

BiNHum Harvesting and Indexing Toolkit

Contents

Pre-requisites.....	1
Basics	1
Web-services	1
Configuration.....	2
Edit the application.properties file	2
Edit applicationContext-security.xml.....	2
Before starting the app.....	3
Starting B-HIT.....	3
Overview.....	3
Add a datasource.....	4
Retrieve datasets	4
Perform inventory	5
Harvest data.....	5
Process harvested records.....	5

Pre-requisites

Basics

Tested with Java 8

Can be deployed with Tomcat (tested with Apache Tomcat/6.0.39 JRE 1.8.0_05-b13) or Jetty (tested with Jetty 8.1.14), or can be run inside Eclipse (tested with Eclipse Luna and Kepler)

MySQL Database v5.5

Web-services

Local gisraphy installation (coordinates check) or maybe online version, see <http://www.gisraphy.com/free-access.htm>

CoordinatesKDTree webservice (based on <https://github.com/AReallyGoodName/OfflineReverseGeocode>)

Configuration

Edit the application.properties file

Quality tests	
• Quality tests: activated / deactivated. The tab won't be displayed if turned OFF	qualityOnOff=OFF / ON
• Where the CoordinatesKDTree Web service is running	coordinatesWebservice=http://localhost:8080/CoordinatesWS
A local folder for temporary file manipulation	temporaryFolder=/home/user/test
Database properties	
• DB credentials	dataSource.username=USER dataSource.password=PASS
• DB name	dataSource.name=DATABASE_NAME
• MySQL-Server IP + DB name	dataSource.url=jdbc:mysql://DATABASE-IP:3306/DATABASE_NAME?autoReconnect=true&useUnicode=true&characterEncoding=UTF8&characterSetResults=UTF8&MVCC=true
Application properties	
• Application server IP (ie. Localhost)	dataSource.servername=APPLICATIONSERVER_IP
• default URL with port (ie. http://localhost:8040/) The port must match your Tomcat/Jetty configuration !	baseUrl=http://APPLICATIONSERVER:PORT/
• directory (with all permissions for Tomcat/Jetty) where to store the harvested files	harvest.directory=/home/USER/PROJECT/

Edit applicationContext-security.xml

Edit the admin password, md5 encoded. The default password corresponds to « banana! »

```
<user name="admin" password="bb7a307e32b93a931da89d0a214dd47f" authorities="ROLE_ADMIN" />
```

Before starting the app

Run the database creation script (schemaOnly.sql).

Starting B-HIT

Based on your configuration, and the name chosen (eclipse config / war file name), open your favorite web-browser to <http://localhost:8040/Bibhum/datasource/list.html>.

Overview

The screenshot shows the B-HIT application interface. At the top, there is a header bar with the GBIF logo, the text "free and open access to biodiversity data", and the HUMBOLDT-RING logo. On the right side of the header is a "login" form with fields for "username" and "password". Below the header is a navigation menu with tabs: Datasources (highlighted in red), Associated Data, Extra units, Jobs, Console, Report, Datasource Management, Data quality summary, and Data viewer.

The main content area is titled "BioDatasource List" and contains a "menu" button. It provides instructions for adding new datasources and managing existing ones. A table below lists available methods, including columns for Provider Name, Datasource, URL, Target, Harvested, Started, Last Inventory, Last Inv. processed, Last harvesting, Last harv. processed, and Country. The table shows 0 entries.

At the bottom of the page, there are buttons for "schedule", "add bioDatasource", and links for "Previous" and "Next". The footer indicates the page is based on Version 1.48, from GBIF | © 2012- BINHum.

Main menu:

- Datasources : main entry point to add datasources/datasets, launch metadata operations, inventory+harvesting+processing data
- Associated datasources: entry point for associated data (relationships), to harvest and process associated data only
- Extra units: entry point for single units retrieval, based on a list of unit IDs
- Jobs: overview of waiting and running jobs
- Console: overview of log events
- Report: generation of reports (statistics)
- Datasource management: to hide or delete datasets
- Data quality : to launch quality tests, exports test results and display results
- Data viewer: to display data stored in the database, either the raw data or the improved data from the quality tests

Add a datasource

Click on the button « add bioDatasource »

BioDatasource Detail

Configure your BioDatasource

BioDatasource Name:	<input type="text"/>	Save	cancel
Provider Name (abbrev):	<input type="text"/>		
=Accesspoint URL:	<input type="text"/>		
Provider Full Name:	<input type="text"/>		
Provider Website URL:	<input type="text"/>		
Provider Address:	<input type="text"/>		
Factory class:	<input type="text" value="DiGIR Metadata Factory"/>		
Country:	<input type="text"/>		

