

# DSpace Statistics

## DSpace Statistics

DSpace 1.6 and newer versions uses the Apache SOLR application underlying the statistics. SOLR enables performant searching and adding to vast amounts of (usage) data.

Unlike previous versions, enabling statistics in DSpace does not require additional installation or customization. All the necessary software is included.

- 1 [What is exactly being logged ?](#)
- 2 [Web user interface for DSpace statistics](#)
  - 2.1 [Home page](#)
  - 2.2 [Community home page](#)
  - 2.3 [Collection home page](#)
  - 2.4 [Item home page](#)
- 3 [Usage Event Logging and Usage Statistics Gathering](#)
- 4 [Configuration settings for Statistics](#)
  - 4.1 [Upgrade Process for Statistics.](#)
- 5 [Older setting that are not related to the new 1.6 Statistics](#)
- 6 [Statistics Administration](#)
  - 6.1 [Converting older DSpace logs into SOLR usage data](#)
  - 6.2 [Statistics Client Utility](#)
- 7 [Statistics differences between DSpace 1.6.x and 1.7.0](#)
  - 7.1 [SOLR optimization added](#)
  - 7.2 [SOLR Autocommit](#)
- 8 [Custom Reporting](#)
  - 8.1 [Resources](#)
  - 8.2 [Examples](#)
    - 8.2.1 [Top downloaded items by a specific user](#)

## What is exactly being logged ?

Each time a page or file gets requested, this request is being logged. The logging happens at the server side, and doesn't require a javascript like Google Analytics does, to provide usage data.

Definition of which fields are to be stored happens in the file `dspace/solr/statistics/conf/schema.xml`.

The fields, stored in a usage event by default are:

```
<field name="type" type="integer" indexed="true" stored="true" required="true" />
<field name="id" type="integer" indexed="true" stored="true" required="true" />
<field name="ip" type="string" indexed="true" stored="true" required="false" />
<field name="time" type="date" indexed="true" stored="true" required="true" />
<field name="epersonid" type="integer" indexed="true" stored="true" required="false" />
<field name="continent" type="string" indexed="true" stored="true" required="false"/>
<field name="country" type="string" indexed="true" stored="true" required="false"/>
<field name="countryCode" type="string" indexed="true" stored="true" required="false"/>
<field name="city" type="string" indexed="true" stored="true" required="false"/>
<field name="longitude" type="float" indexed="true" stored="true" required="false"/>
<field name="latitude" type="float" indexed="true" stored="true" required="false"/>
<field name="owningComm" type="integer" indexed="true" stored="true" required="false" multiValued="true"/>
<field name="owningColl" type="integer" indexed="true" stored="true" required="false" multiValued="true"/>
<field name="owningItem" type="integer" indexed="true" stored="true" required="false" multiValued="true"/>
<field name="dns" type="string" indexed="true" stored="true" required="false"/>
<field name="userAgent" type="string" indexed="true" stored="true" required="false"/>
<field name="isBot" type="boolean" indexed="true" stored="true" required="false"/>
```

The combination of [type](#) and `id` determine which resource (either community, collection, item page or file download) has been requested.

## Web user interface for DSpace statistics

In the XMLUI, statistics can be accessed from the lower end of the navigation menu. In the JSPUI, a view statistics button appears on the bottom of pages for which statistics are available.

If you are not seeing these links or buttons, it's likely that they are only enabled for administrators in your installation. Change the configuration parameter `"statistics.item.authorization.admin"` to false in order to make statistics visible for all repository visitors.

## Home page

Starting from the repository homepage, the statistics page displays the top 10 most popular items of the entire repository.

## Community home page

The following statistics are available for the community home pages:

- Total visits of the current community home page
- Visits of the community home page over a timespan of the last 7 months
- Top 10 country from where the visits originate
- Top 10 cities from where the visits originate

## Collection home page

The following statistics are available for the collection home pages:

- Total visits of the current collection home page
- Visits of the collection home over a timespan of the last 7 months
- Top 10 country from where the visits originate
- Top 10 cities from where the visits originate

## Item home page

The following statistics are available for the item home pages:

- Total visits of the item
- Total visits for the bitstreams attached to the item
- Visits of the item over a timespan of the last 7 months
- Top 10 country views from where the visits originate
- Top 10 cities from where the visits originate

## Usage Event Logging and Usage Statistics Gathering

The DSpace Statistics Implementation is a Client/Server architecture based on Solr for collecting usage events in the JSPUI and XMLUI user interface applications of DSpace. Solr runs as a separate webapplication and an instance of Apache Http Client is utilized to allow parallel requests to log statistics events into this Solr instance.

