



gbif-indexingtoolkit

The GBIF Harvesting and Indexing Toolkit (HIT)

[Project Home](#)[Downloads](#)[Wiki](#)[Issues](#)[Source](#)[Export to GitHub](#)

Search

Current pages



for

FAQ

Frequently Asked Questions

Featured

Updated May 5, 2011 by [kleg...@gmail.com](#)

A list of frequently asked questions

Questions are organised into 3 different categories:

- [GENERAL questions](#)
 - [What is the HIT?](#)
 - [Who are its end-users?](#)
 - [What could it be used for?](#)
 - [How does it work?](#)
 - [What protocols does it support?](#)
 - [What are some of its main features?](#)
 - [What are its current limitations?](#)
- [INSTALLATION questions](#)
 - [What should I do when getting a java.security exception?](#)
 - [What should I do when getting an java.lang.OutOfMemoryError PermGen space?](#)
- [OPERATION/USAGE questions](#)
 - [Why do I get a Connection timed out error?](#)
 - [Why do I harvest LESS than 100% of the target records?](#)
 - [Why is there a problem with the ranges in my name ranges file?](#)
 - [Why do I harvest MORE than 100% of the target records?](#)

GENERAL questions

What is the HIT?

It is a simple to use, extensible open source framework that allows you to easily manage data harvesting and quickly build specific indexes of harvested data.

Who are its end-users?

Data aggregators, or anybody who wants to build a central index from a network of distributed resources.

What could it be used for?

A Node or thematic community would use it to populate their own central index on top of which a portal could reside. GBIF will continue to use it to manage its internal indexing operations.

How does it work?

1. The HIT communicates with a central registry, gathering information about data publishers and their available access-points.
2. The HIT's operator then decides which datasets to schedule for harvesting.
3. Harvested records get saved to intermediary text files.
4. These records are then synchronised with the central index.

What protocols does it support?

Currently, DiGIR, TAPIR, BioCASE, and Darwin Core Archive are all supported. The HIT is entirely customisable, however, and could be extended to include other protocols.

What are some of its main features?

- A user-friendly dashboard from which to manage all harvesting activities.

[Datasources](#)
[Jobs](#)
[Console](#)
[Registry](#)

BioDatasource List

An overview of all BioDatasources managed locally, divided into 2 categories: metadata updaters and operators. Metadata updaters gather information about the number of resources behind a given access point and create a new operator for each one. The operators are then used to manage and perform actions against that individual resource located at the given access point.

Data provider:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

List: All Deleted Metadata Updaters: ALL DIGIR BioCASE TAPIR Operators: ALL DIGIR BioCASE TAPIR

Select: All None

Datasources	Provider	URL	Target	Max Har	Harvested	Dropped	Harvesting Start
<input type="checkbox"/> Lichen DB (www3.bgbm.org) - Lichen Herbarium	Botanic Garden an	http://www3.bgbm.org/biocase/lycoper	84891	0	0	0	2009-07-21 11:17:03
<input type="checkbox"/> typecollection of fossil sponges and corals of Indo	NLBI	http://145.18.182.102/tacirink/tacir.php	3709	3709	3709	0	
<input type="checkbox"/> Europa (Atlas) - atlas-europa	GBIF-Sweden	http://www.gbif.se/tacir/tacir.php?alt=eu	251	248	248	3	
<input type="checkbox"/> Birds (GBIF-SE-Andatabanken) - sq-birds	GBIF-Sweden	http://www.gbif.se/tacir/tacir.php?sq-birds	0	1475958	1475958	-1475958	
<input type="checkbox"/> FishBase - fishbase	FishBase	http://www.gbif.se/tacir/tacir.php/fishbase	1054983	971184	971184	83799	
<input type="checkbox"/> GBIF-Spain (Baray.cak.es) - Estacion Biologica Do	GBIF-Spain	http://baray.cak.es/6005050n/DIGIR.php	23902	0	0	23902	
<input type="checkbox"/> Forster herbarium, Göttingen (GOET) (web-test2)	Universitätsherbari	http://web-test2.uni-goettingen.de/bioc	0	0	0	0	
<input type="checkbox"/> Tylos herbarium, Göttingen (GOET) (web-test2)	Universitätsherbari	http://web-test2.uni-goettingen.de/bioc	0	0	0	0	
<input type="checkbox"/> Israel Nature and Parks Authority (www3.bgbm.org)	Israel Nature and F	http://www3.bgbm.org/biocase/lycoper	433463	429852	429852	3611	
<input type="checkbox"/> http://www.geo-antennivall.de (www.geo-antenniv	GEO-Tag der Arten	http://www.geo-antennivall.de/biocase)	0	0	0	0	
<input type="checkbox"/> http://www.geo-antennivall.de (www.geo-antenniv	GEO-Tag der Arten	http://www.geo-antennivall.de/biocase)	0	0	0	0	
<input type="checkbox"/> http://www.geo-antennivall.de (www.geo-antenniv	GEO-Tag der Arten	http://www.geo-antennivall.de/biocase)	21	21	21	0	
<input type="checkbox"/> http://www.geo-antennivall.de (www.geo-antenniv	GEO-Tag der Arten	http://www.geo-antennivall.de/biocase)	0	0	0	0	

- Multi-threaded execution, allowing several harvests to be carried out concurrently.

