

Tutorial: Building a HULFT Integrate Project for Screen Scraping

Introduction to HULFT Integrate

[HULFT Integrate](#) is a flexible and powerful tool for connecting to a wide variety of data sources for reading and writing information, and provides a comprehensive set of features for manipulating, massaging, and enriching that data.

HULFT Integrate utilizes [adapters](#) to connect with data sources. Given a typical use case involving the movement and manipulation of data, we probably have the appropriate adapters available out-of-the-box to meet your requirement.

HULFT Integrate Project: Screen Scraping Use Case

A customer had purchased a new system from a software vendor, yet still needed to extract data from the system for integration and subsequent reporting needs. The particular system in question utilized a database for a backend, but the license restricted the customer from accessing the database directly. The customer was told that an API could be built within 12 months that would allow them to extract the data they needed, but waiting up to 12 months wasn't an option. The customer needed an immediate, if not temporary, solution.

The customer decided to use a screen scraping, or web scraping approach, since the user interface was rendered in a web browser, as are many modern systems. There are many reasons why screen scraping is a poor approach from a technical perspective, and not recommended for a long-term solution. In this case, the customer was looking for a short term solution until the software vendor delivered a true API.

With this knowledge, the customer approached HULFT. Due to the flexibility of HULFT Integrate, HULFT could help with both the short term goal of screen scraping as well as the long term goal of transitioning to a true API for integration.

Building the Project Step-by-Step

Let's use an example scenario where we want to extract the top ten currency exchange rates from a website, so we can utilize that information for further processing. Here is the website we want to extract the information from, with a highlight around the top ten exchange rates.

The screenshot shows the X-RATES website interface. The main content area displays a table titled "RATES TABLE" with the subtitle "1 US Dollar Rates Table". The table lists the top 10 currencies and their exchange rates against the US Dollar as of June 12, 2019, 20:47 UTC. The table is highlighted with a red box. The data is as follows:

Top 10	1.00 USD	inv. 1.00 USD
US Dollar	0.885947	1.128736
Euro	0.788233	1.268661
British Pound	69.435786	0.014402
Indian Rupee	1.443617	0.692705
Australian Dollar	1.334142	0.749545
Canadian Dollar	1.366981	0.731539
Singapore Dollar	0.995634	1.004385
Swiss Franc	4.158569	0.240467
Malaysian Ringgit	108.518444	0.009215
Japanese Yen	6.917702	0.144557

Below the main table, there is an "Alphabetical order" section with a table of exchange rates for various currencies. The data is as follows:

Alphabetical order	1.00 USD ▲ ▼	inv. 1.00 USD ▲ ▼
US Dollar ▲	43.669327	0.022899
Argentine Peso	1.443617	0.692705
Australian Dollar	0.376000	2.659574
Bahraini Dinar	10.844004	0.092217
Botswana Pula	3.865076	0.258727
Brazilian Real	0.788233	1.268661
British Pound	1.366981	0.731539
Bruneian Dollar	1.732761	0.577114
Bulgarian Lev	1.334142	0.749545
Canadian Dollar		

On the right side of the page, there is a section titled "Percent Change in the Last 24 Hours" with the following data:

