

Business Logic Layer

- 1 Core Classes
 - 1.1 The Configuration Manager
 - 1.2 Constants
 - 1.3 Context
 - 1.4 Email
 - 1.5 LogManager
 - 1.6 Utils
- 2 Content Management API
 - 2.1 Other Classes
 - 2.2 Modifications
 - 2.3 What's In Memory?
 - 2.4 Dublin Core Metadata
 - 2.5 Support for Other Metadata Schemas
 - 2.6 Packager Plugins
- 3 Plugin Manager
 - 3.1 Concepts
 - 3.2 Using the Plugin Manager
 - 3.2.1 Types of Plugin
 - 3.2.2 Self-Named Plugins
 - 3.2.3 Obtaining a Plugin Instance
 - 3.2.4 Lifecycle Management
 - 3.2.5 Getting Meta-Information
 - 3.3 Implementation
 - 3.3.1 PluginManager Class
 - 3.3.2 SelfNamedPlugin Class
 - 3.3.3 Errors and Exceptions
 - 3.4 Configuring Plugins
 - 3.4.1 Configuring Singleton (Single) Plugins
 - 3.4.2 Configuring Sequence of Plugins
 - 3.4.3 Configuring Named Plugins
 - 3.4.4 Configuring the Reusable Status of a Plugin
 - 3.5 Validating the Configuration
 - 3.6 Use Cases
 - 3.6.1 Managing the MediaFilter plugins transparently
 - 3.6.2 A Singleton Plugin
 - 3.6.3 Plugin that Names Itself
 - 3.6.4 Stackable Authentication
- 4 Workflow System
- 5 Administration Toolkit
- 6 E-person/Group Manager
- 7 Authorization
 - 7.1 Special Groups
 - 7.2 Miscellaneous Authorization Notes
- 8 Handle Manager/Handle Plugin
- 9 Search
 - 9.1 Current Lucene Implementation
 - 9.2 Indexed Fields
 - 9.3 Harvesting API
- 10 Browse API
 - 10.1 Using the API
 - 10.2 Index Maintenance
 - 10.3 Caveats
- 11 Checksum checker
- 12 OpenSearch Support
- 13 Embargo Support
 - 13.1 What is an Embargo?
 - 13.2 Embargo Model and Life-Cycle

Core Classes

The *org.dspace.core* package provides some basic classes that are used throughout the DSpace code.

The Configuration Manager

The configuration manager is responsible for reading the main *dspace.cfg* properties file, managing the 'template' configuration files for other applications such as Apache, and for obtaining the text for e-mail messages.

The system is configured by editing the relevant files in `[dspace]/config`, as described in the configuration section.

When editing configuration files for applications that DSpace uses, such as Apache Tomcat, you may want to edit the copy in `[dspace-source]` and then run `ant update` or `ant overwrite_configs` rather than editing the 'live' version directly! This will ensure you have a backup copy of your modified configuration files, so that they are not accidentally overwritten in the future.

The *ConfigurationManager* class can also be invoked as a command line tool:

- `[dspace]/bin/dspace dsprop property.name` This writes the value of *property.name* from *dspace.cfg* to the standard output, so that shell scripts can access the DSpace configuration. If the property has no value, nothing is written.

Constants

This class contains constants that are used to represent types of object and actions in the database. For example, authorization policies can relate to objects of different types, so the *resourcepolicy* table has columns *resource_id*, which is the internal ID of the object, and *resource_type_id*, which indicates whether the object is an item, collection, bitstream etc. The value of *resource_type_id* is taken from the *Constants* class, for example *Constants.ITEM*.

Here are a some of the most commonly used constants you might come across:

DSpace types

- Bitstream: 0
- Bundle: 1
- Item: 2
- Collection: 3
- Community: 4
- Site: 5
- Group: 6
- Eperson: 7

DSpace actions

- Read: 0
- Write: 1
- Delete: 2
- Add: 3
- Remove: 4

Refer to the org.dspace.core.Constants for all of the Constants.

Context

The *Context* class is central to the DSpace operation. Any code that wishes to use the any API in the business logic layer must first create itself a *Context* object. This is akin to opening a connection to a database (which is in fact one of the things that happens.)

A context object is involved in most method calls and object constructors, so that the method or object has access to information about the current operation. When the context object is constructed, the following information is automatically initialized:

- A connection to the database. This is a transaction-safe connection. i.e. the 'auto-commit' flag is set to false.
- A cache of content management API objects. Each time a content object is created (for example *Item* or *Bitstream*) it is stored in the *Context* object. If the object is then requested again, the cached copy is used. Apart from reducing database use, this addresses the problem of having two copies of the same object in memory in different states.
The following information is also held in a context object, though it is the responsibility of the application creating the context object to fill it out correctly:
- The current authenticated user, if any
- Any 'special groups' the user is a member of. For example, a user might automatically be part of a particular group based on the IP address they are accessing DSpace from, even though they don't have an e-person record. Such a group is called a 'special group'.
- Any extra information from the application layer that should be added to log messages that are written within this context. For example, the Web UI adds a session ID, so that when the logs are analyzed the actions of a particular user in a particular session can be tracked.
- A flag indicating whether authorization should be circumvented. This should only be used in rare, specific circumstances. For example, when first installing the system, there are no authorized administrators who would be able to create an administrator account!As noted above, the public API is *trusted*, so it is up to applications in the application layer to use this flag responsibly.
Typical use of the context object will involve constructing one, and setting the current user if one is authenticated. Several operations may be performed using the context object. If all goes well, *complete* is called to commit the changes and free up any resources used by the context. If anything has gone wrong, *abort* is called to roll back any changes and free up the resources.

You should always *abort* a context if *any* error happens during its lifespan; otherwise the data in the system may be left in an inconsistent state. You can also *commit* a context, which means that any changes are written to the database, and the context is kept active for further use.

Email

Sending e-mails is pretty easy. Just use the configuration manager's *getEmail* method, set the arguments and recipients, and send.

The e-mail texts are stored in `[dspace]/config/emails`. They are processed by the standard *java.text.MessageFormat*. At the top of each e-mail are listed the appropriate arguments that should be filled out by the sender. Example usage is shown in the *org.dspace.core.Email* Javadoc API documentation.

LogManager

The log manager consists of a method that creates a standard log header, and returns it as a string suitable for logging. Note that this class does not actually write anything to the logs; the log header returned should be logged directly by the sender using an appropriate Log4J call, so that information about where the logging is taking place is also stored.

The level of logging can be configured on a per-package or per-class basis by editing `[dspace]/config/log4j.properties`. You will need to stop and restart Tomcat for the changes to take effect.

A typical log entry looks like this:

```
2002-11-11 08:11:32,903 INFO org.dspace.app.webui.servlet.DSpaceServlet @ anonymous:session_id=BD84E7C194C2CF4BD0EC3A6CAD0142BB:
view_item:handle=172.1.1/1686
```

This is breaks down like this:

Date and time, milliseconds	2002-11-11 08:11:32,903
Level (<i>FATAL</i> , <i>WARN</i> , <i>INFO</i> or <i>DEBUG</i>)	<i>INFO</i>
Java class	<i>org.dspace.app.webui.servlet.DSpaceServlet</i>
	@
User email or <i>anonymous</i>	<i>anonymous</i>
	:
Extra log info from context	<i>session_id=BD84E7C194C2CF4BD0EC3A6CAD0142BB</i>
	:
Action	<i>view_item</i>
	:
Extra info	<i>handle=172.1.1/1686</i>

The above format allows the logs to be easily parsed and analyzed. The `[dspace]/bin/log-reporter` script is a simple tool for analyzing logs. Try:

```
[dspace]/bin/log-reporter --help
```

It's a good idea to 'nice' this log reporter to avoid an impact on server performance.

Utils

Utils contains miscellaneous utility method that are required in a variety of places throughout the code, and thus have no particular 'home' in a subsystem.

Content Management API

The content management API package *org.dspace.content* contains Java classes for reading and manipulating content stored in the DSpace system. This is the API that components in the application layer will probably use most.

Classes corresponding to the main elements in the DSpace data model (*Community*, *Collection*, *Item*, *Bundle* and *Bitstream*) are sub-classes of the abstract class *DSpaceObject*. The *Item* object handles the Dublin Core metadata record.