[Datasources](#)
[Jobs](#)
[Console](#)
[Registry](#)

Job List

A view of all operations that have been scheduled and are awaiting execution. (The order of the queue going from top to bottom).

Name	Location	Created	Started	Next Fire
synchroise	Birds (GBIF-SE-Andatabanken) -> http://www.gbif.se/tacir/tacir.php/sq-birds	2009-09-16T16:45:11	2009-09-16T16:45:11	
synchroise	Lepidoptera (Observations) -> http://www.gbif.se/tacir/tacir.php/logn-nrm	2009-09-16T16:45:11		
synchroise	Gothenburg Herbarium - General (GBIF-IB-GB-Herbarium) -> http://www.gbif.se/tacir/tacir.php/gn-gen	2009-09-16T16:45:11		
synchroise	Swedish Board of Fisheries (Fiskeriverket) - Swedish Electrofishing Registry -> http://www.gbif.se/tacir	2009-09-16T16:45:11		
synchroise	Lund Botanical Museum (LD) -> http://www.gbif.se/tacir/tacir.php/ld-general	2009-09-16T16:45:11		
synchroise	Nordic Herbarium (N) -> http://www.gbif.se/tacir/tacir.php/nordicarb	2009-09-16T16:45:11		
synchroise	Bird Ringing Centre in Sweden (NRM) -> http://www.gbif.se/tacir/tacir.php/nrm-vingsbirds	2009-09-16T16:45:11		
synchroise	Mammals of the Gothenburg Natural History Museum -> http://www.gbif.se/tacir/tacir.php/gnm-mamm	2009-09-16T16:45:11		
synchroise	Fishes of the Gothenburg Natural History Museum -> http://www.gbif.se/tacir/tacir.php/gnm-fishes	2009-09-16T16:45:11		
synchroise	Gothenburg Herbarium - Types (GBIF-IB-GB-Herbarium) -> http://www.gbif.se/tacir/tacir.php/gn-types	2009-09-16T16:45:11		
synchroise	Invertebrates (GBIF-SE-SMBN) -> http://www.gbif.se/tacir/tacir.php/lev	2009-09-16T16:45:11		
synchroise	NRM-Fishes -> http://www.gbif.se/tacir/tacir.php/nrm-fishes	2009-09-16T16:45:11		
synchroise	Herbarium of Oskarshamnen (OHN) -> http://www.gbif.se/tacir/tacir.php/ohn	2009-09-16T16:45:11		

Version 1.0Beta-SNAPSHOT | [Project Home](#) | [Bug Report](#) | © 2009 GBIF

- Harvesting and indexing operations entirely decoupled from each other, allowing for increased performance and flexibility.
- Controlled workflows.
- Comprehensive logging, permitting the user to quickly identify errors.

[Datasources](#)
[Jobs](#)
[Console](#)
[Registry](#)

LogEvent List

A view of all log messages being generated by the application, automatically refreshed every few seconds.

Console mode: [verbose](#) | [normal](#) | [minimal](#)

```

2009-06-09 10:29:59.0 Finished process harvested
2009-06-09 10:29:59.0 Header line of harvested tab file has been written successfully
2009-06-09 10:29:59.0 Start process harvested
2009-06-09 10:29:57.0 Finished harvest
2009-06-09 10:29:57.0 Success harvesting range [ Acaena - Vulpia megalura ]
2009-06-09 10:29:57.0 END_OF_RECORDS : true
2009-06-09 10:29:57.0 RECORD_COUNT : 859
2009-06-09 10:29:56.0 Writing to file: /tmp/Amphibian/argbf/central.gov.ar/tnsbot/search_response.000
2009-06-09 10:29:39.0 Writing to file: /tmp/Amphibian/argbf/central.gov.ar/tnsbot/search_request.000
2009-06-09 10:29:39.0 Start harvesting range [ Acaena - Vulpia megalura ]
2009-06-09 10:29:39.0 Start harvest
2009-06-09 10:29:29.0 Finished process inventoried
2009-06-09 10:29:29.0 Name ranges file has been written
2009-06-09 10:29:29.0 Start process inventoried
2009-06-09 10:29:22.0 Finished inventory
2009-06-09 10:29:22.0 END_OF_RECORDS : true
2009-06-09 10:29:22.0 RECORD_COUNT : 419
2009-06-09 10:29:22.0 Writing to file: /tmp/Amphibian/argbf/central.gov.ar/tnsbot/inventory_response.000
2009-06-09 10:29:19.0 Writing to file: /tmp/Amphibian/argbf/central.gov.ar/tnsbot/inventory_request.000
2009-06-09 10:29:19.0 Start inventory
  
```

- Data validation.