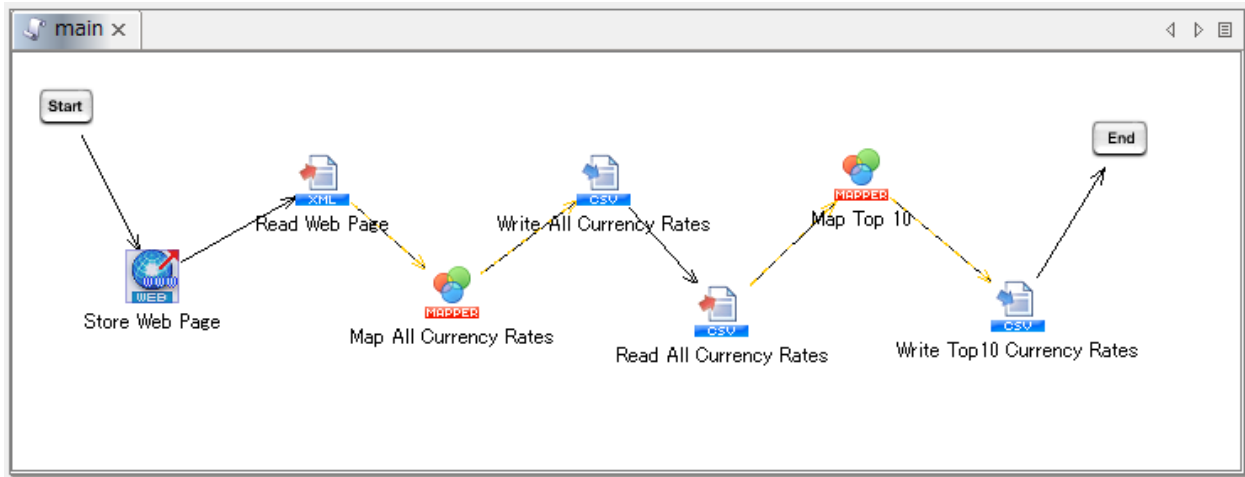
Percent Change in the Last 24 Hours		
EUR/USD	-0.34959%	USD/JPY +0.01216%
GBP/USD	-0.28034%	USD/CHF +0.35048%
USD/CAD	+0.43716%	EUR/JPY -0.33747%
AUD/USD	-0.49178%	CNY/USD -0.08972%

Below this section, there is a "Useful Links" section with links to "Euro Information", "FAQ", and "Feedback".

To create the HULFT Integrate project to accomplish this task, we need to leverage several built-in adapters from the palette;

- Crawl – Allows connection to a web site and retrieval of underlying HTML
- XML File Read – Allows us to read the HTML created by the Crawl adapter
- CSV File Write – Allows creation of a CSV text file containing the scraped data
- CSV File Read – Read a temporary CSV file created as part of the overall flow

Here is a screenshot of the completed HULFT Integrate project detailing the integration flow:



And here is a screenshot of the final CSV output showing the top ten currency exchange rates:

	A	B	C	D
1	Currency	1 USD in Cur	1 Cur in USD	
2	Euro	0.885873	1.128829	
3	British Pound	0.788187	1.268735	
4	Indian Rupee	69.39764	0.01441	
5	Australian Dollar	1.44354	0.692742	
6	Canadian Dollar	1.334085	0.749577	
7	Singapore Dollar	1.366872	0.731597	
8	Swiss Franc	0.995662	1.004357	
9	Malaysian Ringgit	4.158503	0.240471	
10	Japanese Yen	108.522481	0.009215	
11	Chinese Yuan Renminbi	6.917778	0.144555	
12				

Now let's walk through each of the icons in the integration flow and explain their purpose:

1. Store Web Page (Crawl adapter)

This adapter connects to the specific web URL to retrieve the web page and store it into a local file

Properties for Crawl(File/Parameter) operation

Crawl(File/Parameter) operation

Set properties for Crawl(File/Parameter) operation.

Name:

Required settings | Authentication | Cookie | Comment

Destination:

Path:

Method:

Encoding:

