

# DSpace 1.4.1 beta 1 System Documentation

Robert Tansley  
Mick Bass  
Margret Branschofsky  
Grace Carpenter  
Greg McClellan  
David Stuve  
and many other helping hands

09-November-2006

# Contents

<b>List of Tables</b>	<b>5</b>
<b>List of Figures</b>	<b>6</b>
<b>1 Introduction</b>	<b>8</b>
<b>2 Functional Overview</b>	<b>10</b>
2.1 Data Model	10
2.2 Plugin Manager	12
2.3 Metadata	13
2.4 Packager Plugins	13
2.5 Crosswalk Plugins	14
2.6 E-people and Groups	14
2.6.1 E-Person	14
2.6.2 Groups	14
2.7 Authentication	15
2.8 Authorization	15
2.9 Ingest Process and Workflow	16
2.9.1 Workflow Steps	17
2.10 Supervision and Collaboration	18
2.11 Handles	18
2.12 Bitstream 'Persistent' Identifiers	19
2.13 Storage Resource Broker (SRB) Support	19
2.14 Search and Browse	20
2.15 HTML Support	20
2.16 OAI Support	21
2.17 OpenURL Support	21
2.18 Creative Commons Support	21
2.19 Subscriptions	22
2.20 History	22
2.21 Import and Export	22
2.22 Registration	22
2.23 Statistical Reports	22
2.24 Checksum Checker	23

<b>3</b>	<b>Installation</b>	<b>24</b>
3.1	Prerequisite Software	24
3.2	Quick Installation Steps	25
3.3	Advanced Installation	27
3.3.1	'cron' Jobs	28
3.3.2	DSpace over HTTPS	28
3.3.3	The Handle Server	32
3.4	Windows Installation	33
3.4.1	Pre-requisite Software	33
3.4.2	Installation Steps	33
3.5	Known Bugs	35
3.6	Common Problems	35
<b>4</b>	<b>Updating a DSpace Installation</b>	<b>37</b>
4.1	Updating From 1.4 to 1.4.x	37
4.2	Updating From 1.3.x to 1.4.x	38
4.3	Updating From 1.3.1 to 1.3.2	40
4.4	Updating From 1.2.x to 1.3.x	41
4.5	Updating From 1.2.1 to 1.2.2	42
4.6	Updating From 1.2 to 1.2.1	43
4.7	Updating From 1.1 (or 1.1.1) to 1.2	45
4.8	Updating From 1.1 to 1.1.1	48
4.9	Updating From 1.0.1 to 1.1	48
<b>5</b>	<b>Configuration and Customization</b>	<b>51</b>
5.1	The dspace.cfg Configuration Properties File	51
5.2	Wording of E-mail Messages	53
5.3	The Metadata and Bitstream Format Registries	53
5.3.1	Metadata Format Registries	53
5.3.2	Bitstream Format Registry	54
5.4	The Default Submission License	54
5.4.1	Possible Points in a License	54
5.5	Activating Additional OAI-PMH Crosswalks	55
5.5.1	METS	55
5.5.2	Crosswalk Plugins	55
5.5.3	DIDL	55
5.6	Configuration Files for Other Applications	56
5.7	Customizing the Web User Interface	56
5.8	Customizing the Simple Item Display Metadata	57
5.8.1	Customizing the Simple Item Display Metadata for Individual Collections	58
5.9	Custom Authentication Code	58
5.9.1	Authentication by Password	59
5.9.2	LDAP Authentication	59
5.9.3	X.509 Certificate Authentication	59
5.9.4	Example of a Custom Authentication Method	60
5.10	Configuring LDAP Authentication	60
5.11	Configuring Lucene Search Indexes	61
5.12	Configuring System Statistical Reports	62
5.12.1	Configuration File	62

5.12.2	Customising Shell Scripts	62
5.13	MediaFilters	64
5.14	Displaying Image Item Preview	65
5.15	Displaying Image Thumbnails	65
5.15.1	Browse and Search Results Page Thumbnails	65
5.15.2	Configuring Thumbnail Link Behaviour	66
5.15.3	Item Display Page Thumbnails	66
5.16	Displaying Community and Collection Item Counts	66
5.17	Lucene Analyzer	66
5.18	On-line Help About File Formats	67
5.19	View Item Licence	67
5.20	Configuring RSS Syndication	67
5.21	Configuring Item Recommendations	68
5.22	Configuring Controlled Vocabularies	68
5.23	Configuring the checksum checker	70
5.23.1	Checker Execution Mode	70
5.23.2	Checker Reporting	71
5.23.3	Checker Results Pruning	71
5.24	Configuring Packager Plugins	72
5.25	Configuring Crosswalk Plugins	72
5.25.1	Configurable MODS dissemination crosswalk	72
5.25.2	Configurable Qualified Dublin Core (QDC) dissemination crosswalk	73
5.25.3	XSLT-based crosswalks	74
5.25.4	DSpace Intermediate Metadata (DIM) format	74
5.25.5	Testing XSLT Crosswalks	75
<b>6</b>	<b>Directories and Files</b>	<b>76</b>
6.1	Source Directory Layout	76
6.2	Contents of Web Application	78
6.3	Log Files	78
<b>7</b>	<b>Architecture</b>	<b>80</b>
7.1	Overview	80
7.2	Storage Layer	81
7.2.1	RDBMS	81
7.2.2	Bitstream Store	83
7.2.3	Backup	85
7.3	Business Logic Layer	86
7.3.1	Core Classes	86
7.3.2	Content Management API	89
7.3.3	Plugin Manager	93
7.3.4	Workflow System	101
7.3.5	Administration Toolkit	102
7.3.6	E-person/Group Manager	102
7.3.7	Authorization	103
7.3.8	Handle Manager/Handle Plugin	104
7.3.9	Search	105
7.3.10	Browse API	106
7.3.11	History Recorder	109

7.3.12	Checksum checker	111
7.4	Application Layer	111
7.4.1	Web User Interface	111
7.4.2	OAI-PMH Data Provider	119
7.4.3	Package Importer and Exporter	122
7.4.4	Item Importer and Exporter	123
7.4.5	Transferring Items Between DSpace Instances	125
7.4.6	Registering (Not Importing) Bitstreams	126
7.4.7	METS Tools	128
7.4.8	Media Filters	129
7.4.9	Sub-Community Management	130
<b>8</b>	<b>Version History</b>	<b>132</b>
8.1	Changes in DSpace 1.4.1	132
8.1.1	General Improvements	132
8.1.2	Bug fixes	132
8.2	Changes in DSpace 1.4	133
8.2.1	General Improvements	133
8.2.2	Bug fixes	134
8.3	Changes in DSpace 1.3.2	135
8.3.1	General Improvements	135
8.3.2	Bug fixes	135
8.4	Changes in DSpace 1.3.1	135
8.4.1	Bug fixes	135
8.5	Changes in DSpace 1.3	135
8.5.1	General Improvements	135
8.5.2	Bug fixes	136
8.6	Changes in DSpace 1.2.2	137
8.6.1	General Improvements	137
8.6.2	Bug fixes	137
8.6.3	Changes in JSPs	137
8.7	Changes in DSpace 1.2.1	138
8.7.1	General Improvements	138
8.7.2	Bug fixes	138
8.7.3	Changed JSPs	139
8.8	Changes in DSpace 1.2	139
8.8.1	General Improvements	139
8.9	Administration	140
8.9.1	Import/Export/OAI	140
8.9.2	Miscellaneous	140
8.9.3	JSP file changes between 1.1 and 1.2	140
8.10	Changes in DSpace 1.1.1	144
8.10.1	Bug fixes	144
8.10.2	Improvements	144
8.11	Changes in DSpace 1.1	144
<b>A</b>	<b>List of Abbreviations</b>	<b>146</b>
<b>B</b>	<b>Default Dublin Core Metadata Registry</b>	<b>149</b>

*CONTENTS*

5

**C Default Bitstream Format Registry**

**152**

# List of Tables

2.1	Objects in the DSpace Data Model . . . . .	12
6.1	DSpace Log File Locations . . . . .	79
7.1	Source Code Packages . . . . .	81
7.2	Example Logging Entry . . . . .	89
7.3	Dublin Core API . . . . .	92
7.4	Workflow System - primary classes . . . . .	101
7.5	Authorization - primary classes . . . . .	103
7.6	Indexed Fields . . . . .	106
7.7	Locations of Web UI Source Files . . . . .	112
B.1	Default Dublin Core Registry . . . . .	151
C.1	Default Bitstream Format Registry . . . . .	153

# List of Figures

- 2.1 Data Model Diagram . . . . . 11
- 2.2 Ingest Process . . . . . 16
- 2.3 Submission Workflow in DSpace . . . . . 17
  
- 7.1 DSpace System Architecture . . . . . 80
- 7.2 Flow of Control During HTTP Request Processing . . . . . 115



# Chapter 1

## Introduction

DSpace is an open source software platform that enables organisations to:

- capture and describe digital material using a submission workflow module, or a variety of programmatic ingest options
- distribute an organisation's digital assets over the web through a search and retrieval system
- preserve digital assets over the long term

This system documentation includes a functional overview [functional overview](#) of the system, which is a good introduction to the capabilities of the system, and should be readable by non-technical folk. Everyone should read this section first because it introduces some terminology used throughout the rest of the documentation. For people actually running a DSpace service, there is an [installation guide](#), and sections on [configuration](#) and the [directory structure](#). Note that as of DSpace 1.2, the administration user interface guide is now on-line help available from within the DSpace system.

Finally, for those interested in the details of how DSpace works, and those potentially interested in modifying the code for their own purposes, there is a [detailed architecture and design section](#).

Other good sources of information are:

- The DSpace Public API Javadocs. Build these with the `public_api` Ant target.
- The DSpace Wiki contains stacks of useful information about the DSpace platform and the work people are doing with it. You are strongly encouraged to visit this site and add information about your own work. Useful Wiki areas are:
  - The [DSpace Wiki](#) contains stacks of useful information about the DSpace platform and the work people are doing with it. You are strongly encouraged to visit this site and add information about your own work. Useful Wiki areas are:
    - [A list of DSpace resources \(Web sites, mailing lists etc.\)](#)
    - [Technical FAQ](#)
    - [A list of projects using DSpace](#)
    - [Guidelines for contributing back to DSpace](#)
  - [dspace.org](#) has announcements and contains useful information about bringing up an instance of DSpace at your organization.
- The University of Tennessee's Jason Simms has written some additional [installation notes](#).

- The [dspace-tech e-mail list](#) on SourceForge is the recommended place to ask questions, since a growing community of DSpace developers and users is on hand on that list to help with any questions you might have. The e-mail archive of that list is a useful resource.
- The [dspace-devel e-mail list](#), for those developing with the DSpace with a view to contributing to the core DSpace code.

## Chapter 2

# Functional Overview

### 2.1 Data Model

The way data is organized in DSpace is intended to reflect the structure of the organization using the DSpace system. Each DSpace site is divided into communities; these typically correspond to a laboratory, research center or department. As of DSpace version 1.2, these communities can be organized into an hierarchy. Communities contain collections, which are groupings of related content. A collection may appear in more than one community.

Each collection is composed of items, which are the basic archival elements of the archive. Each item is owned by one collection. Additionally, an item may appear in additional collections; however every item has one and only one owning collection.

Items are further subdivided into named bundles of bitstreams. Bitstreams are, as the name suggests, streams of bits, usually ordinary computer files. Bitstreams that are somehow closely related, for example HTML files and images that compose a single HTML document, are organised into bundles.