Each class generally has one or more static *find* methods, which are used to instantiate content objects. Constructors do not have public access and are just used internally. The reasons for this are:

- "Constructing" an object may be misconstrued as the action of creating an object in the DSpace system, for example one might expect something like:

```
Context dsContent = new Context();
Item myItem = new Item(context, id)
```

- to construct a brand new item in the system, rather than simply instantiating an in-memory instance of an object in the system.
- *find* methods may often be called with invalid IDs, and return *null* in such a case. A constructor would have to throw an exception in this case. A *null* return value from a static method can in general be dealt with more simply in code.
- If an instantiation representing the same underlying archival entity already exists, the *find* method can simply return that same instantiation to avoid multiple copies and any inconsistencies which might result.

Collection, *Bundle* and *Bitstream* do not have *create* methods; rather, one has to create an object using the relevant method on the container. For example, to create a collection, one must invoke *createCollection* on the community that the collection is to appear in:

```
Context context = new Context();
Community existingCommunity = Community.find(context, 123);
Collection myNewCollection = existingCommunity.createCollection();
```

The primary reason for this is for determining authorization. In order to know whether an e-person may create an object, the system must know which container the object is to be added to. It makes no sense to create a collection outside of a community, and the authorization system does not have a policy for that.

Item_s are first created in the form of an implementation of `_InProgressSubmission`. An `InProgressSubmission` represents an item under construction; once it is complete, it is installed into the main archive and added to the relevant collection by the `InstallItem` class. The `org.dspace.content` package provides an implementation of `InProgressSubmission` called `WorkspaceItem`; this is a simple implementation that contains some fields used by the Web submission UI. The `org.dspace.workflow` also contains an implementation called `WorkflowItem` which represents a submission undergoing a workflow process.

In the previous chapter there is an overview of the item ingest process which should clarify the previous paragraph. Also see the section on the workflow system.

`Community` and `BitstreamFormat` do have static `create` methods; one must be a site administrator to have authorization to invoke these.

Other Classes

Classes whose name begins `DC` are for manipulating Dublin Core metadata, as explained below.

The `FormatIdentifier` class attempts to guess the bitstream format of a particular bitstream. Presently, it does this simply by looking at any file extension in the bitstream name and matching it up with the file extensions associated with bitstream formats. Hopefully this can be greatly improved in the future!

The `ItemIterator` class allows items to be retrieved from storage one at a time, and is returned by methods that may return a large number of items, more than would be desirable to have in memory at once.

The `ItemComparator` class is an implementation of the standard `java.util.Comparator` that can be used to compare and order items based on a particular Dublin Core metadata field.

Modifications

When creating, modifying or for whatever reason removing data with the content management API, it is important to know when changes happen in-memory, and when they occur in the physical DSpace storage.

Primarily, one should note that no change made using a particular `org.dspace.core.Context` object will actually be made in the underlying storage unless `complete` or `commit` is invoked on that `Context`. If anything should go wrong during an operation, the context should always be aborted by invoking `abort`, to ensure that no inconsistent state is written to the storage.

Additionally, some changes made to objects only happen in-memory. In these cases, invoking the `update` method lines up the in-memory changes to occur in storage when the `Context` is committed or completed. In general, methods that change any metadata field only make the change in-memory; methods that involve relationships with other objects in the system line up the changes to be committed with the context. See individual methods in the API Javadoc.

Some examples to illustrate this are shown below:

<pre>Context context = new Context(); Bitstream b = Bitstream.find(context, 1234); b.setName("newfile.txt"); b.update(); context.complete();</pre>	<p>Will change storage</p>
<pre>Context context = new Context(); Bitstream b = Bitstream.find(context, 1234); b.setName("newfile.txt"); b.update(); context.abort();</pre>	<p>Will not change storage (context aborted)</p>
<pre>Context context = new Context(); Bitstream b = Bitstream.find(context, 1234); b.setName("newfile.txt"); context.complete();</pre>	<p>The new name will not be stored since <code>update</code> was not invoked</p>

```
Context context = new Context();
Bitstream bs = Bitstream.find(context, 1234);
Bundle bnd = Bundle.find(context, 5678);
bnd.add(bs);
context.complete();
```

The bitstream **will** be included in the bundle, since *update* doesn't need to be called

What's In Memory?

Instantiating some content objects also causes other content objects to be loaded into memory.

Instantiating a *Bitstream* object causes the appropriate *BitstreamFormat* object to be instantiated. Of course the *Bitstream* object does not load the underlying bits from the bitstream store into memory!

Instantiating a *Bundle* object causes the appropriate *Bitstream* objects (and hence *_BitstreamFormat_s*) to be instantiated.

Instantiating an *Item* object causes the appropriate *Bundle* objects (etc.) and hence *_BitstreamFormat_s* to be instantiated. All the Dublin Core metadata associated with that item are also loaded into memory.

The reasoning behind this is that for the vast majority of cases, anyone instantiating an item object is going to need information about the bundles and bitstreams within it, and this methodology allows that to be done in the most efficient way and is simple for the caller. For example, in the Web UI, the servlet (controller) needs to pass information about an item to the viewer (JSP), which needs to have all the information in-memory to display the item without further accesses to the database which may cause errors mid-display.

You do not need to worry about multiple in-memory instantiations of the same object, or any inconsistencies that may result; the *Context* object keeps a cache of the instantiated objects. The *find* methods of classes in *org.dspace.content* will use a cached object if one exists.

It may be that in enough cases this automatic instantiation of contained objects reduces performance in situations where it is important; if this proves to be true the API may be changed in the future to include a *loadContents* method or somesuch, or perhaps a Boolean parameter indicating what to do will be added to the *find* methods.

When a *Context* object is completed, aborted or garbage-collected, any objects instantiated using that context are invalidated and should not be used (in much the same way an AWT button is invalid if the window containing it is destroyed).

Dublin Core Metadata

The *Metadatum* class is a simple container that represents a single Dublin Core-like element, optional qualifier, value and language. Note that since DSpace 1.4 the *MetadataValue* and associated classes are preferred (see Support for Other Metadata Schemas). The other classes starting with *DC* are utility classes for handling types of data in Dublin Core, such as people's names and dates. As supplied, the DSpace registry of elements and qualifiers corresponds to the [Library Application Profile](#) for Dublin Core. It should be noted that these utility classes assume that the values will be in a certain syntax, which will be true for all data generated within the DSpace system, but since Dublin Core does not always define strict syntax, this may not be true for Dublin Core originating outside DSpace.

Below is the specific syntax that DSpace expects various fields to adhere to:

Element	Qualifier	Syntax	Helper Class
<i>date</i>	Any or unqualified	ISO 8601 in the UTC time zone, with either year, month, day, or second precision. Examples: <i>_2000 2002-10 2002-08-14 1999-01-01T14:35:23Z _</i>	<i>DCDate</i>
<i>contributor</i>	Any or unqualified	In general last name, then a comma, then first names, then any additional information like "Jr.". If the contributor is an organization, then simply the name. Examples: <i>_Doe, John Smith, John Jr. van Dyke, Dick Massachusetts Institute of Technology _</i>	<i>DCPersonName</i>
<i>language</i>	<i>iso</i>	A two letter code taken ISO 639, followed optionally by a two letter country code taken from ISO 3166. Examples: <i>_en fr en_US _</i>	<i>DCLanguage</i>
<i>relation</i>	<i>ispartofseries</i>	The series name, following by a semicolon followed by the number in that series. Alternatively, just free text. <i>_MIT-TR; 1234 My Report Series; ABC-1234 NS1234 _</i>	<i>DCSeriesNumber</i>

Support for Other Metadata Schemas

To support additional metadata schemas a new set of metadata classes have been added. These are backwards compatible with the DC classes and should be used rather than the DC specific classes wherever possible. Note that hierarchical metadata schemas are not currently supported, only flat schemas (such as DC) are able to be defined.

The *MetadataField* class describes a metadata field by schema, element and optional qualifier. The value of a *MetadataField* is described by a *MetadataValue* which is roughly equivalent to the older *Metadatum* class. Finally the *MetadataSchema* class is used to describe supported schemas. The DC schema is supported by default. Refer to the javadoc for method details.

Packager Plugins

The Packager plugins let you *ingest* a package to create a new DSpace Object, and *disseminate* a content Object as a package. A package is simply a data stream; its contents are defined by the packager plugin's implementation.

To ingest an object, which is currently only implemented for Items, the sequence of operations is:

1. Get an instance of the chosen *PackageIngester* plugin.
2. Locate a Collection in which to create the new Item.
3. Call its *ingest* method, and get back a *WorkspaceItem*.
The packager also takes a *PackageParameters* object, which is a property list of parameters specific to that packager which might be passed in from the user interface.