Parameter

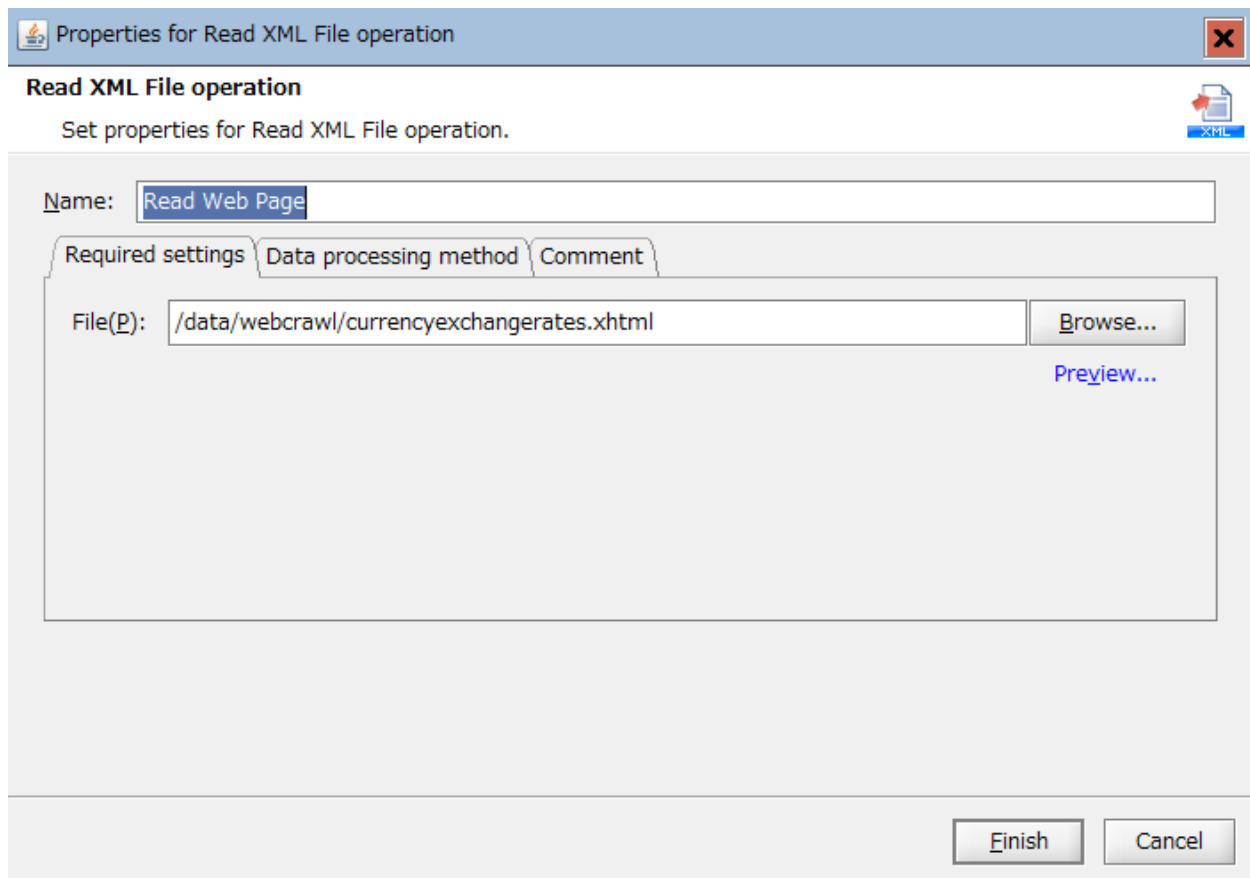
Header name	Value
-------------	-------

Output file:

XHTML format:

