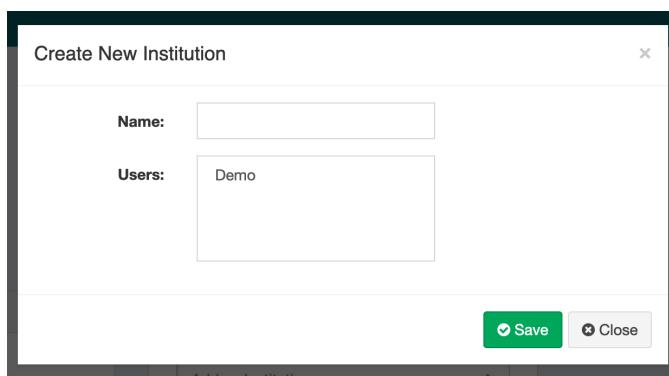


FIGURE 3.3: Modal to create new institution.

3.2 Institutional features

The configuration of a data source is illustrated in Figures 3.4 and 3.5. Firstly is created the data source instance in the institution, and then, the details of this connection are inserted in the modal represented in Figure 3.5.

The connection to the remote server follows a similar flow. This configuration is illustrated in Figures 3.6 and 3.7.

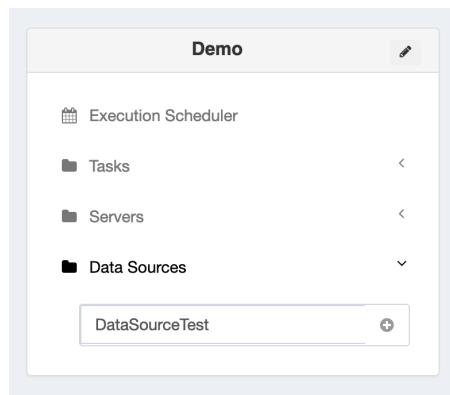


FIGURE 3.4: Creation of connection to a data source (first step).

Edit Data Source ×

Institution:	<input type="text" value="Demo"/>
Connection Name:	<input type="text" value="DataSourceTest"/>
Connection Type:	<input type="text"/>
Connection Method:	<input type="text"/>
Host Name:	<input type="text"/>
Port Number:	<input type="text"/>
Database Name:	<input type="text"/>
Username:	<input type="text"/>
Password:	<input type="text"/>

✖ Delete ✔ Save ✖ Close

FIGURE 3.5: Modal to insert information regarding the new connection to a data source (second step).

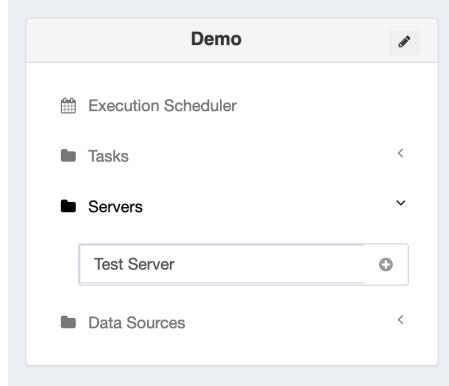


FIGURE 3.6: Creation of connection to a remote server (first step).

Edit Remote Server

Institution:	Demo
Name:	Test Server
Host Name:	
Port Number:	
Username:	
Password:	

FIGURE 3.7: Modal to insert information regarding the new connection to a remote server (second step).

Finally, the ETL tasks are also created for the institutions. This feature redirects the users to the ETL Task Editor, in which they can design and implement an ETL pipeline using the web interface. The creation of a new task can be done by using the option represented in Figure 3.8.

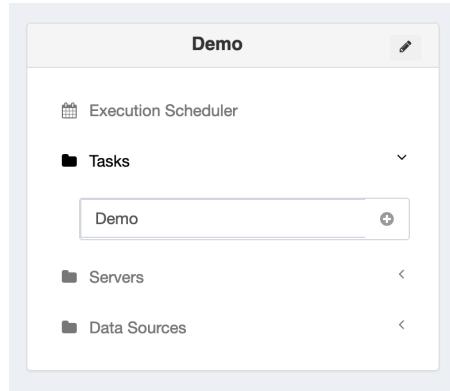


FIGURE 3.8: Creation of a new ETL task (first step).

3.3 ETL Task Editor

The ETL Task Editor allows users to create an ETL pipeline by dragging and dropping ETL components. Figure 3.9 shows an overview of this web editor.