## Configuration settings for Statistics

In the dspace.cfg file review the following fields to make sure they are uncommented:

Property:	<code>solr.log.server</code>
Example Value:	<code>solr.log.server = <a href="http://127.0.0.1/solr/statistics">http://127.0.0.1/solr/statistics</a></code>
Informational Note:	<p>Is used by the SolrLogger Client class to connect to the Solr server over http and perform updates and queries. In most cases, this can (and should) be set to localhost (or 127.0.0.1).</p> <p>To determine the correct path, you can use a tool like <code>wget</code> to see where Solr is responding on your server. For example, you'd want to send a query to Solr like the following:</p> <pre>wget http://127.0.0.1/solr/statistics/select?q=*:*</pre> <p>Assuming you get an HTTP 200 OK response, then you should set <code>solr.log.server</code> to the '/statistics' URL of 'http://127.0.0.1/solr/statistics' (essentially removing the "/select?q=:" query off the end of the responding URL.)</p>
Property:	<code>solr.spiderips.urls</code>

Example Value:	<p><code>solr.spiderips.urls =</code></p> <pre> http://iplists.com/google.txt, \ http://iplists.com/inktomi.txt, \ http://iplists.com/lycos.txt, \ http://iplists.com/infoseek.txt, \ http://iplists.com/altavista.txt, \ http://iplists.com/excite.txt, \ http://iplists.com/misc.txt, \ http://iplists.com/non_engines.txt </pre>
Informational Note:	<p>List of URLs to download spiders files into [dspace]/config/spiders. These files contain lists of known spider IPs and are utilized by the SolrLogger to flag usage events with an "isBot" field, or ignore them entirely.</p> <p>The "stats-util" command can be used to force an update of spider files, regenerate "isBot" fields on indexed events, and delete spiders from the index. For usage, run:</p> <pre>dspace stats-util -h</pre> <p>from your [dspace]/bin directory</p>
Property:	<code>solr.dbfile</code>
Example Value:	<code>solr.dbfile = \${dspace.dir}/config/GeoLiteCity.dat</code>
Informational Note:	The following refers to the GeoLiteCity database file utilized by the LocationUtils to calculate the location of client requests based on IP address. During the Ant build process (both fresh_install and update) this file will be downloaded from <a href="http://www.maxmind.com/app/geolitecity">http://www.maxmind.com/app/geolitecity</a> if a new version has been published or it is absent from your [dspace]/config directory.
Property:	<code>solr.resolver.timeout</code>
Example Value:	<code>solr.resolver.timeout = 200</code>
Informational Note:	Timeout in milliseconds for DNS resolution of origin hosts/IPs. Setting this value too high may result in solr exhausting your connection pool.
Property:	<code>useProxies</code>
Example Value:	<code>useProxies = true</code>
Informational Note:	Will cause Statistics logging to look for X-Forward URI to detect clients IP that have accessed it through a Proxy service (e.g. the Apache mod_proxy). Allows detection of client IP when accessing DSpace. [Note: This setting is found in the DSpace Logging section of dspace.cfg]
Property:	<code>statistics.item.authorization.admin</code>
Example Value:	<code>statistics.item.authorization.admin = true</code>
Informational Note:	When set to true, only general administrators, collection and community administrators are able to access the statistics from the web user interface. As a result, the links to access statistics are hidden for non logged-in admin users. Setting this property to "false" will display the links to access statistics to anyone, making them publicly available.
Property:	<code>solr.statistics.logBots</code>
Example Value:	<code>solr.statistics.logBots = true</code>
Informational Note:	<p>When this property is set to false, and IP is detected as a spider, the event is not logged.</p> <p>When this property is set to true, the event will be logged with the "isBot" field set to true.</p> <p>(see <code>solr.statistics.query.filter.*</code> for query filter options)</p>

Property:	solr.statistics.query.filter.spiderIp
Example Value:	solr.statistics.query.filter.spiderIp = false
Informational Note:	If true, statistics queries will filter out spider IPs -- use with caution, as this often results in extremely long query strings.
Property:	solr.statistics.query.filter.isBot
Example Value:	solr.statistics.query.filter.isBot = true
Informational Note:	If true, statistics queries will filter out events flagged with the "isBot" field. This is the recommended method of filtering spiders from statistics.