In practice, most items tend to have these named bundles:

**ORIGINAL** the bundle with the original, deposited bitstreams

**THUMBNAILS** thumbnails of any image bitstreams

**TEXT** extracted full-text from bitstreams in ORIGINAL, for indexing

**LICENSE** contains the deposit license that the submitter granted the host organization; in other words, specifies the rights that the hosting organization have

**CC\_LICENSE** contains the distribution license, if any (a [Creative Commons](#) license) associated with the item. This license specifies what end users downloading the content can do with the content

Each bitstream is associated with one Bitstream Format. Because preservation services may be an important aspect of the DSpace service, it is important to capture the specific formats of files that users submit. In DSpace, a bitstream format is a unique and consistent way to refer to a particular file format. An integral part of a bitstream format is an either implicit or explicit notion of how material in that format can be interpreted. For example, the interpretation for bitstreams encoded in the JPEG standard for still image compression is defined explicitly in the Standard ISO/IEC 10918-1. The interpretation of bitstreams in Microsoft Word 2000 format is defined implicitly, through reference to the Microsoft Word 2000 application. Bitstream formats can be more specific than MIME types or file suffixes. For example, application/ms-word and .doc span multiple versions of the Microsoft Word application, each of which produces bitstreams with presumably different characteristics.

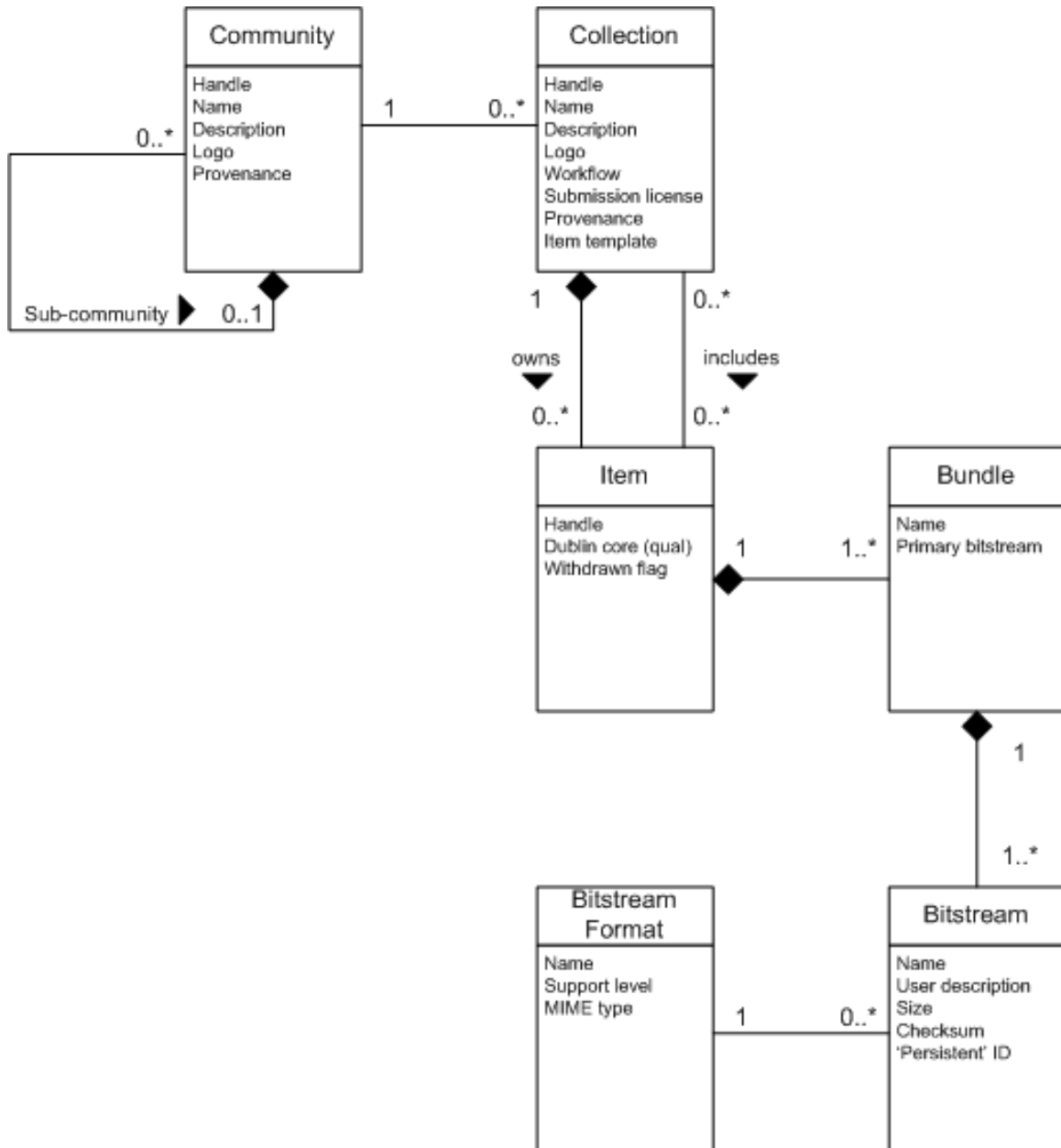


Figure 2.1: Data Model Diagram

Each bitstream format additionally has a support level, indicating how well the hosting institution is likely to be able to preserve content in the format in the future. There are three possible support levels that bitstream formats may be assigned by the hosting institution. The host institution should determine the exact meaning of each support level, after careful consideration of costs and requirements. MIT Libraries' interpretation is shown below:

### MIT Libraries' Definitions of Bitstream Format Support Levels

**Supported** The format is recognized, and the hosting institution is confident it can make bitstreams of this format useable in the future, using whatever combination of techniques (such as migration, emulation, etc.) is appropriate given the context of need.

**Known** The format is recognized, and the hosting institution will promise to preserve the bitstream as-is, and allow it to be retrieved. The hosting institution will attempt to obtain enough information to enable the format to be upgraded to the 'supported' level.

**Unsupported** The format is unrecognized, but the hosting institution will undertake to preserve the bitstream as-is and allow it to be retrieved.

Each item has one qualified Dublin Core metadata record. Other metadata might be stored in an item as a serialized bitstream, but we store Dublin Core for every item for interoperability and ease of discovery. The Dublin Core may be entered by end-users as they submit content, or it might be derived from other metadata as part of an ingest process.

Items can be removed from DSpace in one of two ways: They may be 'withdrawn', which means they remain in the archive but are completely hidden from view. In this case, if an end-user attempts to access the withdrawn item, they are presented with a 'tombstone,' that indicates the item has been removed. For whatever reason, an item may also be 'expunged' if necessary, in which case all traces of it are removed from the archive.

Object	Example
Community Collection	Laboratory of Computer Science; Oceanographic Research Center
Item	LCS Technical Reports; ORC Statistical Data Sets
Bundle	A technical report; a data set with accompanying description; a video recording of a lecture
Bitstream	A group of HTML and image bitstreams making up an HTML document
Bitstream Format	A single HTML file; a single image file; a source code file
	Microsoft Word version 6.0; JPEG encoded image format

Table 2.1: Objects in the DSpace Data Model

## 2.2 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 lifecycle of a plugin.

A plugin is defined by a Java interface. The consumer of a plugin asks for its plugin by interface. A Plugin is an instance of any class that implements the plugin interface. It is interchangeable with other implementations, so that any of them may be "plugged in".

The mediafilter is a simple example of a plugin implementation. Refer to the [Business Logic Layer](#) for more details on Plugins.

## 2.3 Metadata

Broadly speaking, DSpace holds three sorts of metadata about archived content:

### **Descriptive Metadata**

DSpace can support multiple flat metadata schemas for describing an item. A qualified Dublin Core metadata schema loosely based on the [Library Application Profile](#) set of elements and qualifiers is provided by default. The [set of elements and qualifiers](#) used by MIT Libraries comes pre-configured with the DSpace source code. However, you can configure multiple schemas and select metadata fields from a mix of configured schemas to describe your items.

Other descriptive metadata about items (e.g. metadata described in a hierarchical schema) may be held in serialized bitstreams. Communities and collections have some simple descriptive metadata (a name, and some descriptive prose), held in the DBMS.

### **Administrative Metadata**

This includes preservation metadata, provenance and authorization policy data. Most of this is held within DSpace's relation DBMS schema. Provenance metadata (prose) is stored in Dublin Core records. Additionally, some other administrative metadata (for example, bitstream byte sizes and MIME types) is replicated in Dublin Core records so that it is easily accessible outside of DSpace.

### **Structural Metadata**

This includes information about how to present an item, or bitstreams within an item, to an end-user, and the relationships between constituent parts of the item. As an example, consider a thesis consisting of a number of TIFF images, each depicting a single page of the thesis. Structural metadata would include the fact that each image is a single page, and the ordering of the TIFF images/pages. Structural metadata in DSpace is currently fairly basic; within an item, bitstreams can be arranged into separate bundles as [described above](#). A bundle may also optionally have a primary bitstream. This is currently used by the HTML support to indicate which bitstream in the bundle is the first HTML file to send to a browser.

In addition to some basic technical metadata, bitstreams also have a 'sequence ID' that uniquely identifies it within an item. This is used to produce a 'persistent' bitstream identifier for each bitstream.

Additional structural metadata can be stored in serialized bitstreams, but DSpace does not currently understand this natively.

## 2.4 Packager Plugins

Packagers are software modules that translate between DSpace Item objects and a self-contained external representation, or "package". A Package Ingestor interprets, or ingests, the package and creates an Item. A Package Disseminator writes out the contents of an Item in the package format.

A package is typically an archive file such as a "zip" or "tar" file, including a manifest document which contains metadata and a description of the package contents. The <http://www.imsglobal.org/content/packaging/> IMS Content Package is a typical packaging standard. A package might also be a single document or media file that contains its own metadata, such as a PDF document with embedded descriptive metadata.

Package ingesters and package disseminators are each a type of named plugin (see [Plugin Manager](#)), so it is easy to add new packagers specific to the needs of your site. You do not have to supply both an ingester and disseminator for each format; it is perfectly acceptable to just implement one of them.

Most packager plugins call upon Crosswalk plugins to translate the metadata between DSpace's object model and the package format.

## 2.5 Crosswalk Plugins

Crosswalks are software modules that translate between DSpace object metadata and a specific external representation. An Ingestion Crosswalk interprets the external format and crosswalks it to DSpace's internal data structure, while a Dissemination Crosswalk does the opposite.

For example, a MODS ingestion crosswalk translates descriptive metadata from the MODS format to the metadata fields on a DSpace Item. A MODS dissemination crosswalk generates a MODS document from the metadata on a DSpace Item.

Crosswalk plugins are named plugins see [Plugin Manager](#), so it is easy to add new crosswalks. You do not have to supply both an ingester and disseminator for each format; it is perfectly acceptable to just implement one of them.

There is also a special pair of crosswalk plugins which use XSL stylesheets to translate the external metadata to or from an internal DSpace format. You can add and modify XSLT crosswalks simply by editing the DSpace configuration and the stylesheets, which are stored in files in the DSpace installation directory.

The [Packager plugins](#) and [OAI-PMH](#) server make use of crosswalk plugins.

## 2.6 E-people and Groups

Although many of DSpace's functions such as document discovery and retrieval can be used anonymously, some features (and perhaps some documents) are only available to certain "privileged" users. E-People and Groups are the way DSpace identifies application users for the purpose of granting privileges. This identity is bound to a session of a DSpace application such as the Web UI or one of the command-line batch programs. Both E-People and Groups are granted privileges by the authorization system described below.

### 2.6.1 E-Person

DSpace hold the following information about each e-person:

- E-mail address
- First and last names
- Whether the user is able to log in to the system via the Web UI, and whether they must use an X509 certificate to do so;
- A password (encrypted), if appropriate
- A list of collections for which the e-person wishes to be notified of new items
- Whether the e-person 'self-registered' with the system; that is, whether the system created the e-person record automatically as a result of the end-user independently registering with the system, as opposed to the e-person record being generated from the institution's personnel database, for example.
- The network ID for the corresponding LDAP record

### 2.6.2 Groups

Groups are another kind of entity that can be granted permissions in the authorization system. A group is usually an explicit list of E-People; anyone identified as one of those E-People also gains the privileges granted to the group.

However, an application session can be assigned membership in a group without being identified as an E-Person. For example, some sites use this feature to identify users of a local network so they can read restricted

materials not open to the whole world. Sessions originating from the local network are given membership in the "LocalUsers" group and gain the corresponding privileges. Administrators can also use groups as "roles" to manage the granting of privileges more efficiently.