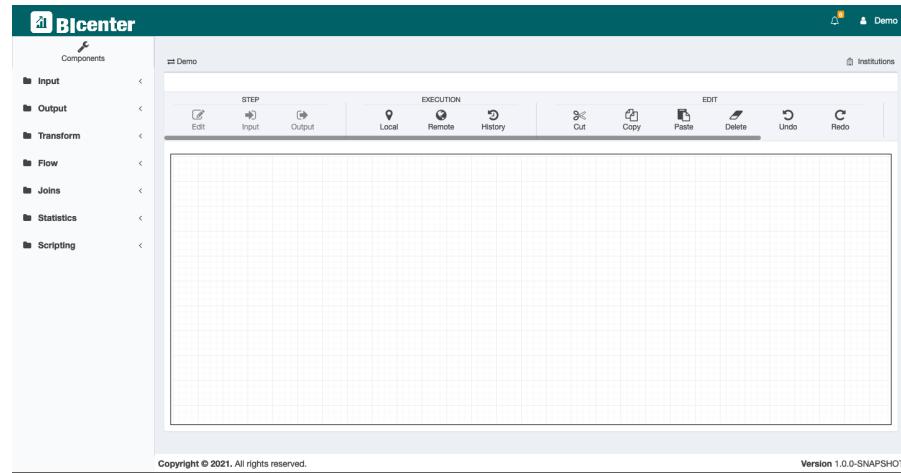


FIGURE 3.9: Overview of the ETL Task Editor.

Figure shows an example of an ETL pipeline implemented in the ETL Task Editor using the four most common ETL components.

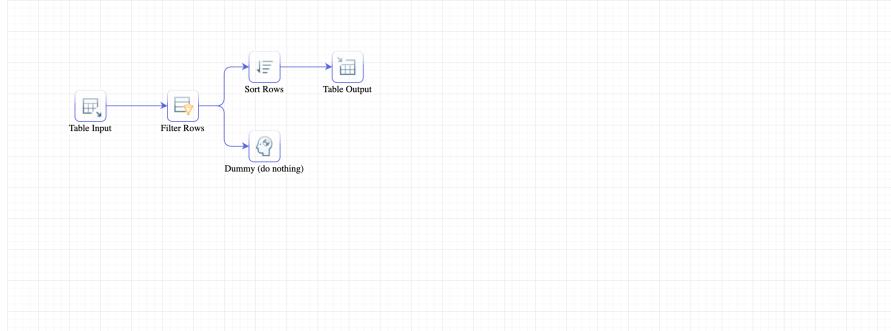


FIGURE 3.10: Simple example of an ETL pipeline represented on this editor

The components available to implement these ETL pipelines are available on the left menu. Part of this menu is represented in Figure , which is expanded the components for input, output and transformations.

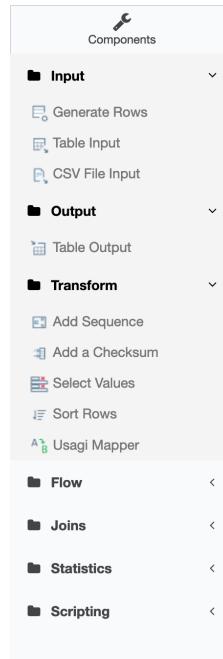


FIGURE 3.11: Overview of part of the ETL components currently deployed in the system.

The navigation bar illustrated in Figure allows the uses to operate over

the ETL Task. This menu is divided into sub-menus: 1) step, which contains the features to configure each ETL component in the pipeline and analyse the input and output of each step; 2) execution, which defined if the ETL task will be executed locally or in a remote server; edit, contain the editing usual features, such as cut, copy, paste, delete, undo and redo; and 4) select, which is related with the drawing features.



FIGURE 3.12: Main features of the ETL Task Editor.



4

Guidelines for Developers

In this chapter, it is described the main guidelines for contributing to BIcenter.

4.1 Branching Convention

After some consideration, we decided to implement naming conventions in the branches being created. We followed a scheme similar to the one referred at <http://stackoverflow.com/a/6065944>

Each branch name is composed of the following:

category/*name*

The possible categories are listed below.

Category	Description
bug	Bug fixing
imp	Improvement on already existing features
new	New features being added
wip	Works in progress - Big features that take long to implement and will probably hang there
junk	Throwaway branch created to experimentation

Category	Description
bug	Bug fixing
imp	Improvement on already existing features
new	New features being added
wip	Works in progress - Big features that take long to implement and will probably hang there
junk	Throwaway branch created to experimentation

The name should be concise, and directly represent what the branch solves.