Here is an example package ingestion code fragment:

```
Collection collection = find target collection
    InputStream source = ...;
    PackageParameters params = ...;
    String license = null;

PackageIngester sip = (PackageIngester) PluginManager
    .getNamedPlugin(PackageIngester.class, packageType);

WorkspaceItem wi = sip.ingest(context, collection, source, params, license);
```

Here is an example of a package dissemination:

```
OutputStream destination = ...;
    PackageParameters params = ...;
    DSpaceObject dso = ...;

PackageIngester dip = (PackageDisseminator) PluginManager
    .getNamedPlugin(PackageDisseminator.class, packageType);

dip.disseminate(context, dso, params, destination);
```

Plugin Manager

The PluginManager is a very simple component container. It creates and organizes components (plugins), and helps select a plugin in the cases where there are many possible choices. It also gives some limited control over the life cycle of a plugin.

Concepts

The following terms are important in understanding the rest of this section:

- **Plugin Interface** A Java interface, the defining characteristic of a plugin. The consumer of a plugin asks for its plugin by interface.
- **Plugin** a.k.a. Component, this is an instance of a class that implements a certain interface. It is interchangeable with other implementations, so that any of them may be "plugged in", hence the name. A Plugin is an instance of any class that implements the plugin interface.
- **Implementation class** The actual class of a plugin. It may implement several plugin interfaces, but must implement at least one.
- **Name** Plugin implementations can be distinguished from each other by name, a short String meant to symbolically represent the implementation class. They are called "named plugins". Plugins only need to be named when the caller has to make an active choice between them.
- **SelfNamedPlugin class** Plugins that extend the *SelfNamedPlugin* class can take advantage of additional features of the Plugin Manager. Any class can be managed as a plugin, so it is not necessary, just possible.
- **Reusable** Reusable plugins are only instantiated once, and the Plugin Manager returns the same (cached) instance whenever that same plugin is requested again. This behavior can be turned off if desired.

Using the Plugin Manager

Types of Plugin

The Plugin Manager supports three different patterns of usage:

1. **Singleton Plugins** There is only one implementation class for the plugin. It is indicated in the configuration. This type of plugin chooses an implementation of a service, for the entire system, at configuration time. Your application just fetches the plugin for that interface and gets the configured-in choice. See the `getSinglePlugin()` method.
2. **Sequence Plugins** You need a sequence or series of plugins, to implement a mechanism like Stackable Authentication or a pipeline, where each plugin is called in order to contribute its implementation of a process to the whole. The Plugin Manager supports this by letting you configure a sequence of plugins for a given interface. See the `getPluginSequence()` method.
3. **Named Plugins** Use a named plugin when the application has to choose one plugin implementation out of many available ones. Each implementation is bound to one or more names (symbolic identifiers) in the configuration. The name is just a string to be associated with the combination of implementation class and interface. It may contain any characters except for comma (,) and equals (=). It may contain embedded

spaces. Comma is a special character used to separate names in the configuration entry. Names must be unique within an interface: No plugin classes implementing the same interface may have the same name. Think of plugin names as a controlled vocabulary – for a given plugin interface, there is a set of names for which plugins can be found. The designer of a Named Plugin interface is responsible for deciding what the name means and how to derive it; for example, names of metadata crosswalk plugins may describe the target metadata format. See the `getNamedPlugin()` method and the `getPluginNames()` methods.

Self-Named Plugins

Named plugins can get their names either from the configuration or, for a variant called self-named plugins, from within the plugin itself.

Self-named plugins are necessary because one plugin implementation can be configured itself to take on many "personalities", each of which deserves its own plugin name. It is already managing its own configuration for each of these personalities, so it makes sense to allow it to export them to the Plugin Manager rather than expecting the plugin configuration to be kept in sync with its own configuration.

An example helps clarify the point: There is a named plugin that does crosswalks, call it *CrosswalkPlugin*. It has several implementations that crosswalk some kind of metadata. Now we add a new plugin which uses XSL stylesheet transformation (XSLT) to crosswalk many types of metadata – so the single plugin can act like many different plugins, depending on which stylesheet it employs.

This XSLT-crosswalk plugin has its own configuration that maps a Plugin Name to a stylesheet – it has to, since of course the Plugin Manager doesn't know anything about stylesheets. It becomes a self-named plugin, so that it reads its configuration data, gets the list of names to which it can respond, and passes those on to the Plugin Manager.

When the Plugin Manager creates an instance of the XSLT-crosswalk, it records the Plugin Name that was responsible for that instance. The plugin can look at that Name later in order to configure itself correctly for the Name that created it. This mechanism is all part of the `SelfNamedPlugin` class which is part of any self-named plugin.

Obtaining a Plugin Instance

The most common thing you will do with the Plugin Manager is obtain an instance of a plugin. To request a plugin, you must always specify the plugin interface you want. You will also supply a name when asking for a named plugin.

A sequence plugin is returned as an array of `_Object_s` since it is actually an ordered list of plugins.

See the `getSinglePlugin()`, `getPluginSequence()`, `getNamedPlugin()` methods.

Lifecycle Management

When *PluginManager* fulfills a request for a plugin, it checks whether the implementation class is reusable; if so, it creates one instance of that class and returns it for every subsequent request for that interface and name. If it is not reusable, a new instance is always created.

For reasons that will become clear later, the manager actually caches a separate instance of an implementation class for each name under which it can be requested.

You can ask the *PluginManager* to forget about (decache) a plugin instance, by releasing it. See the `PluginManager.releasePlugin()` method. The manager will drop its reference to the plugin so the garbage collector can reclaim it. The next time that plugin/name combination is requested, it will create a new instance.

Getting Meta-Information

The *PluginManager* can list all the names of the Named Plugins which implement an interface. You may need this, for example, to implement a menu in a user interface that presents a choice among all possible plugins. See the `getPluginNames()` method.

Note that it only returns the plugin name, so if you need a more sophisticated or meaningful "label" (i.e. a key into the I18N message catalog) then you should add a method to the plugin itself to return that.

Implementation

Note: The *PluginManager* refers to interfaces and classes internally only by their names whenever possible, to avoid loading classes until absolutely necessary (i.e. to create an instance). As you'll see below, self-named classes still have to be loaded to query them for names, but for the most part it can avoid loading classes. This saves a lot of time at start-up and keeps the JVM memory footprint down, too. As the Plugin Manager gets used for more classes, this will become a greater concern.

The only downside of "on-demand" loading is that errors in the configuration don't get discovered right away. The solution is to call the `checkConfiguration()` method after making any changes to the configuration.

PluginManager Class

The *PluginManager* class is your main interface to the Plugin Manager. It behaves like a factory class that never gets instantiated, so its public methods are static.

Here are the public methods, followed by explanations:

- ```
static Object getSinglePlugin(Class interface)
 throws PluginConfigurationException;
```

Returns an instance of the singleton (single) plugin implementing the given interface. There must be exactly one single plugin configured for this interface, otherwise the *PluginConfigurationError* is thrown. Note that this is the only "get plugin" method which throws an exception. It is typically used at initialization time to set up a permanent part of the system so any failure is fatal. See the *plugin.single* configuration key for configuration details.

```
• static Object[] getPluginSequence(Class intface);
```

Returns instances of all plugins that implement the interface *intface*, in an *Array*. Returns an empty array if there are no matching plugins. The order of the plugins in the array is the same as their class names in the configuration's value field. See the *plugin.sequence* configuration key for configuration details.

```
• static Object getNamedPlugin(Class intface, String name);
```

Returns an instance of a plugin that implements the interface *intface* and is bound to a name matching name. If there is no matching plugin, it returns null. The names are matched by *String.equals()*. See the *plugin.named* and *plugin.selfnamed* configuration keys for configuration details.

```
• static void releasePlugin(Object plugin);
```

Tells the Plugin Manager to let go of any references to a reusable plugin, to prevent it from being given out again and to allow the object to be garbage-collected. Call this when a plugin instance must be taken out of circulation.

```
• static String[] getAllPluginNames(Class intface);
```

Returns all of the names under which a named plugin implementing the interface *intface* can be requested (with *getNamedPlugin()*). The array is empty if there are no matches. Use this to populate a menu of plugins for interactive selection, or to document what the possible choices are. The names are NOT returned in any predictable order, so you may wish to sort them first. Note: Since a plugin may be bound to more than one name, the list of names this returns does not represent the list of plugins. To get the list of unique implementation classes corresponding to the names, you might have to eliminate duplicates (i.e. create a Set of classes).

```
• static void checkConfiguration();
```

Validates the keys in the DSpace *ConfigurationManager* pertaining to the Plugin Manager and reports any errors by logging them. This is intended to be used interactively by a DSpace administrator, to check the configuration file after modifying it. See the section about validating configuration for details.