What are its current limitations?

- It can only index into a database mirroring the GBIF data structure. Adopters could add extensions allowing for indexing into other structures.
- It does not contain time-based job scheduling (Due in 1.1 release candidate)

INSTALLATION questions

What should I do when getting a java.security exception?

On platforms such as Ubuntu, the Java security manager can be quite restrictive. To make them less restrictive, and hence to allow the application to be able to run, it's necessary to modify the catalina.policy file, which is how the security policies implemented by the Java SecurityManager are configured. This file is located in the tomcat conf directory. Open this file, and add the following snippet of code (replacing the \${tomcat.directory} with the actual path where your tomcat installation is located).

```
// HIT webapp permission
grant codebase "file:${tomcat.directory}/webapps/hit/-" {
    permission java.io.FilePermission "${tomcat.directory}/webapps/hit/-",
"read,write,delete";
    permission java.util.PropertyPermission "*", "read,write";
    permission java.util.logging.LoggingPermission "control";
    permission java.lang.RuntimePermission "getClassLoader";
    permission java.lang.RuntimePermission "preferences";
    permission java.lang.RuntimePermission "shutdownHooks";
    permission java.lang.reflect.ReflectPermission "suppressAccessChecks";
};
```

Some additional properties other users have added to get things running include:

```
permission java.net.SocketPermission "*:3306", "connect, resolve";
permission java.lang.RuntimePermission "getProtectionDomain";
permission java.util.PropertyPermission "cglib.debugLocation", "read";
```

Then restart Tomcat, and the exception should disappear. Otherwise, create a new issue or post a question to the mailing list.

What should I do when getting an outOfMemoryError permGen space?

The problem here is that the PermGen space is too small. To fix this, it is necessary to increase its size by playing with the flags - XX:PermSize=128m and -XX:MaxPermSize=256m (the default value is 64M).

Please see the section entitled [Tomcat configuration](#) for help with setting this JVM environment variable.

Hopefully this fixes the problem, otherwise, try incrementing it in increments of 128MB until the error goes away (assuming you actually have sufficient memory to allocate)

What should I do when getting a SEVERE: Exception starting filter struts java.lang.NullPointerException at com.opensymphony.xwork2.util.FileManager\$FileRevision.needsReloading

This error first appeared on Windows XP. It seems like this happens when tomcat is installed in a directory that has spaces, such as C:\Program Files\apache-tomcat-6.0.20\ Moving the installation to C:\apache-tomcat-6.0.20 fixed the problem, and allowed the application to run.

OPERATION/USAGE questions

Why do I get a connection timed out error?

- Answer: It could be because your connection to the internet has been lost.
- Answer: It could be because a port, or ip has been blocked on your machine. It is a good idea to try pinging the publisher you are trying to connect to from the same machine where the HIT is installed. It is also a good idea to check your firewall settings.

Why do I harvest LESS than 100% of the target records?

- Answer: It could be because there was a problem constructing the name ranges file. To ensure that the proper name ranges were constructed:
 - examine the name ranges file for any peculiarities - [relates to this FAQ](#)
- Answer: It could be because there was a problem parsing some of the XML responses. To double check, you could do the following:
 - check the logs to see that there were no parsing errors
 - check the actual search response(s), to see that they in fact contain records and that these correspond to the appropriate name range.

Why is there a problem with the ranges in my name ranges file?

- Answer: The most probable answer, is that there was a problem with the inventory response in that it did not contain proper scientific names. To double check: it would be advisable to:
 - check that the inventory response(s) was actually returned
 - check that there were no XML parsing errors

- check the logs to see that no illegal characters had to be replaced from any of the names
- check the logs to see that there were no invalid scientific names as these may have been skipped

Why do I harvest MORE than 100% of the target records?

- Answer: One source of inflated record count in Darwin-Core archives can be illegal line terminating characters. A record containing such a character would break in two and appear to the parser as two lines with an insufficient number of columns. Consequently these two lines would be replaced by blank lines but still appear in the record count turning a single line into two.
 - Search for lines containing line terminating characters *inside* the records and remove these

► [Sign in](#) to add a comment

[Terms](#) - [Privacy](#) - [Project Hosting Help](#)

Powered by [Google Project Hosting](#)