## 2.7 Authentication

Authentication is when an application session positively identifies itself as belonging to an E-Person and/or Group. In DSpace 1.4, it is implemented by a mechanism called Stackable Authentication: the DSpace configuration declares a "stack" of authentication methods. An application (like the Web UI) calls on the Authentication Manager, which tries each of these methods in turn to identify the E-Person to which the session belongs, as well as any extra Groups. The E-Person authentication methods are tried in turn until one succeeds. Every authenticator in the stack is given a chance to assign extra Groups. This mechanism offers the following advantages:

- Separates authentication from the Web user interface so the same authentication methods are used for other applications such as non-interactive Web Services
- Improved modularity: The authentication methods are all independent of each other. Custom authentication methods can be "stacked" on top of the default DSpace username/password method.
- Cleaner support for "implicit" authentication where username is found in the environment of a Web request, e.g. in an X.509 client certificate.

## 2.8 Authorization

DSpace's authorization system is based on associating actions with objects and the lists of EPeople who can perform them. The associations are called Resource Policies, and the lists of EPeople are called Groups. There are two special groups: 'administrators', who can do anything in a site, and 'anonymous', which is a list that contains all users. Assigning a policy for an action on an object to anonymous means giving everyone permission to do that action. (For example, most objects in DSpace sites have a policy of 'anonymous' READ.) Permissions must be explicit - lack of an explicit permission results in the default policy of 'deny'. Permissions also do not 'commute'; for example, if an e-person has READ permission on an item, they might not necessarily have READ permission on the bundles and bitstreams in that item. Currently Collections, Communities and Items are discoverable in the browse and search systems regardless of READ authorization. The following actions are possible:

### **Community**

ADD/REMOVE add or remove collections or sub-communities

### **Collection**

ADD/REMOVE add or remove items (ADD = permission to submit items)

DEFAULT\_ITEM\_READ inherited as READ by all submitted items

DEFAULT\_BITSTREAM\_READ inherited as READ by bitstreams of all submitted items

COLLECTION\_ADMIN collection admins can edit items in a collection, withdraw items, map other items into this collection.

### **Item**

ADD/REMOVE add or remove bundles

READ can view item (item metadata is always viewable)

WRITE can modify item

**Bundle** ADD/REMOVE add or remove bitstreams to a bundle

### **Bitstream**



READ view bitstream

WRITE modify bitstream

Note that there is no 'DELETE' action. In order to 'delete' an object (e.g. an item) from the archive, one must have REMOVE permission on all objects (in this case, collection) that contain it. The 'orphaned' item is automatically deleted.

Policies can apply to individual e-people or groups of e-people.

## 2.9 Ingest Process and Workflow

Rather than being a single subsystem, ingesting is a process that spans several. Below is a simple illustration of the current ingesting process in DSpace.

The batch item importer is an application, which turns an external SIP (an XML metadata document with

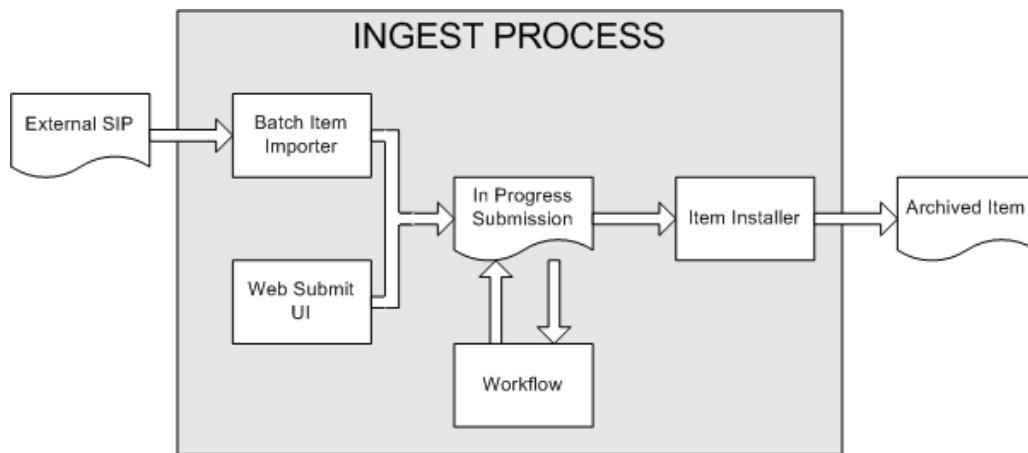


Figure 2.2: Ingest Process

some content files) into an "in progress submission" object. The Web submission UI is similarly used by an end-user to assemble an "in progress submission" object.

Depending on the policy of the collection to which the submission is targeted, a workflow process may be started. This typically allows one or more human reviewers or 'gatekeepers' to check over the submission and ensure it is suitable for inclusion in the collection.

When the Batch Ingester or Web Submit UI completes the InProgressSubmission object, and invokes the next stage of ingest (be that workflow or item installation), a provenance message is added to the Dublin Core which includes the filenames and checksums of the content of the submission. Likewise, each time a workflow changes state (e.g. a reviewer accepts the submission), a similar provenance statement is added. This allows us to track how the item has changed since a user submitted it. (The **History system** is also invoked, but provenance is easier for us to access at the moment.)

Once any workflow process is successfully and positively completed, the InProgressSubmission object is consumed by an "item installer", that converts the InProgressSubmission into a fully blown archived item in DSpace. The item installer:

- Assigns an accession date
- Adds a "date.available" value to the Dublin Core metadata record of the item
- Adds an issue date if none already present

- Adds a provenance message (including bitstream checksums)
- Assigns a Handle persistent identifier
- Adds the item to the target collection, and adds appropriate authorization policies
- Adds the new item to the search and browse indices

### 2.9.1 Workflow Steps

A collection's workflow can have up to three steps. Each collection may have an associated e-person group for performing each step; if no group is associated with a certain step, that step is skipped. If a collection has no e-person groups associated with any step, submissions to that collection are installed straight into the main archive.

In other words, the sequence is this: The collection receives a submission. If the collection has a group assigned for workflow step 1, that step is invoked, and the group is notified. Otherwise, workflow step 1 is skipped. Likewise, workflow steps 2 and 3 are performed if and only if the collection has a group assigned to those steps. When a step is invoked, the task of performing that workflow step put in the 'task pool' of the associated group. One member of that group takes the task from the pool, and it is then removed from the task pool, to avoid the situation where several people in the group may be performing the same task without realizing it.

The member of the group who has taken the task from the pool may then perform one of three actions:

1. Can accept submission for inclusion, or reject submission.
2. Can edit metadata provided by the user with the submission, but cannot change the submitted files. Can accept submission for inclusion, or reject submission.
3. Can edit metadata provided by the user with the submission, but cannot change the submitted files. Must then commit to archive; may not reject submission.

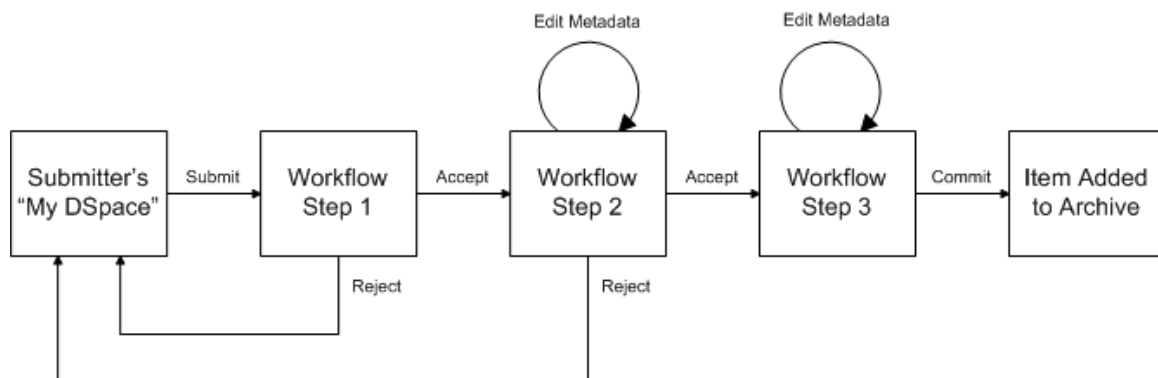


Figure 2.3: Submission Workflow in DSpace

If a submission is rejected, the reason (entered by the workflow participant) is e-mailed to the submitter, and it is returned to the submitter's 'My DSpace' page. The submitter can then make any necessary modifications and re-submit, whereupon the process starts again.

If a submission is 'accepted', it is passed to the next step in the workflow. If there are no more workflow steps with associated groups, the submission is installed in the main archive.

One last possibility is that a workflow can be 'aborted' by a DSpace site administrator. This is accomplished

using the administration UI.

The reason for this apparently arbitrary design is that it was the simplest case that covered the needs of the early adopter communities at MIT. The functionality of the workflow system will no doubt be extended in the future.

## 2.10 Supervision and Collaboration

In order to facilitate, as a primary objective, the opportunity for thesis authors to be supervised in the preparation of their e-thesis, a supervision order system exists to bind groups of other users (thesis supervisors) to an item in someone's pre-submission workspace. The bound group can have system policies associated with it that allow different levels of interaction with the student's item; a small set of default policy groups are provided:

- Full editorial control
- View item contents
- No policies

Once the default set has been applied, a system administrator may modify them as they would any other policy set in DSpace.

This functionality could also be used in situations where researchers wish to collaborate on a particular submission, although there is no particular collaborative workspace functionality.

## 2.11 Handles

Researchers require a stable point of reference for their works. The simple evolution from sharing of citations to emailing of URLs broke when Web users learned that sites can disappear or be reconfigured without notice, and that their bookmark files containing critical links to research results couldn't be trusted long term. To help solve this problem, a core DSpace feature is the creation of persistent identifier for every item, collection and community stored in DSpace. To persist identifier, DSpace requires a storage- and location- independent mechanism for creating and maintaining identifiers. DSpace uses the CNRI Handle System for creating these identifiers. The rest of this section assumes a basic familiarity with the Handle system.

DSpace uses Handles primarily as a means of assigning globally unique identifiers to objects. Each site running DSpace needs to obtain a Handle 'prefix' from [CNRI](#), so we know that if we create identifiers with that prefix, they won't clash with identifiers created elsewhere.

Presently, Handles are assigned to communities, collections, and items. Bundles and bitstreams are not assigned Handles, since over time, the way in which an item is encoded as bits may change, in order to allow access with future technologies and devices. Older versions may be moved to off-line storage as a new standard becomes de facto. Since it's usually the item that is being preserved, rather than the particular bit encoding, it only makes sense to persistently identify and allow access to the item, and allow users to access the appropriate bit encoding from there.

Of course, it may be that a particular bit encoding of a file is explicitly being preserved; in this case, the bitstream could be the only one in the item, and the item's Handle would then essentially refer just to that bitstream. The same bitstream can also be included in other items, and thus would be citable as part of a greater item, or individually.

The Handle system also features a global resolution infrastructure; that is, an end-user can enter a Handle into any service (e.g. Web page) that can resolve Handles, and the end-user will be directed to the object (in the case of DSpace, community, collection or item) identified by that Handle. In order to take advantage of this feature of the Handle system, a DSpace site must also run a 'Handle server' that can accept and resolve

incoming resolution requests. All the code for this is included in the DSpace source code bundle. Handles can be written in two forms:

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

The above represent the same Handle. The first is possibly more convenient to use only as an identifier; however, by using the second form, any Web browser becomes capable of resolving Handles. An end-user need only access this form of the Handle as they would any other URL. It is possible to enable some browsers to resolve the first form of Handle as if they were standard URLs using [CNRI's Handle Resolver plug-in](#), but since the first form can always be simply derived from the second, DSpace displays Handles in the second form, so that it is more useful for end-users.