## SelfNamedPlugin Class

A named plugin implementation must extend this class if it wants to supply its own Plugin Name(s). See Self-Named Plugins for why this is sometimes necessary.

```
abstract class SelfNamedPlugin
{
 // Your class must override this:
 // Return all names by which this plugin should be known.
 public static String[] getPluginNames();

 // Returns the name under which this instance was created.
 // This is implemented by SelfNamedPlugin and should NOT be
 // overridden.
 public String getPluginInstanceName();
}
```

## Errors and Exceptions

```
public class PluginConfigurationError extends Error
{
 public PluginConfigurationError(String message);
}
```

An error of this type means the caller asked for a single plugin, but either there was no single plugin configured matching that interface, or there was more than one. Either case causes a fatal configuration error.

```
public class PluginInstantiationException extends RuntimeException
{
 public PluginInstantiationException(String msg, Throwable cause)
}
}
```

This exception indicates a fatal error when instantiating a plugin class. It should only be thrown when something unexpected happens in the course of instantiating a plugin, e.g. an access error, class not found, etc. Simply not finding a class in the configuration is not an exception.

This is a *RuntimeException* so it doesn't have to be declared, and can be passed all the way up to a generalized fatal exception handler.

## Configuring Plugins

All of the Plugin Manager's configuration comes from the DSpace Configuration Manager, which is a Java Properties map. You can configure these characteristics of each plugin:

1. **Interface:** Classname of the Java interface which defines the plugin, including package name. e.g. *org.dspace.app.mediafilter.FormatFilter*
2. **Implementation Class:** Classname of the implementation class, including package. e.g. *org.dspace.app.mediafilter.PDFFilter*
3. **Names:** (Named plugins only) There are two ways to bind names to plugins: listing them in the value of a *plugin.named.interface* key, or configuring a class in *plugin.selfnamed.interface* which extends the *SelfNamedPlugin* class.
4. **Reusable option:** (Optional) This is declared in a *plugin.reusable* configuration line. Plugins are reusable by default, so you only need to configure the non-reusable ones.

### Configuring Singleton (Single) Plugins

This entry configures a Single Plugin for use with *getSinglePlugin()*:

```
plugin.single.interface = classname
```

For example, this configures the class *org.dspace.checker.SimpleDispatcher* as the plugin for interface *org.dspace.checker.BitstreamDispatcher*:

```
plugin.single.org.dspace.checker.BitstreamDispatcher=org.dspace.checker.SimpleDispatcher
```

### Configuring Sequence of Plugins

This kind of configuration entry defines a Sequence Plugin, which is bound to a sequence of implementation classes. The key identifies the interface, and the value is a comma-separated list of classnames:

```
plugin.sequence.interface = classname, ...
```

The plugins are returned by *getPluginSequence()* in the same order as their classes are listed in the configuration value.

For example, this entry configures Stackable Authentication with three implementation classes:

```
plugin.sequence.org.dspace.eperson.AuthenticationMethod = \
 org.dspace.eperson.X509Authentication, \
 org.dspace.eperson.PasswordAuthentication, \
 edu.mit.dspace.MITSpecialGroup
```

### Configuring Named Plugins

There are two ways of configuring named plugins:

1. **Plugins Named in the Configuration** A named plugin which gets its name(s) from the configuration is listed in this kind of entry: *\_plugin.named.interface = classname = name [ , name.. ] [ classname = name.. ]* The syntax of the configuration value is: classname, followed by an equal-sign and then at least one plugin name. Bind more names to the same implementation class by adding them here, separated by commas. Names may include any character other than comma (,) and equal-sign (=). For example, this entry creates one plugin with the names GIF, JPEG, and image/png, and another with the name TeX:

```
plugin.named.org.dspace.app.mediafilter.MediaFilter = \
 org.dspace.app.mediafilter.JPEGFilter = GIF, JPEG, image/png \
 org.dspace.app.mediafilter.TeXFilter = TeX
```

This example shows a plugin name with an embedded whitespace character. Since comma (,) is the separator character between plugin names, spaces are legal (between words of a name; leading and trailing spaces are ignored). This plugin is bound to the names "Adobe PDF", "PDF", and "Portable Document Format".

```
plugin.named.org.dspace.app.mediafilter.MediaFilter = \
 org.dspace.app.mediafilter.TeXFilter = TeX \
 org.dspace.app.mediafilter.PDFFilter = Adobe PDF, PDF, Portable Document Format
```

NOTE: Since there can only be one key with `plugin.named.` followed by the interface name in the configuration, all of the plugin implementations must be configured in that entry.

2. **Self-Named Plugins** Since a self-named plugin supplies its own names through a static method call, the configuration only has to include its interface and classname: `plugin.selfnamed.interface = classname [ , classname.. ]`. The following example first demonstrates how the plugin class, `_XsltDisseminationCrosswalk` is configured to implement its own names "MODS" and "DublinCore". These come from the keys starting with `crosswalk.dissemination.stylesheet.`. The value is a stylesheet file. The class is then configured as a self-named plugin:

```
crosswalk.dissemination.stylesheet.DublinCore = xwalk/TESTDIM-2-DC_copy.xml
crosswalk.dissemination.stylesheet.MODS = xwalk/mods.xml

plugin.selfnamed.crosswalk.org.dspace.content.metadata.DisseminationCrosswalk = \
 org.dspace.content.metadata.MODSDisseminationCrosswalk, \
 org.dspace.content.metadata.XsltDisseminationCrosswalk
```

NOTE: Since there can only be one key with `plugin.selfnamed.` followed by the interface name in the configuration, all of the plugin implementations must be configured in that entry. The `MODSDisseminationCrosswalk` class is only shown to illustrate this point.

## Configuring the Reusable Status of a Plugin

Plugins are assumed to be reusable by default, so you only need to configure the ones which you would prefer not to be reusable. The format is as follows:

```
plugin.reusable.classname = (true | false)
```

For example, this marks the PDF plugin from the example above as non-reusable:

```
plugin.reusable.org.dspace.app.mediafilter.PDFFilter = false
```

## Validating the Configuration

The Plugin Manager is very sensitive to mistakes in the DSpace configuration. Subtle errors can have unexpected consequences that are hard to detect: for example, if there are two "plugin.single" entries for the same interface, one of them will be silently ignored.

To validate the Plugin Manager configuration, call the `PluginManager.checkConfiguration()` method. It looks for the following mistakes:

- Any duplicate keys starting with "plugin."
- Keys starting `plugin.single`, `plugin.sequence`, `plugin.named`, and `plugin.selfnamed` that don't include a valid interface.
- Classnames in the configuration values that don't exist, or don't implement the plugin interface in the key.
- Classes declared in `plugin.selfnamed` lines that don't extend the `SelfNamedPlugin` class.
- Any name collisions among named plugins for a given interface.
- Named plugin configuration entries without any names.
- Classnames mentioned in `plugin.reusable` keys must exist and have been configured as a plugin implementation class.  
The `PluginManager` class also has a `main()` method which simply runs `checkConfiguration()`, so you can invoke it from the command line to test the validity of plugin configuration changes.

Eventually, someone should develop a general configuration-file sanity checker for DSpace, which would just call `PluginManager.checkConfiguration()`.

## Use Cases

Here are some usage examples to illustrate how the Plugin Manager works.

### Managing the MediaFilter plugins transparently

The existing DSpace 1.3 `MediaFilterManager` implementation has been largely replaced by the Plugin Manager. The `MediaFilter` classes become plugins named in the configuration. Refer to the configuration guide for further details.

### A Singleton Plugin

This shows how to configure and access a single anonymous plugin, such as the `BitstreamDispatcher` plugin:

Configuration:

```
plugin.single.org.dspace.checker.BitstreamDispatcher=org.dspace.checker.SimpleDispatcher
```

The following code fragment shows how dispatcher, the service object, is initialized and used:

```
BitstreamDispatcher dispatcher =
 (BitstreamDispatcher)PluginManager.getSinglePlugin(BitstreamDispatcher
.class);
int id = dispatcher.next();
while (id != BitstreamDispatcher.SENTINEL)
{
 /*
 do some processing here
 */
 id = dispatcher.next();
}
```

### Plugin that Names Itself

This crosswalk plugin acts like many different plugins since it is configured with different XSL translation stylesheets. Since it already gets each of its stylesheets out of the DSpace configuration, it makes sense to have the plugin give PluginManager the names to which it answers instead of forcing someone to configure those names in two places (and try to keep them synchronized).