2. Read Web Page (XML File Read adapter)

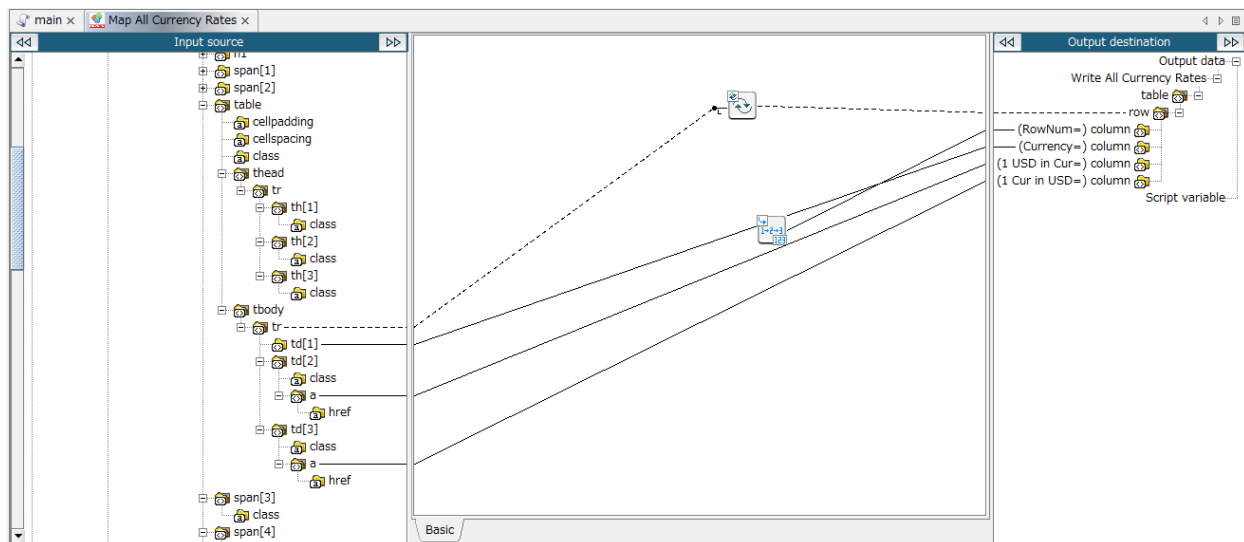
This adapter allows us to read the file created by the previous step



3. Map All Currency Rates

This component allows us to extract the precise fields, and nothing more, from the web page XML file and map to the structure of a CSV output file. This requires a hierarchical structure that's mapped to a relational structure, and shows the power of the mapping component. The schema for the hierarchical structure is automatically created by using the Load Schema feature of the mapper to read the web page in XML format and generate a corresponding schema. Another very powerful feature of HULFT Integrate's Mapper feature.

Due to the nature of the web page design, there is one table representing the top 10 exchange rates plus all the other exchange rates. At this point we capture all the rates and tag each output row with an incrementing row number. Later in the process we reduce the final output to just the top 10 rates (the first 10 rows).



4. Write All Currency Rates (CSV File Write adapter)

Here we take the output from the previous Mapper component and write to a local CSV file.

Properties for Write CSV File operation

Write CSV File operation

Set properties for Write CSV File operation.

Name:

Input data:

Required settings | Write settings | Transaction | Comment

File(P): [Browse...](#)

[Preview...](#)

Delimiter mode: Select from list Enter directly Enter character code

Delimiter:

Column list

Column name	Quotation
RowNum	<input type="checkbox"/>
Currency	<input type="checkbox"/>
1 USD in Cur	<input type="checkbox"/>
1 Cur in USD	<input type="checkbox"/>

[Up\(U\)](#)
[Down\(D\)](#)
[Add\(A\)](#)
[Delete\(R\)](#)

[Update column list...](#)
[Read column names from first row of file\(O\)...](#)
[Read the number of columns from file\(G\)...](#)

[Finish](#) [Cancel](#)

5. Read All Currency Rates (CSV File Read adapter)

This adapter reads all rows from the previously created CSV file containing all exchange rates

Properties for Read CSV File operation

Read CSV File operation

Set properties for Read CSV File operation.

Name:

Required settings | **Read settings** | Multi-thread processing settings | Data processing method | Comment

File(P): [Browse...](#)
[Preview...](#)

Delimiter mode: Select from list Enter directly Enter character code

Delimiter:

Column list

Column name
RowNum
Currency
1 USD in Cur
1 Cur in USD

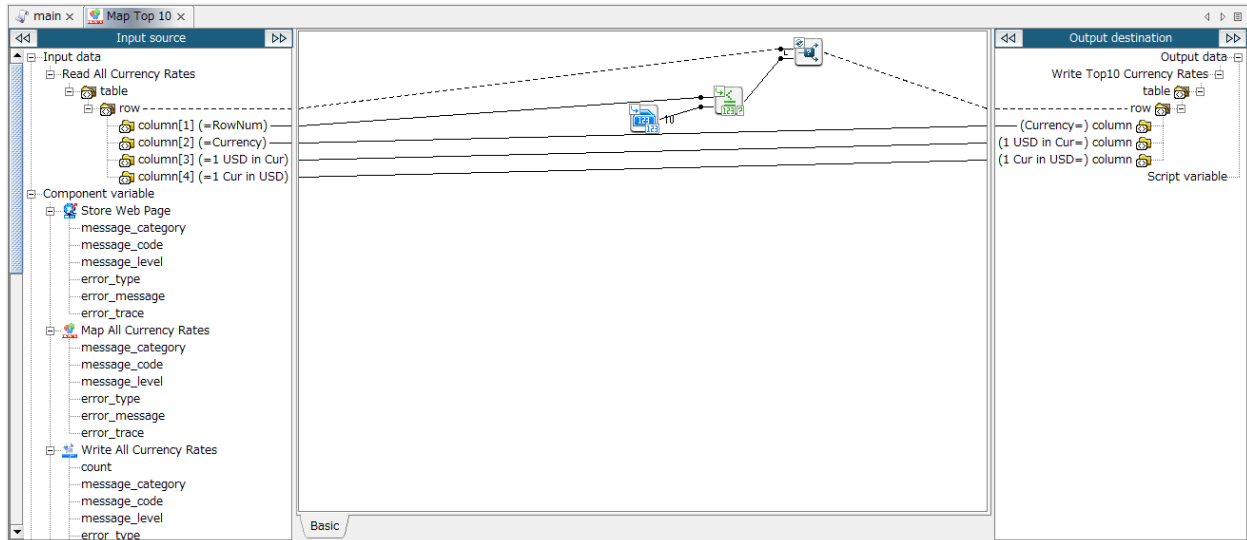
[Up\(U\)](#) [Down\(D\)](#) [Add\(A\)](#) [Delete\(R\)](#)

[Update column list...](#)
[Read column names from first row of file\(O\)...](#)
[Read the number of columns from file\(G\)...](#)

[Finish](#) [Cancel](#)

6. Map Top 10 (Mapper)

This component filters down the number of rows by only processing the first 10 rows and preparing that data for storage



7. Write Top 10 Currency Rates (CSV File Write adapter)

Takes the output from the previous Mapper and stores the exchange rate data into the final output file.

Properties for Write CSV File operation

Write CSV File operation

Set properties for Write CSV File operation.

Name:

Input data:

Required settings | Write settings | Transaction | Comment

File(P): [Browse...](#)
[Preview...](#)

Delim~~i~~ter mode: Select from list Enter directly Enter character code

Delim~~i~~ter:

Column list

Column name	Quotation
Currency	<input type="checkbox"/>
1 USD in Cur	<input type="checkbox"/>
1 Cur in USD	<input type="checkbox"/>

[Up\(U\)](#)
[Down\(D\)](#)
[Add\(A\)](#)
[Delete\(B\)](#)

[Update column list...](#)
[Read column names from first row of file\(O\)...](#)
[Read the number of columns from file\(C\)...](#)

[Finish](#) [Cancel](#)

Summary

HULFT Integrate provides flexibility for your data integration problems. As we have demonstrated in this tutorial, some use cases may not seem to be an immediately obvious choice, but HULFT Integrate is adaptable to a variety of uses and data sources.

Here we showed how HULFT Integrate can be used to accomplish screen scraping, for a case where there were no other legitimate means of extracting data in the short term. As I mentioned, screen scraping is an inadequate choice for long term integration needs due to its inherent fragility. There are always exceptions of course. For example, when screen scraping is a short term solution until a more viable long term option is available. In that case HULFT Integrate can be used to accomplish both the immediate and long term needs for integration.

Other Resources:

HULFT Integrate [product sheet](#).

HULFT Integrate connects to all your diverse data sources, no matter where they reside. Here's [a list of our adapters](#).

More than 10,000 customers across 43 countries trust HULFT. Learn more about our customers [here](#).

Want to learn more?



Visit us at
<https://hulftinc.com>



Call us at
855-815-1518



Email us at
salesop@hulftinc.com