It is important to note that DSpace uses the CNRI Handle infrastructure only at the 'site' level. For example, in the above example, the DSpace site has been assigned the prefix '1721.123'. It is still the responsibility of the DSpace site to maintain the association between a full Handle (including the '4567' local part) and the community, collection or item in question.

## 2.12 Bitstream 'Persistent' Identifiers

As of DSpace 1.2, bitstreams in DSpace also have more persistent identifiers. They are more volatile than Handles, since if the content is moved to a different server or organization, they will no longer work (hence the quotes around 'persistent'). However, they are more easily persisted than the simple URLs based on database primary key previously used. This means that external systems can more reliably refer to specific bitstreams stored in a DSpace instance.

Each bitstream has a sequence ID, unique within an item. This sequence ID is used to create a persistent ID, of the form:

```
dspace url/bitstream/handle/sequence ID/filename
```

For example:

```
https://dspace.myu.edu/bitstream/123.456/789/24/foo.html
```

The above refers to the bitstream with sequence ID 24 in the item with the Handle hdl:123.456/789. The foo.html is really just there as a hint to browsers: Although DSpace will provide the appropriate MIME type, some browsers only function correctly if the file has an expected extension.

## 2.13 Storage Resource Broker (SRB) Support

DSpace offers two means for storing bitstreams. The first is in the file system on the server. The second is using SRB (Storage Resource Broker). Both are achieved using a simple, lightweight API.

SRB is purely an option but may be used in lieu of the server's file system or in addition to the file system. Without going into a full description, SRB is a very robust, sophisticated storage manager that offers essentially unlimited storage and straightforward means to replicate (in simple terms, backup) the content on other local or remote storage resources.

SRB is purely an option but may be used in lieu of the server's file system or in addition to the file system. Without going into a full description, SRB is a very robust, sophisticated storage manager that offers essentially unlimited storage and straightforward means to replicate (in simple terms, backup) the content on other local or remote storage resources.

## 2.14 Search and Browse

DSpace allows end-users to discover content in a number of ways, including:

- Via external reference, such as a Handle
- Searching for one or more keywords in metadata or extracted full-text
- Browsing through title, author, date or subject indices, with optional image thumbnails

Search is an essential component of discovery in DSpace. Users' expectations from a search engine are quite high, so a goal for DSpace is to supply as many search features as possible. DSpace's indexing and search module has a very simple API which allows for indexing new content, regenerating the index, and performing searches on the entire corpus, a community, or collection. Behind the API is the Java freeware search engine **Lucene**. Lucene gives us fielded searching, stop word removal, stemming, and the ability to incrementally add new indexed content without regenerating the entire index. As of DSpace 1.2.1 the Lucene search indexes are configurable, enabling institutions to customise which DSpace metadata fields are indexed.

Another important mechanism for discovery in DSpace is the browse. This is the process whereby the user views a particular index, such as the title index, and navigates around it in search of interesting items. The browse subsystem provides a simple API for achieving this by allowing a caller to specify an index, and a subsection of that index. The browse subsystem then discloses the portion of the index of interest. Indices that may be browsed are item title, item issue date, item author, and subject terms. Additionally, the browse can be limited to items within a particular collection or community.

## 2.15 HTML Support

For the most part, at present DSpace simply supports uploading and downloading of bitstreams as-is. This is fine for the majority of commonly-used file formats – for example PDFs, Microsoft Word documents, spreadsheets and so forth. HTML documents (Web sites and Web pages) are far more complicated, and this has important ramifications when it comes to digital preservation:

- Web pages tend to consist of several files – one or more HTML files that contain references to each other, and stylesheets and image files that are referenced by the HTML files.
- Web pages also link to or include content from other sites, often imperceptibly to the end-user. Thus, in a few year's time, when someone views the preserved Web site, they will probably find that many links are now broken or refer to other sites than are now out of context.  
In fact, it may be unclear to an end-user when they are viewing content stored in DSpace and when they are seeing content included from another site, or have navigated to a page that is not stored in DSpace. This problem can manifest when a submitter uploads some HTML content. For example, the HTML document may include an image from an external Web site, or even their local hard drive. When the submitter views the HTML in DSpace, their browser is able to use the reference in the HTML to retrieve the appropriate image, and so to the submitter, the whole HTML document appears to have been deposited correctly. However, later on, when another user tries to view that HTML, their browser might not be able to retrieve the included image since it may have been removed from the external server. Hence the HTML will seem broken.
- Often Web pages are produced dynamically by software running on the Web server, and represent the state of a changing database underneath it.

Dealing with these issues is the topic of much active research. Currently, DSpace bites off a small, tractable chunk of this problem. DSpace can store and provide on-line browsing capability for self-contained, non-dynamic HTML documents. In practical terms, this means:

- No dynamic content (CGI scripts and so forth)
- All links to preserved content must be relative links, that do not refer to 'parents':
  - diagram.gif is OK
  - image/foo.gif is OK
  - /stylesheet.css is not OK
  - http://somedomain.com/content.html is not OK
- Any 'absolute links' (e.g. <http://somedomain.com/content.html>) are stored 'as is', and will continue to link to the external content (as opposed to relative links, which will link to the copy of the content stored in DSpace.) Thus, over time, the content referred to by the absolute link may change or disappear.

## 2.16 OAI Support

The [Open Archives Initiative](#) has developed a [protocol for metadata harvesting](#). This allows sites to programmatically retrieve or 'harvest' the metadata from several sources, and offer services using that metadata, such as indexing or linking services. Such a service could allow users to access information from a large number of sites from one place.

DSpace exposes the Dublin Core metadata for items that are publicly (anonymously) accessible. Additionally, the collection structure is also exposed via the OAI protocol's 'sets' mechanism. OCLC's open source [OAICat](#) framework is used to provide this functionality.

You can also configure the OAI service to make use of any [crosswalk plugin](#) to offer additional metadata formats, such as MODS.

DSpace's OAI service does support the exposing of deletion information for withdrawn items, but not for items that are 'expunged' (see above). DSpace also supports OAI-PMH resumption tokens.

## 2.17 OpenURL Support

DSpace supports the [OpenURL](#) protocol from [SFX](#), in a rather simple fashion. If your institution has an SFX server, DSpace will display an OpenURL link on every item page, automatically using the Dublin Core metadata. Additionally, DSpace can respond to incoming OpenURLs. Presently it simply passes the information in the OpenURL to the search subsystem. A list of results is then displayed, which usually gives the relevant item (if it is in DSpace) at the top of the list.

## 2.18 Creative Commons Support

DSpace provides support for Creative Commons licenses to be attached to items in the repository. They represent an alternative to traditional copyright. To learn more about Creative Commons, visit [their website](#). Support for the licenses is controlled by a site-wide configuration option, and since license selection involves redirection to the Creative Commons website, additional parameters may be configured to work with a proxy server. If the option is enabled, users may select a Creative Commons license during the submission process, or elect to skip Creative Commons licensing. If a selection is made a copy of the license text and RDF metadata is stored along with the item in the repository. There is also an indication - text and a Creative Commons icon - in the item display page of the web user interface when an item is licensed under Creative Commons.

## 2.19 Subscriptions

As **noted above**, end-users (e-people) may 'subscribe' to collections in order to be alerted when new items appear in those collections. Each day, end-users who are subscribed to one or more collections will receive an e-mail giving brief details of all new items that appeared in any of those collections the previous day. If no new items appeared in any of the subscribed collections, no e-mail is sent. Users can unsubscribe themselves at any time. RSS feeds of new items are also available for collections and communities.

## 2.20 History

While provenance information in the form of prose is very useful, it is not easily programmatically manipulated. The History system captures a time-based record of significant changes in DSpace, in a manner suitable for later 'refactoring' or repurposing.

Currently, the History subsystem is explicitly invoked when significant events occur (e.g., DSpace accepts an item into the archive). The History subsystem then creates RDF data describing the current state of the object. The RDF data is modeled using **Harmony/ABC**, an ontology for describing temporal-based data, and stored in the file system. Some simple indices for unwinding the data are available.

## 2.21 Import and Export

DSpace also includes batch tools to import and export items in a simple directory structure, where the Dublin Core metadata is stored in an XML file. This may be used as the basis for moving content between DSpace and other systems.

There is also a METS-based export tool, which exports items as METS-based metadata with associated bitstreams referenced from the METS file.

## 2.22 Registration

Registration is an alternate means of incorporating items, their metadata, and their bitstreams into DSpace by taking advantage of the bitstreams already being in accessible computer storage. An example might be that there is a repository for existing digital assets. Rather than using the normal interactive **ingest process** or the **batch import** to furnish DSpace the metadata and to upload bitstreams, registration provides DSpace the metadata and the location of the bitstreams. DSpace uses a variation of the import tool to accomplish registration.

## 2.23 Statistical Reports

Various statistical report about the contents and use of your system can be automatically generated by the system. These are generated by analysing DSpace's log files. Statistics can be broken down monthly. The report includes data such as:

- A customisable general summary of activities in the archive, by default including:
  - Number of item views
  - Number of collection visits
  - Number of community visits

– Number of OAI Requests

- Customisable summary of archive contents
- Broken-down list of item viewings
- A full break-down of all system activity
- User logins
- Most popular searches

The results of statistical analysis can be presented on a by-month and an in-total report, and are available via the user interface. The reports can also either be made public or restricted to administrator access only.

## 2.24 Checksum Checker

The purpose of the checker is to verify that the content in a DSpace repository has not become corrupted or been tampered with. The functionality can be invoked on an ad-hoc basis from the command line, or configured via cron or similar. Options exist to support large repositories that cannot be entirely checked in one run of the tool. The tool is extensible to new reporting and checking priority approaches.



# Chapter 3

## Installation

### 3.1 Prerequisite Software

The list below describes the third-party components and tools you'll need to run a DSpace server. These are simply recommendations based on our setup at MIT; since DSpace is built on open source, standards-based tools, there are numerous other possibilities and setups.

Also, please note that the configuration and installation guidelines relating to a particular tool below are here for convenience. You should refer to the documentation for each individual component for complete and up-to-date details. Many of the tools are updated on a frequent basis, and the guidelines below may become out of date.

1. UNIX-like OS (Linux, HP/UX etc)
2. **Java 1.4** or later (standard SDK is fine, you don't need J2EE)
3. **Apache Ant 1.6.2** or later (Java make-like tool)
4. **PostgreSQL 7.3** or later, an open source relational database, or **Oracle 9** or higher.

- **PostgreSQL**

Unicode (specifically UTF-8) support must be enabled. This is enabled by default in 8.0+. For 7.x, be sure to compile with the following options to the 'configure' script:

```
enable-multibyte --enable-unicode --with-java
```

Once installed, you need to enable TCP/IP connections (DSpace uses JDBC). For 7.x, edit postgresql.conf (usually in /usr/local/pgsql/data or /var/lib/pgsql/data), and add this line:

```
tcpip_socket = true
```

For 8.0+, in postgresql.conf uncomment the line starting:

```
listen_addresses = 'localhost'
```

Then tighten up security a bit by editing pg\_hba.conf and adding this line:

```
host dspace dspace 127.0.0.1 255.255.255.255 md5
```

Then restart PostgreSQL.

- **Oracle**

You will need to create a database for DSpace. Make sure that the character set is one of the Unicode character sets. DSpace uses UTF-8 natively, and it is suggested that the Oracle database use the same character set. You will also need to create a user account for DSpace (e.g. dspace,) and ensure that it has permissions to add and remove tables in the database. Refer to the Quick Installation for

more details.

**NOTE:** DSpace uses sequences to generate unique object IDs - beware Oracle sequences, which are said to lose their values when doing a database export/import, say restoring from a backup. Be sure to run the script `etc/update-sequences.sql`.

**ALSO NOTE:** Everything is fully functional, although Oracle limits you to 4k of text in text fields such as item metadata or collection descriptions.

For people interested in switching from Postgres to Oracle, I know of no tools that would do this automatically. You will need to recreate the community, collection, and eperson structure in the Oracle system, and then use the item export and import tools to move your content over.