NOTE: Remember how *getPlugin()* caches a separate instance of an implementation class for every name bound to it? This is why: the instance can look at the name under which it was invoked and configure itself specifically for that name. Since the instance for each name might be different, the Plugin Manager has to cache a separate instance for each name.

Here is the configuration file listing both the plugin's own configuration and the *PluginManager* config line:

```
crosswalk.dissemination.stylesheet.DublinCore = xwalk/TESTDIM-2-DC_copy.xsl
crosswalk.dissemination.stylesheet.MODS = xwalk/mods.xsl

plugin.selfnamed.org.dspace.content.metadata.DisseminationCrosswalk = \
 org.dspace.content.metadata.XsltDisseminationCrosswalk
```

This look into the implementation shows how it finds configuration entries to populate the array of plugin names returned by the *getPluginNames()* method. Also note, in the *getStyleSheet()* method, how it uses the plugin name that created the current instance (returned by *getPluginInstanceName()*) to find the correct stylesheet.

```

public class XsltDisseminationCrosswalk extends SelfNamedPlugin
{

 private final String prefix =
 "crosswalk.dissemination.stylesheet.";

 public static String[] getPluginNames()
 {
 List aliasList = new ArrayList();
 Enumeration pe = ConfigurationManager.propertyNames();

 while (pe.hasMoreElements())
 {
 String key = (String)pe.nextElement();
 if (key.startsWith(prefix))
 aliasList.add(key.substring(prefix.length()));
 }
 return (String[])aliasList.toArray(new
 String[aliasList.size()]);
 }

 // get the crosswalk stylesheet for an instance of the plugin:
 private String getStylesheet()
 {
 return ConfigurationManager.getProperty(prefix +
 getPluginInstanceName());
 }
}

```

## Stackable Authentication

The Stackable Authentication mechanism needs to know all of the plugins configured for the interface, in the order of configuration, since order is significant. It gets a Sequence Plugin from the Plugin Manager. Refer to the Configuration Section on Stackable Authentication for further details.

## Workflow System

The primary classes are:

|                                            |                                                                     |
|--------------------------------------------|---------------------------------------------------------------------|
| <i>org.dspace.content.WorkspaceItem</i>    | contains an Item before it enters a workflow                        |
| <i>org.dspace.workflow.WorkflowItem</i>    | contains an Item while in a workflow                                |
| <i>org.dspace.workflow.WorkflowManager</i> | responds to events, manages the WorkflowItem states                 |
| <i>org.dspace.content.Collection</i>       | contains List of defined workflow steps                             |
| <i>org.dspace.eperson.Group</i>            | people who can perform workflow tasks are defined in EPerson Groups |
| <i>org.dspace.core.Email</i>               | used to email messages to Group members and submitters              |

The workflow system models the states of an Item in a state machine with 5 states (SUBMIT, STEP\_1, STEP\_2, STEP\_3, ARCHIVE.) These are the three optional steps where the item can be viewed and corrected by different groups of people. Actually, it's more like 8 states, with STEP\_1\_POOL, STEP\_2\_POOL, and STEP\_3\_POOL. These pooled states are when items are waiting to enter the primary states.

The WorkflowManager is invoked by events. While an Item is being submitted, it is held by a WorkspaceItem. Calling the start() method in the WorkflowManager converts a WorkspaceItem to a WorkflowItem, and begins processing the WorkflowItem's state. Since all three steps of the workflow are optional, if no steps are defined, then the Item is simply archived.

Workflows are set per Collection, and steps are defined by creating corresponding entries in the List named workflowGroup. If you wish the workflow to have a step 1, use the administration tools for Collections to create a workflow Group with members who you want to be able to view and approve the Item, and the workflowGroup[0] becomes set with the ID of that Group.

If a step is defined in a Collection's workflow, then the WorkflowItem's state is set to that step\_POOL. This pooled state is the WorkflowItem waiting for an EPerson in that group to claim the step's task for that WorkflowItem. The WorkflowManager emails the members of that Group notifying them that there is a task to be performed (the text is defined in config/emails,) and when an EPerson goes to their 'My DSpace' page to claim the task, the WorkflowManager is invoked with a claim event, and the WorkflowItem's state advances from STEP\_x\_POOL to STEP\_x (where x is the corresponding step.) The EPerson can also generate an 'unclaim' event, returning the WorkflowItem to the STEP\_x\_POOL.

Other events the WorkflowManager handles are `advance()`, which advances the WorkflowItem to the next state. If there are no further states, then the WorkflowItem is removed, and the Item is then archived. An EPerson performing one of the tasks can reject the Item, which stops the workflow, rebuilds the WorkspaceItem for it and sends a rejection note to the submitter. More drastically, an `abort()` event is generated by the admin tools to cancel a workflow outright.

## Administration Toolkit

The `org.dspace.administer` package contains some classes for administering a DSpace system that are not generally needed by most applications.

The `CreateAdministrator` class is a simple command-line tool, executed via `[dspace]/bin/dspace create-administrator`, that creates an administrator e-person with information entered from standard input. This is generally used only once when a DSpace system is initially installed, to create an initial administrator who can then use the Web administration UI to further set up the system. This script does not check for authorization, since it is typically run before there are any e-people to authorize! Since it must be run as a command-line tool on the server machine, generally this shouldn't cause a problem. A possibility is to have the script only operate when there are no e-people in the system already, though in general, someone with access to command-line scripts on your server is probably in a position to do what they want anyway!

The `DCType` class is similar to the `org.dspace.content.BitstreamFormat` class. It represents an entry in the Dublin Core type registry, that is, a particular element and qualifier, or unqualified element. It is in the `administer` package because it is only generally required when manipulating the registry itself. Elements and qualifiers are specified as literals in `org.dspace.content.Item` methods and the `org.dspace.content.Metadatum` class. Only administrators may modify the Dublin Core type registry.

The `org.dspace.administer.RegistryLoader` class contains methods for initializing the Dublin Core type registry and bitstream format registry with entries in an XML file. Typically this is executed via the command line during the build process (see `build.xml` in the source.) To see examples of the XML formats, see the files in `config/registries` in the source directory. There is no XML schema, they aren't validated strictly when loaded in.

## E-person/Group Manager

DSpace keeps track of registered users with the `org.dspace.eperson.EPerson` class. The class has methods to create and manipulate an `EPerson` such as get and set methods for first and last names, email, and password. (Actually, there is no `getPassword()` method, an MD5 hash of the password is stored, and can only be verified with the `checkPassword()` method.) There are find methods to find an `EPerson` by email (which is assumed to be unique,) or to find all `EPeople` in the system.

The `EPerson` object should probably be reworked to allow for easy expansion; the current `EPerson` object tracks pretty much only what MIT was interested in tracking - first and last names, email, phone. The access methods are hardcoded and should probably be replaced with methods to access arbitrary name/value pairs for institutions that wish to customize what `EPerson` information is stored.

Groups are simply lists of `EPerson` objects. Other than membership, `Group` objects have only one other attribute: a name. Group names must be unique, so we have adopted naming conventions where the role of the group is its name, such as `COLLECTION_100_ADD`. Groups add and remove `EPerson` objects with `addMember()` and `removeMember()` methods. One important thing to know about groups is that they store their membership in memory until the `update()` method is called - so when modifying a group's membership don't forget to invoke `update()` or your changes will be lost! Since group membership is used heavily by the authorization system a fast `isMember()` method is also provided.

Another kind of Group is also implemented in DSpace, special Groups. The `Context` object for each session carries around a List of Group IDs that the user is also a member of, currently the `MITUser` Group ID is added to the list of a user's special groups if certain IP address or certificate criteria are met.