Enter a name for the datasource, the provider name abbreviated (ie. BGBM), the provider fullname (ie. Botanischer Garten und Botanisches Museum Berlin-Dahlem), the accesspoint (ie. http://ww3.bgbm.org/biocase/pywrapper.cgi?dsa=test_ABCD21), the provider website (ie. http://www.bgbm.org), the provider postal address and country, and the type of datasource (ie. Biocase/digir/Tapir/darwincore archive). Then save. The new datasource is now listed on the http://localhost:8040/Bibhum/datasource/list.html page.

To see the available operations for this datasource, click on the « available methods » checkbox on the left.

Available Methods	Provider Name	Datasource	URL	Target	Harvested	Started	Last Inventory	Last Inv. processed	Last harvesting	Last harv. processed	Country
<input checked="" type="checkbox"/> Metadata update	BGBM	Test ABCD	http://ww3.bgbm.org/biocase/pywrapper.cgi?dsa=test_ABCD21								Germany

Retrieve datasets

In order to trigger this « Metadata update » operation, check the box and click on « Schedule » at the bottom of the page.

A new Job is created and can be followed under the « Job » tab (/Bibhum/job/list.html).

Datasources	Associated Data	Extra units	Jobs	Console	Report	Datasource Management	Data quality summary	Data viewer
-----------------------------	---------------------------------	-----------------------------	---	-------------------------	------------------------	---------------------------------------	--------------------------------------	-----------------------------

Job List

A view of all operations that have been scheduled and are awaiting execution. Please note that the maximum number of operations that can be run in parallel is 500.

<input type="button" value="kill"/>	<input type="checkbox"/> all:	<input type="checkbox"/> reschedule:	<input type="text"/> id:	<input type="button" value="Show / hide columns"/>
Show <input type="text" value="10"/> entries				Search: <input type="text"/>
ID	Name	Description	Created	Started
1	IssueMetadata	BGBM : Test ABCD -> ww3.bgbm.org	2015-03-04 16:05:16	

Showing 1 to 1 of 1 entries

Previous Next

Refresh the page to update the job status :

ID	Name	Description	Created	Started
1	IssueMetadata	BGBM : Test ABCD -> ww3.bgbm.org	2015-03-04 16:05:16	2015-03-04 16:06:10

Going back to the main page, a new line appeared in the list of datasource for the discovered dataset, with its corresponding available methods. These methods depend on the type of provider. ABCD-Archive will enable the methods XXXXXXXXXXXXXXXXXX, DwC-A will enable YYYYYYYYYYYY.

Available Methods	Provider Name	Datasource	URL	Target	Harvested	Started	Last Inventory	Last Inv. processed	Last harvesting	Last harv. processed	Country
<input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/>	BGBM	Test ABCD - Herbarium Berolinense	http://www3.bgbm.org/biocase /pywrapper.cgi?ds=tes_ABCD21	192551	0						Germany
<input type="checkbox"/>	BGBM	Test ABCD	http://www3.bgbm.org/biocase /pywrapper.cgi?ds=tes_ABCD21			04-03-2015					Germany

DEBUGGING HINT: The most common reason that a metadata update fails, is that the accesspoint is wrong.

The target columns contains the theoretical number of units, according to the providers metadata.

Perform inventory

Select the first operation and click on « schedule ». The queries are stored on the disc, in compressed files. The location path is defined in the configuration file, and the directory for each datasource can be found in the MySQL table bio_datasource, column « basedirectory ».

Harvest data

Data will be harvested if and only if a list of units has been retrieved with the Inventory. The outputs from this operation are:

1. One or more search requests (with enumerated extensions corresponding to the order in which they were dispatched, i.e. search_request.000)
2. One or more search responses (with enumerated extensions corresponding to the order in which they were dispatched, i.e. search_response.000). Often, there will only be a single response per request, but sometimes there can be multiple responses for a single request!)

The default range is 200 units per request.

Process harvested records

In this operation, an Operator (BioDatasource) will collect all the search responses, parse them, and write the parsed values to the database. In case of re-harvesting a dataset, the system will check if the data changed since the last harvesting. It will calculate the new checksum of each downloaded file, and will modify data in the DB if and only if the checksum changed. Each checksum is stored in the DB table « sha1responses ».