5. **Jakarta Tomcat 4.x/5.x** or equivalent, such as **Jetty** or **Caucho Resin**.

Note that DSpace will need to run as the same user as Tomcat, so you might want to install and run Tomcat as a user called 'dspace'.

You need to ensure that Tomcat has a) enough memory to run DSpace and b) uses UTF-8 as its default file encoding for international character support. So ensure in your startup scripts (etc) that the following environment variable is set:

```
JAVA_OPTS="-Xmx512M -Xms64M -Dfile.encoding=UTF-8"
```

You also need to alter Tomcat's default configuration to support searching and browsing of multi-byte UTF-8 correctly. You need to add a configuration option to the `<Connector>` element in `[tomcat]/config/server.xml`:

```
URIEncoding="UTF-8"
```

e.g. if you're using the default Tomcat config, it should read:

```
<!-- Define a non-SSL HTTP/1.1 Connector on port 8080 -->
<Connector port="8080"
  maxThreads="150" minSpareThreads="25" maxSpareThreads="75"
  enableLookups="false" redirectPort="8443" acceptCount="100"
  connectionTimeout="20000" disableUploadTimeout="true"
  URIEncoding="UTF-8" />
```

Jetty and Resin are configured for correct handling of UTF-8 by default.

## 3.2 Quick Installation Steps

### But First, a Word on Directories and Path Names

DSpace uses three separate directory trees. Although you don't need to know all the details of them in order to install DSpace, you do need to know they exist and also know how they're referred to in this document:

- the source directory, referred to as `[dspace-source]`
- the install directory, referred to as `[dspace]`
- the web deployment directory. If you're using Tomcat, this will be `[tomcat]/webapps/dspace` (with `[tomcat]` being wherever you installed Tomcat—also known as `$CATALINA_HOME`). This directory is generated by the web server when it unpacks `dspace.war`, and should never be edited.

For details on the contents of these separate directory trees, refer to `directories.html`.

**Note that the source directory and install directory should always be separate!**

1. Create the DSpace user. This needs to be the same user that Tomcat (or Jetty etc) will run as. e.g. as root run:

```
useradd -m dspace
```

2. Download the [latest DSpace source code release](#) and unpack it:

```
gunzip -c dspace-source-1.x.tar.gz | tar -xf -
```

3. Database Setup

**Postgres** `quickInstallationPostgres]`

- (a) Copy the PostgreSQL JDBC driver (.jar file) into `[dspace-source]/lib`. If you compiled PostgreSQL yourself, it'll be in `postgresql-7.x.x/src/interfaces/jdbc/jars/postgresql.jar`. Alternatively you can download it directly from the [PostgreSQL JDBC](#) site. Make sure you get the driver for the version of PostgreSQL you're running and for JDBC2.
- (b) Create a dspace database, owned by the dspace PostgreSQL user:

```
createuser -U postgres -d -A -P dspace
createdb -U dspace -E UNICODE dspace
```

Enter a password for the DSpace database. (This isn't the same as the dspace user's UNIX password.)

**Oracle** (a) Copy the Oracle JDBC driver into `[dspace-source]/lib`.

- (b) Create a database for DSpace. Make sure that the character set is one of the Unicode character sets. DSpace uses UTF-8 natively, and it is suggested that the Oracle database use the same character set. Create a user account for DSpace (e.g. dspace,) and ensure that it has permissions to add and remove tables in the database.

- (c) Edit the `config/dspace.cfg` file in your source directory for the following settings:

```
db.name      = oracle
db.url       = jdbc:oracle:thin:@//host:port/dspace
db.driver    = oracle.jdbc.OracleDriver
```

- (d) Go to `[dspace-source]/etc/oracle` and copy the contents to their parent directory, overwriting the versions in the parent:

```
cd dspace_source/etc/oracle
cp * ..
```

You now have Oracle-specific .sql files in your etc directory, and your `dspace.cfg` is modified to point to your Oracle database.

4. Edit `[dspace-source]/config/dspace.cfg`, in particular you'll need to set these properties:

```
dspace.dir -- must be set to the [dspace] (installation) directory.
dspace.url -- complete URL of this server's DSpace home page.
dspace.hostname -- fully-qualified domain name of web server.
dspace.name -- "Proper" name of your server, e.g. "My Digital Library".
db.password -- the database password you entered in the previous step.
mail.server -- fully-qualified domain name of your outgoing mail server.
```

```
mail.from.address -- the "From:" address to put on email sent by DSpace.  
feedback.recipient -- mailbox for feedback mail.  
mail.admin -- mailbox for DSpace site administrator.  
alert.recipient -- mailbox for server errors/alerts (not essential but very useful!)
```

**Note:** You can interpolate the value of one configuration variable in the value of another one. For example, to set `feedback.recipient` to the same value as `mail.admin`, the line would look like:

```
feedback.recipient = ${mail.admin}
```

See the `dspace.cfg` file for examples.

5. Create the directory for the DSpace installation. As root, run:

```
mkdir [dspace]  
chown dspace [dspace]
```

(Assuming the `dspace` UNIX username.)

6. As the `dspace` UNIX user, compile and install DSpace:

```
cd [dspace-source]  
ant fresh_install
```

The most likely thing to go wrong here is the database connection. See the [common problems](#) section.

7. Copy the DSpace Web application archives (`.war` files) to the appropriate directory in your Tomcat/Jetty/Resin installation. For example:

```
cp [dspace-source]/build/*.war [tomcat]/webapps
```

8. Create an initial administrator account:

```
[dspace]/bin/create-administrator
```

9. Now the moment of truth! Start up (or restart) Tomcat. Visit the base URL of your server, e.g. `http://dspace.myu.edu:8080/dspace`. You should see the DSpace home page. Congratulations!

In order to set up some communities and collections, you'll need to access the administration UI. To do this, append `'admin'` to your server's URL, e.g. `http://dspace.myu.edu:8080/dspace/dspace-admin`.

### 3.3 Advanced Installation

The above installation steps are sufficient to set up a test server to play around with, but there are a few other steps and options you should probably consider before deploying a DSpace production site.

### 3.3.1 'cron' Jobs

A couple of DSpace features require that a script is run regularly – the e-mail subscription feature that alerts users of new items being deposited, and the new 'media filter' tool, that generates thumbnails of images and extracts the full-text of documents for indexing.

To set these up, you just need to run the following command as the dspace UNIX user:

```
crontab -e
```

Then add the following lines:

```
# Send out subscription e-mails at 01:00 every day
0 1 * * * [dspace]/bin/sub-daily
# Run the media filter at 02:00 every day
0 2 * * * [dspace]/bin/filter-media
# Run the checksum checker at 03:00
0 3 * * * [dspace]/bin/checker -lp
# Mail the results to the sysadmin at 04:00
0 4 * * * [dspace]/bin/dsrun org.dspace.checker.DailyReportEmailer -c
```

Naturally you should change the frequencies to suit your environment.

PostgreSQL also benefits from regular 'vacuuming', which optimizes the indices and clears out any deleted data. Become the postgres UNIX user, run `crontab -e` and add (for example):

```
# Clean up the database nightly at 4.20am
20 4 * * * vacuumdb --analyze dspace > /dev/null 2>&1
```

In order that statistical reports are generated regularly and thus kept up to date you should set up the following cron jobs:

```
# Run stat analyses
0 1 * * * [dspace]/bin/stat-general
0 1 * * * [dspace]/bin/stat-monthly
0 2 * * * [dspace]/bin/stat-report-general
0 2 * * * [dspace]/bin/stat-report-monthly
```

Obviously, you should choose execution times which are most useful to you, and you should ensure that the `-report-` scripts run a short while after the analysis scripts to give them time to complete (a run of around 8 months worth of logs can take around 25 seconds to complete); the resulting reports will let you know how long analysis took and you can adjust your cron times accordingly.

Note that [Perl](#) needs to be installed in order to run the statistical reports.

For information on customising the output of this see [configuring system statistical reports](#).

### 3.3.2 DSpace over HTTPS

If your DSpace is configured to have users login with a username and password (as opposed to, say, client Web certificates), then you should consider using HTTPS. Whenever a user logs in with the Web form (e.g. `dspace.myuni.edu/dspace/password-login`) their DSpace password is exposed in plain text on the network. This is a very serious security risk since network traffic monitoring is very common, especially at universities. If the risk seems minor, then consider that your DSpace administrators also login this way and they have ultimate control over the archive.

The solution is to use HTTPS (HTTP over SSL, i.e. Secure Socket Layer, an encrypted transport), which protects your passwords against being captured. You can configure DSpace to require SSL on all "authenticated"

transactions so it only accepts passwords on SSL connections.

The following sections show how to set up the most commonly-used Java Servlet containers to support HTTP over SSL.

### To enable the HTTPS support in Tomcat 5.0:

1. **For Production use:** Follow this procedure to set up SSL on your server. Using a "real" server certificate ensures your users' browsers will accept it without complaints.

In the examples below, \$CATALINA\_BASE is the directory under which your Tomcat is installed.

- (a) Create a Java keystore for your server with the password changeit, and install your server certificate under the alias "tomcat". This assumes the certificate was put in the file server.pem:

```
$JAVA_HOME/bin/keytool -import
                        -noprompt -v
                        -storepass changeit
                        -keystore $CATALINA_BASE/conf/keystore
                        -alias tomcat
                        -file myserver.pem
```

- (b) Install the CA (Certifying Authority) certificate for the CA that granted your server cert, if necessary. This assumes the server CA certificate is in ca.pem:

```
$JAVA_HOME/bin/keytool -import \
                        -noprompt \
                        -storepass changeit \
                        -trustcacerts \
                        -keystore $CATALINA_BASE/conf/keystore \
                        -alias ServerCA -file ca.pem \
```

- (c) Optional – ONLY if you need to accept client certificates for the X.509 certificate stackable authentication module See the [configuration section](#) for instructions on enabling the X.509 authentication method. Load the keystore with the CA (certifying authority) certificates for the authorities of any clients whose certificates you wish to accept. For example, assuming the client CA certificate is in client1.pem:

```
$JAVA_HOME/bin/keytool -import \
                        -noprompt \
                        -storepass changeit \
                        -trustcacerts \
                        -keystore $CATALINA_BASE/conf/keystore \
                        -alias client1 \
                        -file client1.pem
```

- (d) Now add another Connector tag to your server.xml Tomcat configuration file, like the example below. The parts affecting or specific to SSL are shown in bold. (You may wish to change some details such as the port, pathnames, and keystore password)

```
<Connector port="8443"
           maxThreads="150" minSpareThreads="25" maxSpareThreads="75"
           enableLookups="false" disableUploadTimeout="true"
           acceptCount="100" debug="0"
           scheme="https"
```

```

secure="true"
sslProtocol="TLS"
keystoreFile="conf/keystore"
keystorePass="changeit"
clientAuth="true" - ONLY if using client X.509 certs for authentication!
truststoreFile="conf/keystore"
trustedstorePass="changeit" />

```

Also, check that the default Connector is set up to redirect "secure" requests to the same port as your SSL connector, e.g.:

```

<Connector port="8080"
    maxThreads="150" minSpareThreads="25" maxSpareThreads="75"
    enableLookups="false" redirectPort="8443"
    acceptCount="100" debug="0" />

```

## 2. Quick-and-dirty Procedure for Testing:

If you are just setting up a DSpace server for testing, or to experiment with HTTPS, then you don't need to get a real server certificate. You can create a "self-signed" certificate for testing; web browsers will issue warnings before accepting it but they will function exactly the same after that as with a "real" certificate. In the examples below, \$CATALINA\_BASE is the directory under which your Tomcat is installed.