## Authorization

The primary classes are:

|                                                    |                                                          |
|----------------------------------------------------|----------------------------------------------------------|
| <code>org.dspace.authorize.AuthorizeManager</code> | does all authorization, checking policies against Groups |
| <code>org.dspace.authorize.ResourcePolicy</code>   | defines all allowable actions for an object              |
| <code>org.dspace.eperson.Group</code>              | all policies are defined in terms of EPerson Groups      |

The authorization system is based on the classic 'police state' model of security; no action is allowed unless it is expressed in a policy. The policies are attached to resources (hence the name `ResourcePolicy`.) and detail who can perform that action. The resource can be any of the DSpace object types, listed in `org.dspace.core.Constants` (`BITSTREAM`, `ITEM`, `COLLECTION`, etc.) The 'who' is made up of `EPerson` groups. The actions are also in `Constants.java` (`READ`, `WRITE`, `ADD`, etc.) The only non-obvious actions are `ADD` and `REMOVE`, which are authorizations for container objects. To be able to create an `Item`, you must have `ADD` permission in a `Collection`, which contains `Items`. (`Communities`, `Collections`, `Items`, and `Bundles` are all container objects.)

Currently most of the read policy checking is done with items, communities and collections are assumed to be openly readable, but items and their bitstreams are checked. Separate policy checks for items and their bitstreams enables policies that allow publicly readable items, but parts of their content may be restricted to certain groups.

Three new attributes have been introduced in the `ResourcePolicy` class as part of the DSpace 3.0 [Embargo](#) Contribution:

- `rpname`: resource policy name
- `rptype`: resource policy type
- `rpdescription`: resource policy description

While `rpname` and `rpdescription` are fields manageable by the users the `_rptype` is a field managed by the system. It represents a type that a resource policy can assume between the following:

- TYPE\_SUBMISSION: all the policies added automatically during the submission process
- TYPE\_WORKFLOW: all the policies added automatically during the workflow stage
- TYPE\_CUSTOM: all the custom policies added by users
- TYPE\_INHERITED: all the policies inherited by the DSO father.

An custom policy, created for the purpose of creating an embargo could look like:

```
policy_id: 4847
resource_type_id: 2
resource_id: 89
action_id: 0
eperson_id:
epersongroup_id: 0
start_date: 2013-01-01
end_date:
rpname: Embargo Policy
rpdescription: Embargoed through 2012
rptype: TYPE_CUSTOM
```

The *AuthorizeManager* class'

*authorizeAction(Context, object, action)* is the primary source of all authorization in the system. It gets a list of all of the ResourcePolicies in the system that match the object and action. It then iterates through the policies, extracting the EPerson Group from each policy, and checks to see if the EPersonID from the Context is a member of any of those groups. If all of the policies are queried and no permission is found, then an *AuthorizeException* is thrown. An *authorizeAction()* method is also supplied that returns a boolean for applications that require higher performance.

ResourcePolicies are very simple, and there are quite a lot of them. Each can only list a single group, a single action, and a single object. So each object will likely have several policies, and if multiple groups share permissions for actions on an object, each group will get its own policy. (It's a good thing they're small.)

## Special Groups

All users are assumed to be part of the public group (ID=0.) DSpace admins (ID=1) are automatically part of all groups, much like super-users in the Unix OS. The Context object also carries around a List of special groups, which are also first checked for membership. These special groups are used at MIT to indicate membership in the MIT community, something that is very difficult to enumerate in the database! When a user logs in with an MIT certificate or with an MIT IP address, the login code adds this MIT user group to the user's Context.

## Miscellaneous Authorization Notes

Where do items get their read policies? From the their collection's read policy. There once was a separate item read default policy in each collection, and perhaps there will be again since it appears that administrators are notoriously bad at defining collection's read policies. There is also code in place to enable policies that are timed, have a start and end date. However, the admin tools to enable these sorts of policies have not been written.

## Handle Manager/Handle Plugin

The *org.dspace.handle* package contains two classes; *HandleManager* is used to create and look up Handles, and *HandlePlugin* is used to expose and resolve DSpace Handles for the outside world via the CNRI Handle Server code.

Handles are stored internally in the *handle* database table in the form:

```
1721.123/4567
```

Typically when they are used outside of the system they are displayed in either URI or "URL proxy" forms:

```
hdl:1721.123/4567
http://hdl.handle.net/1721.123/4567
```

It is the responsibility of the caller to extract the basic form from whichever displayed form is used.

The *handle* table maps these Handles to resource type/resource ID pairs, where resource type is a value from *org.dspace.core.Constants* and resource ID is the internal identifier (database primary key) of the object. This allows Handles to be assigned to any type of object in the system, though as explained in the functional overview, only communities, collections and items are presently assigned Handles.

*HandleManager* contains static methods for:

- Creating a Handle
- Finding the Handle for a *DSpaceObject*, though this is usually only invoked by the object itself, since *DSpaceObject* has a *getHandle* method
- Retrieving the *DSpaceObject* identified by a particular Handle
- Obtaining displayable forms of the Handle (URI or "proxy URL").

*HandlePlugin* is a simple implementation of the Handle Server's *net.handle.hdl.lib.HandleStorage* interface. It only implements the basic Handle retrieval methods, which get information from the *handle* database table. The CNRI Handle Server is configured to use this plug-in via its *config.dct* file.

Note that since the Handle server runs as a separate JVM to the DSpace Web applications, it uses a separate 'Log4J' configuration, since Log4J does not support multiple JVMs using the same daily rolling logs. This alternative configuration is located at `[dspace]/config/log4j-handle-plugin.properties`. The `[dspace]/bin/start-handle-server` script passes in the appropriate command line parameters so that the Handle server uses this configuration.

## Search

DSpace's search code is a simple API which currently wraps the Lucene search engine. The first half of the search task is indexing, and *org.dspace.search.DSIndexer* is the indexing class, which contains *indexContent()* which if passed an *Item*, *Community*, or *Collection*, will add that content's fields to the index. The methods *unIndexContent()* and *reIndexContent()* remove and update content's index information. The *DSIndexer* class also has a *main()* method which will rebuild the index completely. This can be invoked by the *dspace/bin/index-init* (complete rebuild) or *dspace/bin/index-update* (update) script. The intent was for the *main()* method to be invoked on a regular basis to avoid index corruption, but we have had no problem with that so far.

Which fields are indexed by *DSIndexer*? These fields are defined in *dspace.cfg* in the section "Fields to index for search" as name-value-pairs. The name must be unique in the form *search.index.i* (i is an arbitrary positive number). The value on the right side has a unique value again, which can be referenced in search-form (e.g. title, author). Then comes the metadata element which is indexed. *\** is a wildcard which includes all sub elements. For example:

```
search.index.4 = keyword:dc.subject.*
```

tells the indexer to create a keyword index containing all *dc.subject* element values. Since the wildcard (*\**) character was used in place of a qualifier, all subject metadata fields will be indexed (e.g. *dc.subject.other*, *dc.subject.lcsh*, etc)

By default, the fields shown in the *Indexed Fields* section below are indexed. These are hardcoded in the *DSIndexer* class. If any *search.index.i* items are specified in *dspace.cfg* these are used rather than these hardcoded fields.

The query class *DSQuery* contains the three flavors of *doQuery()* methods, one searches the DSpace site, and the other two restrict searches to Collections and Communities. The results from a query are returned as three lists of handles; each list represents a type of result. One list is a list of Items with matches, and the other two are Collections and Communities that match. This separation allows the UI to handle the types of results gracefully without resolving all of the handles first to see what kind of content the handle points to. The *DSQuery* class also has a *main()* method for debugging via command-line searches.

## Current Lucene Implementation

Currently we have our own Analyzer and Tokenizer classes (*DSAnalyzer* and *DSTokenizer*) to customize our indexing. They invoke the stemming and stop word features within Lucene. We create an *IndexReader* for each query, which we now realize isn't the most efficient use of resources - we seem to run out of filehandles on really heavy loads. (A wildcard query can open many filehandles!) Since Lucene is thread-safe, a better future implementation would be to have a single Lucene *IndexReader* shared by all queries, and then is invalidated and re-opened when the index changes. Future API growth could include relevance scores (Lucene generates them, but we ignore them,) and abstractions for more advanced search concepts such as booleans.