Some examples:

bug/issue234

bug/fixeditdb

new/statistics

junk/tryingbootstrap3

4.2 Adding an REST API endpoint

1. Create the endpoint at `app > controllers`, either on one of the existent JAVA files or a new one if it handles a completely different topic/task, yet to be addressed.
2. Add the respective endpoint URL to:
 - `conf > routes`, so that the system knows that endpoint exists;
 - `app > controllers > Application.java > javascriptRoutes()`, so you can use this URL dynamically at JavaScript, regarding frontend code.

Since the structure of this project includes both Frontend and Backend, when you create a REST API endpoint, you'll probably consume it on some frontend feature. Here's how to do that:

1. At `app > javascripts > services`, there is a collection of JS files that have all the encapsulated logic responsible for communicating with the REST API: one file for each major entity or group of actions. So, in one of these files or a new one, add the function that handles the communication you need with the REST API through your newly developed endpoint.
2. At the JS controller level of your functionality, you call the function specified in the previous step with the necessary arguments and information.

4.2.1 Example

This example assumes that there isn't an endpoint capable of providing information about all the existent institutions. So the purpose is to create such endpoint and provide a function that allows us to consume it on the frontend. Hence, there are 2 phases: the **creation of the REST API**

endpoint and the **creation of the necessary utility to consume it** on the frontend.

4.2.1.1 Creation of REST API endpoint

Because the endpoint we want to create has the purpose of returning the information about all existent institutions, it should be placed at `app > controllers > InstitutionController.java`, with some code like this:

Please check the online documentation to access this code.

After creating the endpoint, we need to make it discoverable (be available to the outside world). To do that, we add its URL to: 1. `conf > routes`:

Please check the online documentation to access this code.

```
2. app > controllers > Application.java >
   javascriptRoutes():
```

Please check the online documentation to access this code.

4.2.1.2 Creation of utility function to consume the REST endpoint

This second phase explains how to create the referred utility function so we can more easily consume the developed endpoint. Since it's an institution-related method, we will be adding the following utility function to `app > javascripts > services > institution.js`.

Please check the online documentation to access this code.

This consumes the endpoint and, on success calls the specified callback function; otherwise (on fail) displays an error message on your browser's console.

After this, the final part is to consume this utility function, which can also be called a service. A simple example is to use the service on your javascript file that bears the “backend” responsibility of your UI component/view, in a similar way to the following code:

Please check the online documentation to access this code.

4.3 Adding web view

If you want to create an entirely new base schema (a view), you first need to add a base template name `<name>.scala.html` to `app > views`. You also need to make sure that you have an endpoint on JAVA code that handles your request and returns this web view (the “**backend**” process is similar to create a REST API endpoint).

If you just want to create a new UI stage, based on an already defined view, follow these steps:

1. Create an endpoint (and register it) that returns the base view you want to use.

2. At `app > assets > javascripts`, create the controller and view associated with your new UI stage (the view JS file handles the stuff UI related and the controller handles the stuff Backend related).
3. At `app > assets > templates`, create a file named `<name>.handlebars` where you’ll insert the UI code that will be dynamically added to the main div container of your base web view.
4. After all these files are ready, you need to make sure the framework knows the path to everything you created, so you’ll need to add the paths to your view and controller javascript files in the file `app > assets > javascripts > main.js`.
5. Finally, you’ll need to add the RegExp handler that will assemble the resources you need (variables and controllers) at `app > assets > application > application.js`.

4.3.1 Example

This example is going to be divided into two big parts: the **creation of a new view from scratch** and the **creation of a UI component** (through handlebars).

Let’s take as an example the creation of a dashboard page, given the fact that we need to create a new base-view from scratch.

NOTE: if you just want to create a new UI component, with handlebars, jump to the second part.

4.3.2 Creation of a new (base) view

To create a new view to use as a base to new UI stages, you first need to add the HTML template. So, add your view structure to a file named `home.scala.html`, for the sake of this example, at `app > views`:

Please check the online documentation to access this code.

After creating the template, you now need to create an endpoint to return this view and register it in the system.

So, to create the endpoint you add a new controller at `app > controllers` or add a function to an existent controller. In this example we will create a new controller called `HomeController.java`:

Please check the online documentation to access this code.

And, then, to register it, add to: `1.conf > routes`:

Please check the online documentation to access this code.

```
2. app > controllers > Application.java >
   javascriptRoutes():
```

Please check the online documentation to access this code.

4.3.3 Creation of a UI component (handlebars)

Once the URL for the base view is configured, the next step is to create the functional files that handle frontend and backend of the UI element. So, we will create a folder `home` at `app > assets > javascripts`. Inside this folder create the files: `homeController.js` and `homeView.js`.