- (a) Optional – ONLY if you don't already have a server certificate. Follow this sub-procedure to request a new, signed server certificate from your Certifying Authority (CA):

- Create a new key pair under the alias name "tomcat". When generating your key, give the Distinguished Name fields the appropriate values for your server and institution. CN should be the fully-qualified domain name of your server host. Here is an example:

```

$JAVA_HOME/bin/keytool -genkey -alias tomcat -keyalg RSA -keysize 1024 \
  -keystore $CATALINA_BASE/conf/keystore \
  -storepass changeit -validity 365 \
  -dname 'CN=dspace.myuni.edu, OU=MIT Libraries, \
  O=Massachusetts Institute of Technology, L=Cambridge, S=MA, C=US'

```

- Then, create a CSR (Certificate Signing Request) and send it to your Certifying Authority. They will send you back a signed Server Certificate. This example command creates a CSR in the file tomcat.csr

```

$JAVA_HOME/bin/keytool -keystore $CATALINA_BASE/conf/keystore \
  -storepass changeit \
  -certreq -alias tomcat \
  -v -file tomcat.csr

```

- Before importing the signed certificate, you must have the CA's certificate in your keystore as a trusted certificate. Get their certificate, and import it with a command like this (for the example mitCA.pem):

```

$JAVA_HOME/bin/keytool -keystore $CATALINA_BASE/conf/keystore \
  -storepass changeit \
  -import \
  -alias mitCA \
  -trustcacerts \
  -file mitCA.pem

```

- Finally, when you get the signed certificate from your CA, import it into the keystore with a command like the following example: (cert is in the file signed-cert.pem)

```
$JAVA_HOME/bin/keytool -keystore $CATALINA_BASE/conf/keystore \
    -storepass changeit \
    -import \
    -alias tomcat \
    -trustcacerts -file signed-cert.pem
```

Since you now have a signed server certificate in your keystore, you can, obviously, skip the next steps of installing a signed server certificate and the server CA's certificate.

- (b) Create a Java keystore for your server with the password changeit, and install your server certificate under the alias "tomcat". This assumes the certificate was put in the file server.pem:

```
$JAVA_HOME/bin/keytool -genkey \
    -alias tomcat \
    -keyalg RSA \
    -keystore $CATALINA_BASE/conf/keystore \
    -storepass changeit
```

When answering the questions to identify the certificate, be sure to respond to "First and last name" with the fully-qualified domain name of your server (e.g. test-dspace.myuni.edu). The other questions are not important.

- (c) Optional – ONLY if you need to accept client certificates for the X.509 certificate stackable authentication module See the [configuration section](#) for instructions on enabling the X.509 authentication method. Load the keystore with the CA (certifying authority) certificates for the authorities of any clients whose certificates you wish to accept. For example, assuming the client CA certificate is in client1.pem:

```
$JAVA_HOME/bin/keytool -import \
    -noprompt \
    -storepass changeit \
    -trustcacerts \
    -keystore $CATALINA_BASE/conf/keystore \
    -alias client1 \
    -file client1.pem
```

- (d) Follow the procedure in the section above to add another Connector tag, for the HTTPS port, to your server.xml file.

### To use SSL on Apache HTTPD with mod\_jk:

If you choose [Apache HTTPD](#) as your primary HTTP server, you can have it forward requests to the [Tomcat](#) servlet container via [Apache Jakarta Tomcat Connector](#). This can be configured to work over SSL as well.

First, you must configure Apache for SSL; for Apache 2.0 see [Apache SSL/TLS](#) Encryption for information about using [mod\\_ssl](#).

**If you are using X.509 Client Certificates for authentication:** add these configuration options to the appropriate httpd configuration file, e.g. ssl.conf, and be sure they are in force for the virtual host and namespace locations dedicated to DSpace:

```
## SSLVerifyClient can be "optional" or "require"
SSLVerifyClient optional
SSLVerifyDepth 10
SSLCACertificateFile path-to-your-client-CA-certificate
SSLOptions StdEnvVars ExportCertData
```



Now consult the [Apache Jakarta Tomcat Connector](#) documentation to configure the `mod_jk` (note: NOT `mod_jk2`) module. Select the AJP 1.3 connector protocol. Also follow the instructions there to configure your Tomcat server to respond to AJP.

To use SSL on Apache HTTPD with `mod_webapp` consult the DSpace 1.3.2 documentation. Apache have deprecated the `mod_webapp` connector and recommend using `mod_jk`.

To use Jetty's HTTPS support consult the documentation for the relevant tool.

### 3.3.3 The Handle Server

First a few facts to clear up some common misconceptions:

- You don't **have** to use CNRI's Handle system. At the moment, you need to change the code a little to use something else (e.g PURLs) but that should change soon.
- You'll notice that while you've been playing around with a test server, DSpace has apparently been creating handles for you looking like `hdl:123456789/24` and so forth. These aren't really Handles, since the global Handle system doesn't actually know about them, and lots of other DSpace test installs will have created the same IDs.  
They're only really Handles once you've registered a prefix with CNRI (see below) and have correctly set up the Handle server included in the DSpace distribution. This Handle server communicates with the rest of the global Handle infrastructure so that anyone that understands Handles can find the Handles your DSpace has created.

If you want to use the Handle system, you'll need to set up a Handle server. This is included with DSpace. Note that this is not required in order to evaluate DSpace; you only need one if you are running a production service. You'll need to obtain a Handle prefix from the central [CNRI Handle](#) site.

A Handle server runs as a separate process that receives TCP requests from other Handle servers, and issues resolution requests to a global server or servers if a Handle entered locally does not correspond to some local content. The Handle protocol is based on TCP, so it will need to be installed on a server that can broadcast and receive TCP on port 2641.

The Handle server code is included with the DSpace code in `[dspace-source]/lib/handle.jar`. Note: The latest version of the `handle.jar` file is not included in the release due to licensing conditions changing between the provided version and later versions. It is recommended you read the [new license conditions](#) and decide whether you wish to update your installation's `handle.jar`. If you decide to update, you should replace the existing `handle.jar` in `[dspace-source]/lib` with the new version and rebuild your war files.

A script exists to create a simple Handle configuration - simply run `[dspace]/bin/make-handle-config` after you've set the appropriate parameters in `dspace.cfg`. You can also create a Handle configuration directly by following the [installation instructions on handle.net](#), but with these changes:

- Instead of running:

```
java -cp /hs/bin/handle.jar net.handle.server.SimpleSetup /hs/svr_1
```

as directed in the [Handle Server Administration Guide](#), you should run

```
[dspace]/bin/dsrun net.handle.server.SimpleSetup [dspace]/handle-server
```

ensuring that `[dspace]/handle-server` matches whatever you have in `dspace.cfg` for the `handle.dir` property.

- Edit the resulting `[dspace]/handle-server/config.dct` file to include the following lines in the `"server_config"` clause:

```
"storage_type" = "CUSTOM"  
"storage_class" = "org.dspace.handle.HandlePlugin"
```

This tells the Handle server to get information about individual Handles from the DSpace code.

Whichever approach you take, start the Handle server with `[dspace]/bin/start-handle-server`, as the DSpace user. Once the configuration file has been generated, you will need to go to <http://hdl.handle.net/4263537/5014> to upload the generated `sitebndl.zip` file. The upload page will ask you for your contact information. An administrator will then create the naming authority/prefix on the root service (known as the Global Handle Registry), and notify you when this has been completed. You will not be able to continue the handle server installation until you receive further information concerning your naming authority.

Note that since the DSpace code manages individual Handles, administrative operations such as Handle creation and modification aren't supported by DSpace's Handle server.

If you need to update the handle prefix on items created before the CNRI registration process you can run the `[dspace]/bin/update-handle-prefix` script. You may need to do this if you loaded items prior to CNRI registration (e.g. setting up a demonstration system prior to migrating it to production). The script takes the current and new prefix as parameters. For example:

```
[dspace]/bin/update-handle-prefix 123456789 1303
```

will change any handles currently assigned prefix 123456789 to prefix 1303, so for example handle 123456789/23 will be updated to 1303/23 in the database.

## 3.4 Windows Installation

### 3.4.1 Pre-requisite Software

You'll need to install this pre-requisite software:

- **Java SDK 1.4** or later (standard SDK is fine, you don't need J2EE)
- **PostgreSQL 8.x for Windows**. This comes with an installer application now, so Cygwin is no longer required. Make sure the ODBC + JDBC options are selected, as well as the pgAdmin III tool
- **Apache Ant 1.6.x**. Unzip the package in C:  
and add C:  
`apache-ant-1.6.2`  
bin to the PATH environment variable. For Ant to work properly, you should ensure that `JAVA_HOME` is set.
- **Jakarta Tomcat 5.x+**

### 3.4.2 Installation Steps

1. Download the DSpace source from [SourceForge](#) and untar it ([WinZip](#) will do this)
2. Copy the PostgreSQL JDBC driver across to the DSpace source tree. The drivers will be in C:  
Program Files  
PostgreSQL  
8.x  
jdbc  
. The `postgresql-8.x-yyy.jdbc2.jar` is the jar file you need. Copy it to `[dspace-source]/lib`.

3. Ensure the PostgreSQL service is running, and then run pgAdmin III (Start -> PostgreSQL 8.0 -> pgAdmin III). Connect to the local database as the postgres user and:
  - Create a 'Login Role' (user) called dspace with the password dspace
  - Create a database called dspace owned by the user dspace, with UTF-8 encoding

4. Update paths in [dspace-source]  
config

dspace.cfg. Note: Use forward slashes / for path separators, though you can still use drive letters, e.g.:

```
dspace.dir = C:/DSpace
```

Make sure you change all of the parameters with file paths to suit, specifically:

```
dspace.dir
config.template.log4j.properties
config.template.log4j-handle-plugin.properties
config.template.oaicat.properties
assetstore.dir
history.dir
log.dir
upload.temp.dir
report.dir
handle.dir
```

5. Create the directory for the DSpace installation (e.g. C:  
DSpace)
6. Run:

```
ant fresh_install
```

7. Create an administrator account, e.g. assuming C:  
dspace is where your DSpace installation is:

```
C:\dspace\bin\dsrun org.dspace.administer.CreateAdministrator
```

and enter the required information

8. Copy the .war Web application files from [dspace-source]  
build to Tomcat's webapps dir, which should be somewhere like C:  
Program Files  
Apache Software Foundation  
Tomcat 5.5  
webapps

9. Start the Tomcat service

10. Browse <http://localhost:8080/dspace>. You should see the DSpace home page

### 3.5 Known Bugs

In any software project of the scale of DSpace, there will be bugs. Sometimes, a stable version of DSpace includes known bugs. We do not always wait until every known bug is fixed before a release. If the software is sufficiently stable and an improvement on the previous release, and the bugs are minor and have known workarounds, we release it to enable the community to take advantage of those improvements.

The known bugs in a release are documented in the `KNOWN_BUGS` file in the source package.

Please see the DSpace [bug tracker](#) for further information on current bugs, and to find out if the bug has subsequently been fixed. This is also where you can report any further bugs you find.

### 3.6 Common Problems

In an ideal world everyone would follow the above steps and have a fully functioning DSpace. Of course, in the real world it doesn't always seem to work out that way. This section lists common problems that people encounter when installing DSpace, and likely causes and fixes. This is likely to grow over time as we learn about users' experiences.

**Database errors occur when you run `ant fresh_install`** There are two common errors that occur. If your error looks like this—

```
[java] 2004-03-25 15:17:07,730 INFO
        org.dspace.storage.rdbms.InitializeDatabase @ Initializing Database
[java] 2004-03-25 15:17:08,816 FATAL
        org.dspace.storage.rdbms.InitializeDatabase @ Caught exception:
[java] org.postgresql.util.PSQLException: Connection refused.
        Check that the hostname and port are correct and
        that the postmaster is accepting TCP/IP connections.
[java] at org.postgresql.jdbc1.AbstractJdbc1Connection.openConnection
        (AbstractJdbc1Connection.java:204)
[java] at org.postgresql.Driver.connect(Driver.java:139)
```

it usually means you haven't yet added the relevant configuration parameter to your PostgreSQL configuration (see above), or perhaps you haven't restarted PostgreSQL after making the change. Also, make sure that the `db.username` and `db.password` properties are correctly set in `[dspace-source]/config/dspace.cfg`. An easy way to check that your DB is working OK over TCP/IP is to try this on the command line:

```
psql -U dspace -W -h localhost
```

Enter the dspace database password, and you should be dropped into the psql tool with a `dspace=>` prompt.