## Indexed Fields

The *DSIndexer* class shipped with DSpace indexes the Dublin Core metadata in the following way:

| Search Field | Taken from Dublin Core Fields                                    |
|--------------|------------------------------------------------------------------|
| Authors      | <i>contributor.creator.description.statementofresponsibility</i> |
| Titles       | <i>title.*</i>                                                   |
| Keywords     | <i>subject.*</i>                                                 |
| Abstracts    | <i>description.abstractdescription.tableofcontents</i>           |
| Series       | <i>relation.ispartofseries</i>                                   |
| MIME types   | <i>format.mimetype</i>                                           |
| Sponsors     | <i>description.sponsorship</i>                                   |
| Identifiers  | <i>identifier.*</i>                                              |

## Harvesting API

The *org.dspace.search* package also provides a 'harvesting' API. This allows callers to extract information about items modified within a particular timeframe, and within a particular scope (all of DSpace, or a community or collection.) Currently this is used by the Open Archives Initiative metadata harvesting protocol application, and the e-mail subscription code.

The *Harvest.harvest* is invoked with the required scope and start and end dates. Either date can be omitted. The dates should be in the ISO8601, UTC time zone format used elsewhere in the DSpace system.

*HarvestedItemInfo* objects are returned. These objects are simple containers with basic information about the items falling within the given scope and date range. Depending on parameters passed to the *harvest* method, the *containers* and *item* fields may have been filled out with the IDs of communities and collections containing an item, and the corresponding *Item* object respectively. Electing not to have these fields filled out means the harvest operation executes considerably faster.

In case it is required, *Harvest* also offers a method for creating a single *HarvestedItemInfo* object, which might make things easier for the caller.

## Browse API

The browse API maintains indexes of dates, authors, titles and subjects, and allows callers to extract parts of these:

- **Title:** Values of the Dublin Core element **title** (unqualified) are indexed. These are sorted in a case-insensitive fashion, with any leading article removed. For example: "The DSpace System" would appear under 'D' rather than 'T'.
- **Author:** Values of the **contributor** (any qualifier or unqualified) element are indexed. Since *contributor* values typically are in the form 'last name, first name', a simple case-insensitive alphanumeric sort is used which orders authors in last name order. Note that this is an index of *authors*, and not *items by author*. If four items have the same author, that author will appear in the index only once. Hence, the index of authors may be greater or smaller than the index of titles; items often have more than one author, though the same author may have authored several items. The author indexing in the browse API does have limitations:
  - Ideally, a name that appears as an author for more than one item would appear in the author index only once. For example, 'Doe, John' may be the author of tens of items. However, in practice, author's names often appear in slightly differently forms, for example:

```
Doe, John
Doe, John Stewart
Doe, John S.
```

Currently, the above three names would all appear as separate entries in the author index even though they may refer to the same author. In order for an author of several papers to be correctly appear once in the index, each item must specify *exactly* the same form of their name, which doesn't always happen in practice.

- Another issue is that two authors may have the same name, even within a single institution. If this is the case they may appear as one author in the index. These issues are typically resolved in libraries with *authority control records*, in which are kept a 'preferred' form of the author's name, with extra information (such as date of birth/death) in order to distinguish between authors of the same name. Maintaining such records is a huge task with many issues, particularly when metadata is received from faculty directly rather than trained library catalogers.
- **Date of Issue:** Items are indexed by date of issue. This may be different from the date that an item appeared in DSpace; many items may have been originally published elsewhere beforehand. The Dublin Core field used is **date.issued**. The ordering of this index may be reversed so 'earliest first' and 'most recent first' orderings are possible. Note that the index is of *items by date*, as opposed to an index of *dates*. If 30 items have the same issue date (say 2002), then those 30 items all appear in the index adjacent to each other, as opposed to a single 2002 entry. Since dates in DSpace Dublin Core are in ISO8601, all in the UTC time zone, a simple alphanumeric sort is sufficient to sort by date, including dealing with varying granularities of date reasonably. For example:

```
2001-12-10
2002
2002-04
2002-04-05
2002-04-09T15:34:12Z
2002-04-09T19:21:12Z
2002-04-10
```

- **Date Accessioned:** In order to determine which items most recently appeared, rather than using the date of issue, an item's accession date is used. This is the Dublin Core field **date.accessioned**. In other aspects this index is identical to the date of issue index.
- **Items by a Particular Author:** The browse API can perform is to extract items by a particular author. They do not have to be primary author of an item for that item to be extracted. You can specify a scope, too; that is, you can ask for items by author X in collection Y, for example. This particular flavor of browse is slightly simpler than the others. You cannot presently specify a particular subset of results to be returned. The API call will simply return all of the items by a particular author within a certain scope. Note that the author of the item must *exactly* match the author passed in to the API; see the explanation about the caveats of the author index browsing to see why this is the case.
- **Subject:** Values of the Dublin Core element **subject** (both unqualified and with any qualifier) are indexed. These are sorted in a case-insensitive fashion.

## Using the API

The API is generally invoked by creating a *BrowseScope* object, and setting the parameters for which particular part of an index you want to extract. This is then passed to the relevant *Browse* method call, which returns a *BrowseInfo* object which contains the results of the operation. The parameters set in the *BrowseScope* object are:

- How many entries from the index you want
- Whether you only want entries from a particular community or collection, or from the whole of DSpace
- Which part of the index to start from (called the *focus* of the browse). If you don't specify this, the start of the index is used
- How many entries to include before the *focus* entry

To illustrate, here is an example:

- We want **7** entries in total

- We want entries from collection *x*
- We want the focus to be 'Really'
- We want **2** entries included before the focus.

The results of invoking `Browse.getItemsByTitle` with the above parameters might look like this:

```
Rabble-Rousing Rabbis From Sardinia
 Reality TV: Love It or Hate It?
FOCUS> The Really Exciting Research Video
 Recreational Housework Addicts: Please Visit My House
 Regional Television Variation Studies
 Revenue Streams
 Ridiculous Example Titles: I'm Out of Ideas
```

Note that in the case of title and date browses, *Item* objects are returned as opposed to actual titles. In these cases, you can specify the 'focus' to be a specific item, or a partial or full literal value. In the case of a literal value, if no entry in the index matches exactly, the closest match is used as the focus. It's quite reasonable to specify a focus of a single letter, for example.

Being able to specify a specific item to start at is particularly important with dates, since many items may have the same issue date. Say 30 items in a collection have the issue date 2002. To be able to page through the index 20 items at a time, you need to be able to specify exactly which item's 2002 is the focus of the browse, otherwise each time you invoked the browse code, the results would start at the first item with the issue date 2002.

Author browses return *String* objects with the actual author names. You can only specify the focus as a full or partial literal *String*.

Another important point to note is that presently, the browse indexes contain metadata for all items in the main archive, regardless of authorization policies. This means that all items in the archive will appear to all users when browsing. Of course, should the user attempt to access a non-public item, the usual authorization mechanism will apply. Whether this approach is ideal is under review; implementing the browse API such that the results retrieved reflect a user's level of authorization may be possible, but rather tricky.

## Index Maintenance

The browse API contains calls to add and remove items from the index, and to regenerate the indexes from scratch. In general the content management API invokes the necessary browse API calls to keep the browse indexes in sync with what is in the archive, so most applications will not need to invoke those methods.

If the browse index becomes inconsistent for some reason, the `InitializeBrowse` class is a command line tool (generally invoked using the `[dSPACE]/bin/dSPACE index-init` command) that causes the indexes to be regenerated from scratch.

## Caveats

Presently, the browse API is not tremendously efficient. 'Indexing' takes the form of simply extracting the relevant Dublin Core value, normalizing it (lower-casing and removing any leading article in the case of titles), and inserting that normalized value with the corresponding item ID in the appropriate browse database table. Database views of this table include collection and community IDs for browse operations with a limited scope. When a browse operation is performed, a simple `SELECT` query is performed, along the lines of:

```
SELECT item_id FROM ItemsByTitle ORDER BY sort_title OFFSET 40 LIMIT 20
```

There are two main drawbacks to this: Firstly, `LIMIT` and `OFFSET` are PostgreSQL-specific keywords. Secondly, the database is still actually performing dynamic sorting of the titles, so the browse code as it stands will not scale particularly well. The code does cache `BrowseInfo` objects, so that common browse operations are performed quickly, but this is not an ideal solution.

## Checksum checker

Checksum checker is used to verify every item within DSpace. While DSpace calculates and records the checksum of every file submitted to it, the checker can determine whether the file has been changed. The idea being that the earlier you can identify a file has changed, the more likely you would be able to record it (assuming it was not a wanted change).

`org.dspace.checker.CheckerCommand` class, is the class for the checksum checker tool, which calculates checksums for each bitstream whose ID is in the `most_recent_checksum` table, and compares it against the last calculated checksum for that bitstream.

## OpenSearch Support

DSpace is able to support OpenSearch. For those not acquainted with the standard, a very brief introduction, with emphasis on what possibilities it holds for current use and future development.