## Upgrade Process for Statistics.

Example of rebuild and redeploy DSpace (only if you have configured your distribution in this manner)

First approach the traditional DSpace build process for updating

```
cd [dspace-source]/dspace
mvn package
cd [dspace-source]/dspace/target/dspace-<version>-build.dir
ant -Dconfig=[dspace]/config/dspace.cfg update
cp -R [dspace]/webapps/* [TOMCAT]/webapps
```

The last step is only used if you are not mounting *[dspace]/webapps* directly into your Tomcat, Resin or Jetty host (the recommended practice) If you only need to build the statistics, and don't make any changes to other web applications, you can replace the copy step above with:

```
cp -R dspace/webapps/solr TOMCAT/webapps
```

*Again, only if you are not mounting [dspace]/webapps directly into your Tomcat, Resin or Jetty host (the recommended practice)*

Restart your webapps (Tomcat/Jetty/Resin)

## Older setting that are not related to the new 1.6 Statistics

The following Dspace.cfg fields are only applicable to the older statistics solution.

```
##### Statistical Report Configuration Settings #####

# should the stats be publicly available?  should be set to false if you only
# want administrators to access the stats, or you do not intend to generate
# any
report.public = false

# directory where live reports are stored
report.dir = ${dspace.dir}/reports/
```

These fields are not used by the new 1.6 Statistics, but are only related to the Statistics from previous DSpace releases

## Statistics Administration

### Converting older DSpace logs into SOLR usage data

If you have upgraded from a previous version of DSpace, converting older log files ensures that you carry over older usage stats from before the upgrade.

### Statistics Client Utility

The command line interface (CLI) scripts can be used to clean the usage database from additional spider traffic and other maintenance tasks.

## Statistics differences between DSpace 1.6.x and 1.7.0

### SOLR optimization added

If required, the solr server can be optimized by running

```
{dspace.dir}/bin/stats-util -o
```

. More information on how these solr server optimizations work can be found here: [http://wiki.apache.org/solr/SolrPerformanceFactors#Optimization\\_Considerations](http://wiki.apache.org/solr/SolrPerformanceFactors#Optimization_Considerations).

### SOLR Autocommit

In DSpace 1.6.x, each solr event was committed to the solr server individually. For high load DSpace installations, this would result in a huge load of small solr commits resulting in a very high load on the solr server.

This has been resolved in dspace 1.7 by only committing usage events to the solr server every 15 minutes. This will result in a delay of the storage of a usage event of maximum 15 minutes. If required, this value can be altered by changing the maxTime property in the

```
{dspace.dir}/solr/statistics/conf/solrconfig.xml.
```

## Custom Reporting

When the web user interface does not offer you the statistics you need, you can greatly expand the reports by querying the SOLR index directly.

### Resources

- <http://www.lucidimagination.com/Community/Hear-from-the-Experts/Articles/Faceted-Search-Solr>
- <http://my.safaribooksonline.com/9781847195883/Cover>

### Examples

#### Top downloaded items by a specific user

Query:

```
http://localhost:8080/solr/statistics/select?indent=on&version=2.2&start=0&rows=10&fl=*&2Cscore&qt=standard&wt=standard&explainOther=&hl.fl=&facet=true&facet.field=epersonid&q=type:0
```

Explained:

facet.field=epersonid — You want to group by epersonid, which is the user id.  
type:0 — Interested in bitstreams only

```
<lst name="facet_counts">
  <lst name="facet_fields">
    <lst name="epersonid">
      <int name="66">1167</int>

<int name="117">251</int>

<int name="52">42</int>

<int name="19">36</int>

<int name="88">20</int>

<int name="112">18</int>

<int name="110">9</int>

<int name="96">0</int>

</lst>
  </lst>
</lst>
```