Another common error looks like this:

```
[java] 2004-03-25 16:37:16,757 INFO
        org.dspace.storage.rdbms.InitializeDatabase @ Initializing Database
[java] 2004-03-25 16:37:17,139 WARN
        org.dspace.storage.rdbms.DatabaseManager @ Exception initializing DB pool
[java] java.lang.ClassNotFoundException: org.postgresql.Driver
[java] at java.net.URLClassLoader$1.run(URLClassLoader.java:198)
[java] at java.security.AccessController.doPrivileged(Native Method)
[java] at java.net.URLClassLoader.findClass(URLClassLoader.java:186)
```

This means that the PostgreSQL JDBC driver is not present in [dspace-source]/lib. See [above](#).

**Tomcat doesn't shut down** If you're trying to tweak Tomcat's configuration but nothing seems to make a difference to the error you're seeing, you might find that Tomcat hasn't been shutting down properly, perhaps because it's waiting for a stale connection to close gracefully which won't happen. To see if this is the case, try:

```
ps -ef | grep java
```

and look for Tomcat's Java processes. If they stay around after running Tomcat's shutdown.sh script, trying killing them (with -9 if necessary), then starting Tomcat again.

**Database connections don't work, or accessing DSpace takes forever** If you find that when you try to access a DSpace Web page and your browser sits there connecting, or if the database connections fail, you might find that a 'zombie' database connection is hanging around preventing normal operation. To see if this is the case, try:

```
ps -ef | grep postgres
```

You might see some processes like this

```
dspace 16325 1997 0 Feb 14 ?
      0:00 postgres: dspace dspace 127.0.0.1 idle in transaction
```

This is normal—DSpace maintains a 'pool' of open database connections, which are re-used to avoid the overhead of constantly opening and closing connections. If they're 'idle' it's OK; they're waiting to be used. However sometimes, if something went wrong, they might be stuck in the middle of a query, which seems to prevent other connections from operating, e.g.:

```
dspace 16325 1997 0 Feb 14 ?
      0:00 postgres: dspace dspace 127.0.0.1 SELECT
```

This means the connection is in the middle of a SELECT operation, and if you're not using DSpace right that instant, it's probably a 'zombie' connection. If this is the case, try killing the process, and stopping and restarting Tomcat.

**You've made changes to the code** or to the JSP's and rebuilt DSpace successfully, but when you run Tomcat you don't see any of your changes in DSpace. After you've rebuilt DSpace and copied dspace.war from your [dspace-source]/build directory into your [tomcat]/webapps directory, you must also delete the existing [tomcat]/webapps/dspace directory before re-starting Tomcat. Otherwise Tomcat will continue to use the old code.

## Chapter 4

# Updating a DSpace Installation

This section describes how to update a DSpace installation from one version to the next. Details of the differences between the functionality of each version are given in the [Version History](#) section.

### 4.1 Updating From 1.4 to 1.4.x

The changes in 1.4.1 are only code changes so the update is simply a matter of rebuilding the wars. In the notes below [dspace] refers to the install directory for your existing DSpace installation, and [dspace-1.4.1-source] to the source directory for DSpace 1.4.1. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Get the new DSpace 1.4.1 source code from the DSpace page on [SourceForge](#) and unpack it somewhere. Do not unpack it on top of your existing installation!!
2. Copy the PostgreSQL driver JAR to the source tree. For example:

```
cd [dspace]/lib
cp postgresql.jar [dspace-1.4.1-source]/lib
```

3. **Note:** Licensing conditions for the handle.jar file have changed. As a result, the latest version of the handle.jar file is not included in this distribution. It is recommended you read the [new license conditions](#) and decide whether you wish to update your installation's handle.jar. If you decide to update, you should replace the existing handle.jar in [dspace-1.4.1-source]/lib with the new version.
4. Take down Tomcat (or whichever servlet container you're using).
5. Your 'localized' JSPs (those in jsp/local) now need to be maintained in the source directory. If you have locally modified JSPs in your [dspace]/jsp/local directory, you will need to merge the changes in the new 1.4.1 versions into your locally modified ones. You can use the diff command to compare your JSPs against the 1.4.1 versions to do this. You can also check against the [DSpace CVS](#).
6. In [dspace-1.4.1-source] run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

7. Copy the .war Web application files in [dspace-1.4.1-source]/build to the webapps sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dSPACE-1.4.1-source]/build/*.war [tomcat]/webapps
```

If you're using Tomcat, you need to delete the directories corresponding to the old .war files. For example, if dspace.war is installed in [tomcat]/webapps/dspace.war, you should delete the [tomcat]/webapps/dspace directory. Otherwise, Tomcat will continue to use the old code in that directory.

8. Restart Tomcat.

## 4.2 Updating From 1.3.x to 1.4.x

1. First and foremost, make a complete backup of your system, including:
  - A snapshot of the database
  - The asset store ([dSPACE]/assetstore by default)
  - Your configuration files and localized JSPs

2. Download the latest **DSpace 1.4.x** source bundle and unpack it in a suitable location (not over your existing DSpace installation or source tree!)

3. Copy the PostgreSQL driver JAR to the source tree. For example

```
cd [dSPACE]/lib
cp postgresql.jar [dSPACE-1.4.x-source]/lib
```

4. **Note:** Licensing conditions for the handle.jar file have changed. As a result, the latest version of the handle.jar file is not included in this distribution. It is recommended you read the [new license conditions](#) and decide whether you wish to update your installation's handle.jar. If you decide to update, you should replace the existing handle.jar in [dSPACE-1.4.1-source]/lib with the new version.

5. Take down Tomcat (or whichever servlet container you're using).

6. Your DSpace configuration will need some updating:

- In dspace.cfg, paste in the following lines for the new stackable authentication feature, the new method for managing Media Filters, and the Checksum Checker.

```
##### Stackable Authentication Methods #####
# Stack of authentication methods
# (See org.dspace.eperson.AuthenticationManager)
plugin.sequence.org.dspace.eperson.AuthenticationMethod = \
    org.dspace.eperson.PasswordAuthentication

##### Example of configuring X.509 authentication
##### (to use it, add org.dspace.eperson.X509Authentication to stack)
## method 1, using keystore
#authentication.x509.keystore.path = /var/local/tomcat/conf/keystore
#authentication.x509.keystore.password = changeit
## method 2, using CA certificate
#authentication.x509.ca.cert = ${dSPACE.dir}/config/mitClientCA.der
## Create e-persons for unknown names in valid certificates?
#authentication.x509.autoregister = true
```

```

#### Media Filter plugins (through PluginManager) ####
plugin.sequence.org.dspace.app.mediafilter.MediaFilter = \
    org.dspace.app.mediafilter.PDFFilter, org.dspace.app.mediafilter.HTMLFilter, \
    org.dspace.app.mediafilter.WordFilter, org.dspace.app.mediafilter.JPEGFilter
# to enable branded preview: remove last line above, and uncomment 2 lines below
# org.dspace.app.mediafilter.WordFilter, org.dspace.app.mediafilter.JPEGFilter, \
# org.dspace.app.mediafilter.BrandedPreviewJPEGFilter
filter.org.dspace.app.mediafilter.PDFFilter.inputFormats = Adobe PDF
filter.org.dspace.app.mediafilter.HTMLFilter.inputFormats = HTML, Text
filter.org.dspace.app.mediafilter.WordFilter.inputFormats = Microsoft Word
filter.org.dspace.app.mediafilter.JPEGFilter.inputFormats = GIF, JPEG, image/png
filter.org.dspace.app.mediafilter.BrandedPreviewJPEGFilter.inputFormats = \
GIF, JPEG, image/png

#### Settings for Item Preview ####
webui.preview.enabled = false
# max dimensions of the preview image
webui.preview.maxwidth = 600
webui.preview.maxheight = 600
# the brand text
webui.preview.brand = My Institution Name
# an abbreviated form of the above text, this will be used
# when the preview image cannot fit the normal text
webui.preview.brand.abbrev = MyOrg
# the height of the brand
webui.preview.brand.height = 20
# font settings for the brand text
webui.preview.brand.font = SansSerif
webui.preview.brand.fontpoint = 12
#webui.preview.dc = rights

#### Checksum Checker Settings ####
# Default dispatcher in case none specified
plugin.single.org.dspace.checker.BitstreamDispatcher=org.dspace.checker.SimpleDispatcher
# Standard interface implementations. You shouldn't need to tinker with these.
plugin.single.org.dspace.checker.ReporterDAO=org.dspace.checker.ReporterDAOImpl
# check history retention
checker.retention.default=10y
checker.retention.CHECKSUM_MATCH=8w

```

- If you have customised advanced search fields (search.index.n fields, note that you now need to include the schema in the values. Dublin Core is specified as dc. So for example, if in 1.3.2 you had:

```
search.index.1 = title:title.alternative
```

That needs to be changed to:

```
search.index.1 = title:dc.title.alternative
```

- If you use LDAP or X509 authentication, you'll need to add org.dspace.eperson.LDAPAuthentication or org.dspace.eperson.X509Authentication respectively. See also [configuring custom authentication code](#).



- If you have custom Media Filters, note that these are now configured through `dspace.cfg` (instead of `mediafilter.cfg` which is obsolete.)
- Also, take a look through the default `dspace.cfg` file supplied with DSpace 1.4.x, as this contains configuration options for various new features you might like to use. In general, these new features default to 'off' and you'll need to add configuration properties as described in the default 1.4.x `dspace.cfg` to activate them.

7. Your 'localized' JSPs (those in `jsp/local`) now need to be maintained in the source directory. If you have locally modified JSPs in your `[dspace]/jsp/local` directory, you will need to merge the changes in the new 1.4.x versions into your locally modified ones. You can use the `diff` command to compare your JSPs against the 1.4.x versions to do this. You can also check against the [DSpace CVS](#).

8. In `[dspace-1.4.x-source]` run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

9. The database schema needs updating. SQL files containing the relevant file are provided. If you've modified the schema locally, you may need to check over this and make alterations.

**For PostgreSQL** `[dspace-1.4.x-source]/etc/database_schema_13-14.sql` contains the SQL commands to achieve this for PostgreSQL. To apply the changes, go to the source directory, and run:

```
psql -f etc/database_schema_13-14.sql [DSpace database name] -h localhost
```

**For Oracle** `[dspace-1.4.x-source]/etc/oracle/database_schema_13-14.sql` should be run on the DSpace database to update the schema.

10. Rebuild the search indices:

```
[dspace]/bin/index-all
```

11. Copy the `.war` Web application files in `[dspace-1.4-source]/build` to the `webapps` sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dspace-1.4-source]/build/*.war [tomcat]/webapps
```

If you're using Tomcat, you need to delete the directories corresponding to the old `.war` files. For example, if `dspace.war` is installed in `[tomcat]/webapps/dspace.war`, you should delete the `[tomcat]/webapps/dspace` directory. Otherwise, Tomcat will continue to use the old code in that directory.

12. Restart Tomcat.

### 4.3 Updating From 1.3.1 to 1.3.2

The changes in 1.3.2 are only code changes so the update is simply a matter of rebuilding the wars.

In the notes below `[dspace]` refers to the install directory for your existing DSpace installation, and `[dspace-1.3.2-source]` to the source directory for DSpace 1.3.2. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Get the new DSpace 1.3.2 source code from the DSpace page on SourceForge and unpack it somewhere. Do not unpack it on top of your existing installation!!