OpenSearch is a small set of conventions and documents for describing and using 'search engines', meaning any service that returns a set of results for a query. It is nearly ubiquitous, but also nearly invisible, in modern web sites with search capability. If you look at the page source of Wikipedia, Facebook, CNN, etc you will find buried a link element declaring OpenSearch support. It is very much a lowest-common-denominator abstraction (think Google box), but does provide a means to extend its expressive power. This first implementation for DSpace supports *none* of these extensions, many of which are of potential value, so it should be regarded as a foundation, not a finished solution. So the short answer is that DSpace appears as a 'search-engine' to OpenSearch-aware software.

Another way to look at OpenSearch is as a RESTful web service for search, very much like SRW/U, but considerably simpler. This comparative loss of power is offset by the fact that it is widely supported by web tools and players: browsers understand it, as do large metasearch tools.

### How Can It Be Used

- **Browser Integration** Many recent browsers (IE7+, FF2+) can detect, or 'autodiscover', links to the document describing the search engine. Thus you can easily add your or other DSpace instances to the drop-down list of search engines in your browser. This list typically appears in the upper right corner of the browser, with a search box. In Firefox, for example, when you visit a site supporting OpenSearch, the color of the drop-down list widget changes color, and if you open it to show the list of search engines, you are offered an opportunity to add the site to the list. IE works nearly the same way but instead labels the web sites 'search providers'. When you select a DSpace instance as the search engine and enter a search, you are simply sent to the regular search results page of the instance.
- **Flexible, interesting RSS Feeds** Because one of the formats that OpenSearch specifies for its results is RSS (or Atom), you can turn any search query into an RSS feed. So if there are keywords highly discriminative of content in a collection or repository, these can be turned into a URL that a feed reader can subscribe to. Taken to the extreme, one could take any search a user makes, and dynamically compose an RSS feed URL for it in the page of returned results. To see an example, if you have a DSpace with OpenSearch enabled, try:

```
http://dspace.mysite.edu/open-search/?query=<your query>
```

The default format returned is Atom 1.0, so you should see an Atom document containing your search results.

- You can extend the syntax with a few other parameters, as follows:

| Parameter | Values                                                                     |
|-----------|----------------------------------------------------------------------------|
| format    | atom, rss, html                                                            |
| scope     | handle of a collection or community to restrict the search to              |
| rpp       | number indicating the number of results per page (i.e. per request)        |
| start     | number of page to start with (if paginating results)                       |
| sort_by   | number indicating sorting criteria (same as DSpace advanced search values) |

Multiple parameters may be specified on the query string, using the "&" character as the delimiter, e.g.:

```
http://dspace.mysite.edu/open-search/?query=<your query>&format=rss&scope=123456789/1
```

- **Cheap metasearch** Search aggregators like A9 (Amazon) recognize OpenSearch-compliant providers, and so can be added to metasearch sets using their UIs. Then your site can be used to aggregate search results with others.

Configuration is through the `dspace.cfg` file. See [OpenSearch Support](#) for more details.

## Embargo Support

### What is an Embargo?

An embargo is a temporary access restriction placed on content, commencing at time of accession. Its scope or duration may vary, but the fact that it eventually expires is what distinguishes it from other content restrictions. For example, it is not unusual for content destined for DSpace to come with permanent restrictions on use or access based on license-driven or other IP-based requirements that limit access to institutionally affiliated users. Restrictions such as these are imposed and managed using standard administrative tools in DSpace, typically by attaching specific policies to Items or Collections, Bitstreams, etc. The embargo functionality introduced in 1.6, however, includes tools to automate the imposition and removal of restrictions in managed timeframes.

### Embargo Model and Life-Cycle

Functionally, the embargo system allows you to attach 'terms' to an item before it is placed into the repository, which express how the embargo should be applied. What do 'we mean by terms' here? They are really any expression that the system is capable of turning into (1) the time the embargo expires, and (2) a concrete set of access restrictions. Some examples:

"2020-09-12" - an absolute date (i.e. the date embargo will be lifted)  
"6 months" - a time relative to when the item is accessioned  
"forever" - an indefinite, or open-ended embargo  
"local only until 2015" - both a time and an exception (public has no access until 2015, local users OK immediately)  
"Nature Publishing Group standard" - look-up to a policy somewhere (typically 6 months)

These terms are 'interpreted' by the embargo system to yield a specific date on which the embargo can be removed or 'lifted', and a specific set of access policies. Obviously, some terms are easier to interpret than others (the absolute date really requires none at all), and the 'default' embargo logic understands only the most basic terms (the first and third examples above). But as we will see below, the embargo system provides you with the ability to add in your own 'interpreters' to cope with any terms expressions you wish to have. This date that is the result of the interpretation is stored with the item and the embargo system detects when that date has passed, and removes the embargo ("lifts it"), so the item bitstreams become available. Here is a more detailed life-cycle for an embargoed item:

1. **Terms Assignment.** The first step in placing an embargo on an item is to attach (assign) 'terms' to it. If these terms are missing, no embargo will be imposed. As we will see below, terms are carried in a configurable DSpace metadata field, so assigning terms just means assigning a value to a metadata field. This can be done in a web submission user interface form, in a SWORD deposit package, a batch import, etc. - anywhere

metadata is passed to DSpace. The terms are not immediately acted upon, and may be revised, corrected, removed, etc, up until the next stage of the life-cycle. Thus a submitter could enter one value, and a collection editor replace it, and only the last value will be used. Since metadata fields are multivalued, theoretically there can be multiple terms values, but in the default implementation only one is recognized.

2. **Terms interpretation/imposition.** In DSpace terminology, when an item has exited the last of any workflow steps (or if none have been defined for it), it is said to be 'installed' into the repository. At this precise time, the 'interpretation' of the terms occurs, and a computed 'lift date' is assigned, which like the terms is recorded in a configurable metadata field. It is important to understand that this interpretation happens only once, (just like the installation), and cannot be revisited later. Thus, although an administrator can assign a new value to the metadata field holding the terms after the item has been installed, this will have no effect on the embargo, whose 'force' now resides entirely in the 'lift date' value. For this reason, you cannot embargo content already in your repository (at least using standard tools). The other action taken at installation time is the actual imposition of the embargo. The default behavior here is simply to remove the read policies on all the bundles and bitstreams except for the "LICENSE" or "METADATA" bundles. See the section on *Extending Embargo Functionality* for how to alter this behavior. Also note that since these policy changes occur before installation, there is no time during which embargoed content is 'exposed' (accessible by non-administrators). The terms interpretation and imposition together are called 'setting' the embargo, and the component that performs them both is called the embargo 'setter'.
3. **Embargo Period.** After an embargoed item has been installed, the policy restrictions remain in effect until removed. This is not an automatic process, however: a 'lifter' must be run periodically to look for items whose 'lift date' is past. Note that this means the effective removal of an embargo is **not** the lift date, but the earliest date after the lift date that the lifter is run. Typically, a nightly cron-scheduled invocation of the lifter is more than adequate, given the granularity of embargo terms. Also note that during the embargo period, all metadata of the item remains visible. This default behavior can be changed. One final point to note is that the 'lift date', although it was computed and assigned during the previous stage, is in the end a regular metadata field. That means, if there are extraordinary circumstances that require an administrator (or collection editor, anyone with edit permissions on metadata) to change the lift date, they can do so. Thus, they can 'revise' the lift date without reference to the original terms. This date will be checked the next time the 'lifter' is run. One could immediately lift the embargo by setting the lift date to the current day, or change it to 'forever' to indefinitely postpone lifting.
4. **Embargo Lift.** When the lifter discovers an item whose lift date is in the past, it removes (lifts) the embargo. The default behavior of the lifter is to add the resource policies *that would have been added* had the embargo not been imposed. That is, it replicates the standard DSpace behavior, in which an item inherits its policies from its owning collection. As with all other parts of the embargo system, you may replace or extend the default behavior of the lifter (see section V. below). You may wish, e.g. to send an email to an administrator or other interested parties, when an embargoed item becomes available.
5. **Post Embargo.** After the embargo has been lifted, the item ceases to respond to any of the embargo life-cycle events. The values of the metadata fields reflect essentially historical or provenance values. With the exception of the additional metadata fields, they are indistinguishable from items that were never subject to embargo.

More Embargo Details



More details on Embargo configuration, including specific examples can be found in the [Embargo](#) section of the documentation.