The `homeController.js` is responsible to handle the backend operations and should follow a structure similar to this:

Please check the online documentation to access this code.

On the other hand, the `homeView.js` is responsible for frontend interactions and should follow a structure similar to:

Please check the online documentation to access this code.

As can be seen in this last code snippet, we “import” a web (HTML) structure from 'home'. This refers to the file `home.handlebars` that must be placed at `app > assets > templates`. Hence, we will create this exact file in the specified location with the following content:

Please check the online documentation to access this code.

The next step is to tell the system where the newly-created resources are. So, at `app > assets > javascripts > main.js` add the following:

Please check the online documentation to access this code.

Last, but not least, we have to create the RegExp handler, as said earlier, so the system know what controllers and tools to load in the specific UI stage, at `app > assets > application > application.js`:

Please check the online documentation to access this code.

With this, we will be able to have a web page similar to this image, based on an entirely new HTML template:

New HTML view created using a new UI component following the handlebars pattern.

4.4 Adding new PDI step

Adding and registering a new **Pentaho Data Integration** Step is a straightforward process which has been streamlined by the design of the BICenter project. A list of all available PDI Steps can be found in this link¹.

After selecting the Step to be added, the following is the process required to add said step into the system.

1. Register the new step on the `public > editor > diagrameditor.xml` file.
2. Add a new object to `conf > configuration.json` under the subsection pertaining to the new component's type with the desired component properties.

¹<https://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+Steps>

3. If necessary, add extra functionality to the `submitClick` method on the `app > assets > javascripts > step > stepController.js` or the `applyChanges` method on `app > assets > javascripts > services > step.js`.
4. Create a new class on `app > diSdk > step > parser` with the appropriate name.
5. If the component requires pre-processing or extra logic you can alter the `decodeStep` method on `app > diSdk > step > AbstractStep.java`

4.4.1 Example

The following is an example of how to add the CSVFileInput component. For more information about this component's functioning you can check out this link².

1. Firstly, we have to add the following line to `public > editor > diagrameditor.xml`

Please check the online documentation to access this code.

2. Next we'll have to register the component's properties on the `conf > configuration.json` file. As the CSV File Input step is an Input type component, we'll be registering it under the Input Components. We'll be naming this object the same name we used on the `template` field on the prior step. Note that the `shortName` field should have a name that goes in accordance with the naming specified by the Pentaho Kettle library (in our specific case, this can be seen in this link³). Under the `componentProperties` array we can specify the multiple inputs and fields of our component.

```
{
    "name": "CSVInput",
    "description": "CSV File Input",
    "shortName": "csvInput",
    "componentProperties": [
        {

```

²<https://wiki.pentaho.com/display/EAI/CSV+File+Input>

³<https://javadoc.pentaho.com/kettle/org/pentaho/di/trans/steps/csvinput/CsvInputMeta.html>

```
"name": "Step Name",
"shortName": "stepName",
"type": "input"
},
{
  "name": "File Name",
  "shortName": "fileName",
  "type": "fileinput"
},
{
  "name": "Delimiter",
  "shortName": "delimiter",
  "type": "input"
},
{
  "name": "Enclosure",
  "shortName": "enclosure",
  "type": "input"
},
{
  "name": "NIO Buffer Size",
  "shortName": "bufferSize",
  "type": "number"
},
{
  "name": "File Encoding",
  "shortName": "encoding",
  "type": "select",
  "componentMetadatas": [
    {
      "value": "UTF-8",
      "name": "UTF-8"
    },
    {
      "value": "ANSI",
      "name": "ANSI"
    }
  ]
},
{
```

```
        "name": "Lazy Conversion?",  
        "shortName": "lazyConversionActive",  
        "type": "checkbox"  
    }  
]  
}  
}
```

3. As BICenter is already prepared to receive files and all of our CSVInput's parameters, we can skip this step.
4. Now we have to create a new class on the `app > diSdk > step > parser` directory. This can simply be done by copying any of the other classes already present in the folder. Don't worry about the lack of logic in this class, as this serves as a mere extension of the `AbstractStep` class, which itself contains all logic needed to communicate with the PDI, used in order to allow our component to be detected and processed.

Please check the online documentation to access this code.

5. As we want our system to automatically detect the CSV's fields automatically without the users having to manually introduce them themselves, we have to add some logic to the `decodeStep` method on `app > diSdk > step > AbstractStep.java` which will do an initial read of the file in order to extrapolate it's columns' names. This can be done by creating a new method on this class and adding it to the `decodeStep` method, or by injecting the logic directly into the function.

Please check the online documentation to access this code.