2. Copy the PostgreSQL driver JAR to the source tree. For example:

```
cd [dspace]/lib
cp postgresql.jar [dspace-1.3.2-source]/lib
```

3. Take down Tomcat (or whichever servlet container you're using).
4. Your 'localized' JSPs (those in jsp/local) now need to be maintained in the source directory. If you have locally modified JSPs in your [dspace]/jsp/local directory, you will need to merge the changes in the new 1.3.2 versions into your locally modified ones. You can use the diff command to compare the 1.3.1 and 1.3.2 versions to do this.
5. In [dspace-1.3.2-source] run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

6. Copy the .war Web application files in [dspace-1.3.2-source]/build to the webapps sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dspace-1.3.2-source]/build/*.war [tomcat]/webapps
```

If you're using Tomcat, you need to delete the directories corresponding to the old .war files. For example, if dspace.war is installed in [tomcat]/webapps/dspace.war, you should delete the [tomcat]/webapps/dspace directory. Otherwise, Tomcat will continue to use the old code in that directory.

7. Restart Tomcat.

## 4.4 Updating From 1.2.x to 1.3.x

In the notes below [dspace] refers to the install directory for your existing DSpace installation, and [dspace-1.3.x-source] to the source directory for DSpace 1.3.x. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Step one is, of course, to back up all your data before proceeding!! Include all of the contents of [dspace] and the PostgreSQL database in your backup.
2. Get the new DSpace 1.3.x source code from the DSpace page on SourceForge and unpack it somewhere. Do not unpack it on top of your existing installation!!
3. Copy the PostgreSQL driver JAR to the source tree. For example:

```
cd [dspace]/lib
cp postgresql.jar [dspace-1.2.2-source]/lib
```

4. Take down Tomcat (or whichever servlet container you're using).
5. Remove the old version of xerces.jar from your installation, so it is not inadvertently later used:

```
rm [dspace]/lib/xerces.jar
```

6. Install the new config files by moving dstat.cfg and dstat.map from [dspace-1.3.x-source]/config/ to [dspace]/config

7. You need to add new parameters to your `[dspace]/dspace.cfg`:

```
##### 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/reports/
```

8. Build and install the updated DSpace 1.3.x code. Go to the `[dspace-1.3.x-source]` directory, and run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

9. You'll need to make some changes to the database schema in your PostgreSQL database. `[dspace-1.3.x-source]/etc/database_schema_12-13.sql` contains the SQL commands to achieve this. If you've modified the schema locally, you may need to check over this and make alterations.

To apply the changes, go to the source directory, and run:

```
psql -f etc/database_schema_12-13.sql [DSpace database name] -h localhost
```

10. Customise the stat generating statistics as per the instructions in System Statistical Reports

11. Initialise the statistics using:

```
[dspace]/bin/stat-initial
[dspace]/bin/stat-general
[dspace]/bin/stat-report-initial
[dspace]/bin/stat-report-general
```

12. Rebuild the search indices:

```
[dspace]/bin/index-all
```

13. Copy the `.war` Web application files in `[dspace-1.3.x-source]/build` to the `webapps` sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dspace-1.3.x-source]/build/*.war [tomcat]/webapps
```

14. Restart Tomcat.

## 4.5 Updating From 1.2.1 to 1.2.2

The changes in 1.2.2 are only code and config changes so the update should be fairly simple.

In the notes below `[dspace]` refers to the install directory for your existing DSpace installation, and `[dspace-1.2.2-source]` to the source directory for DSpace 1.2.2. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Get the new DSpace 1.2.2 source code from the DSpace page on SourceForge and unpack it somewhere. Do not unpack it on top of your existing installation!!
2. Copy the PostgreSQL driver JAR to the source tree. For example:

```
cd [dSPACE]/lib
cp postgresql.jar [dSPACE-1.2.2-source]/lib
```

3. Take down Tomcat (or whichever servlet container you're using).
4. Your 'localized' JSPs (those in jsp/local) now need to be maintained in the source directory. If you have locally modified JSPs in your [dSPACE]/jsp/local directory, you might like to merge the changes in the new 1.2.2 versions into your locally modified ones. You can use the diff command to compare the 1.2.1 and 1.2.2 versions to do this. Also see the version history for a list of modified JSPs.
5. You need to add a new parameter to your [dSPACE]/dSPACE.cfg for configurable fulltext indexing

```
##### Fulltext Indexing settings #####
# Maximum number of terms indexed for a single field in Lucene.
# Default is 10,000 words - often not enough for full-text indexing.
# If you change this, you'll need to re-index for the change
# to take effect on previously added items.
# -1 = unlimited (Integer.MAX_VALUE)
search.maxfieldlength = 10000
```

6. In [dSPACE-1.2.2-source] run:

```
ant -Dconfig=[dSPACE]/config/dSPACE.cfg update
```

7. Copy the .war Web application files in [dSPACE-1.2.2-source]/build to the webapps sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dSPACE-1.2.2-source]/build/*.war [tomcat]/webapps
```

If you're using Tomcat, you need to delete the directories corresponding to the old .war files. For example, if dSPACE.war is installed in [tomcat]/webapps/dSPACE.war, you should delete the [tomcat]/webapps/dSPACE directory. Otherwise, Tomcat will continue to use the old code in that directory.

8. To finalise the install of the new configurable submission forms you need to copy the file [dSPACE-1.2.2-source]/config/input-forms.xml into [dSPACE]/config.
9. Restart Tomcat.

## 4.6 Updating From 1.2 to 1.2.1

The changes in 1.2.1 are only code changes so the update should be fairly simple.

In the notes below [dSPACE] refers to the install directory for your existing DSpace installation, and [dSPACE-1.2.1-source] to the source directory for DSpace 1.2.1. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Get the new DSpace 1.2.1 source code from the DSpace page on SourceForge and unpack it somewhere. Do not unpack it on top of your existing installation!!

- Copy the PostgreSQL driver JAR to the source tree. For example:

```
cd [dspace]/lib
cp postgresql.jar [dspace-1.2.1-source]/lib
```

- Take down Tomcat (or whichever servlet container you're using).
- Your 'localized' JSPs (those in jsp/local) now need to be maintained in the source directory. If you have locally modified JSPs in your [dspace]/jsp/local directory, you might like to merge the changes in the new 1.2.1 versions into your locally modified ones. You can use the diff command to compare the 1.2 and 1.2.1 versions to do this. Also see the version history for a list of modified JSPs.
- You need to add a few new parameters to your [dspace]/dspace.cfg for browse/search and item thumbnails display, and for configurable DC metadata fields to be indexed.

```
# whether to display thumbnails on browse and search results pages (1.2+)
webui.browse.thumbnail.show = false

# max dimensions of the browse/search thumbs. Must be <= thumbnail.maxwidth
# and thumbnail.maxheight. Only need to be set if required to be smaller than
# dimension of thumbnails generated by mediafilter (1.2+)
#webui.browse.thumbnail.maxheight = 80
#webui.browse.thumbnail.maxwidth = 80

# whether to display the thumb against each bitstream (1.2+)
webui.item.thumbnail.show = true

# where should clicking on a thumbnail from browse/search take the user
# Only values currently supported are "item" and "bitstream"
#webui.browse.thumbnail.linkbehaviour = item

##### Fields to Index for Search #####

# DC metadata elements.qualifiers to be indexed for search
# format: - search.index.[number] = [search field]:element.qualifier
#          - * used as wildcard

###      changing these will change your search results,      ###
###      but will NOT automatically change your search displays  ###

search.index.1 = author:contributor.*
search.index.2 = author:creator.*
search.index.3 = title:title.*
search.index.4 = keyword:subject.*
search.index.5 = abstract:description.abstract
search.index.6 = author:description.statementsofresponsibility
search.index.7 = series:relation.ispartofseries
search.index.8 = abstract:description.tableofcontents
search.index.9 = mime:format.mimetype
search.index.10 = sponsor:description.sponsorship
search.index.11 = id:identifier.*
```

6. In [dspace-1.2.1-source] run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

7. Copy the .war Web application files in [dspace-1.2.1-source]/build to the webapps sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dspace-1.2.1-source]/build/*.war [tomcat]/webapps
```

If you're using Tomcat, you need to delete the directories corresponding to the old .war files. For example, if dspace.war is installed in [tomcat]/webapps/dspace.war, you should delete the [tomcat]/webapps/dspace directory. Otherwise, Tomcat will continue to use the old code in that directory.

8. Restart Tomcat.

## 4.7 Updating From 1.1 (or 1.1.1) to 1.2

The process for upgrading to 1.2 from either 1.1 or 1.1.1 is the same. If you are running DSpace 1.0 or 1.0.1, you need to follow the instructions for upgrading from 1.0.1 to 1.1 to before following these instructions.

Note also that if you've substantially modified DSpace, these instructions apply to an unmodified 1.1.1 DSpace instance, and you'll need to adapt the process to any modifications you've made.

This document refers to the install directory for your existing DSpace installation as [dspace], and to the source directory for DSpace 1.2 as [dspace-1.2-source]. Whenever you see these path references below, be sure to replace them with the actual path names on your local system.

1. Step one is, of course, to back up all your data before proceeding!! Include all of the contents of [dspace] and the PostgreSQL database in your backup.
2. Get the new DSpace 1.2 source code from the DSpace page on SourceForge and unpack it somewhere. Do not unpack it on top of your existing installation!!
3. Copy the required Java libraries that we couldn't include in the bundle to the source tree. For example:

```
cd [dspace]/lib
cp activation.jar servlet.jar mail.jar [dspace-1.2-source]/lib
```

4. Stop Tomcat (or other servlet container.)
5. It's a good idea to upgrade all of the various third-party tools that DSpace uses to their latest versions:
  - Java (note that now version 1.4.0 or later is required)
  - Tomcat (Any version after 4.0 will work; symbolic links are no longer an issue)
  - PostgreSQL (don't forget to build/download an updated JDBC driver .jar file! Also, back up the database first.)
  - Ant
6. You need to add the following new parameters to your [dspace]/dspace.cfg:

```
##### Media Filter settings #####
# maximum width and height of generated thumbnails
thumbnail.maxwidth 80
thumbnail.maxheight 80
```

There are one or two other, optional extra parameters (for controlling the pool of database connections). See the version history for details. If you leave them out, defaults will be used.

Also, to avoid future confusion, you might like to remove the following property, which is no longer required:

```
config.template.oai-web.xml = [dSPACE]/oai/WEB-INF/web.xml
```

7. The layout of the installation directory (i.e. the structure of the contents of [dSPACE]) has changed somewhat since 1.1.1. First up, your 'localized' JSPs (those in jsp/local) now need to be maintained in the source directory. So make a copy of them now!

Once you've done that, you can remove [dSPACE]/jsp and [dSPACE]/oai, these are no longer used. (.war Web application archive files are used instead).

Also, if you're using the same version of Tomcat as before, you need to remove the lines from Tomcat's conf/server.xml file that enable symbolic links for DSpace. These are the <Context> elements you added to get DSpace 1.1.1 working, looking something like this:

```
<Context path="/dSPACE" docBase="dSPACE" debug="0" reloadable="true" crossContext="true">
  <Resources className="org.apache.naming.resources.FileDirContext" allowLinking="true" />
</Context>
```

Be sure to remove the <Context> elements for both the Web UI and the OAI Web applications.

8. Build and install the updated DSpace 1.2 code. Go to the DSpace 1.2 source directory, and run:

```
ant -Dconfig=[dSPACE]/config/dSPACE.cfg update
```

9. Copy the new config files in config to your installation, e.g.:

```
cp [dSPACE-1.2-source]/config/news-* \
[dSPACE-1.2-source]/config/mediafilter.cfg \
dSPACE-1.2-source]/config/dc2mods.cfg [dSPACE]/config
```

10. You'll need to make some changes to the database schema in your PostgreSQL database. [dSPACE-1.2-source]/etc/database\_schema\_11-12.sql contains the SQL commands to achieve this. If you've modified the schema locally, you may need to check over this and make alterations.

To apply the changes, go to the source directory, and run:

```
psql -f etc/database_schema_11-12.sql [DSpace database name] -h localhost
```

11. A tool supplied with the DSpace 1.2 codebase will then update the actual data in the relational database. Run it using:

```
[dSPACE]/bin/dsrun org.dSPACE.administer.Upgrade11To12
```

12. Then rebuild the search indices:

```
[dSPACE]/bin/index-all
```

13. Delete the existing symlinks from your servlet container's (e.g. Tomcat's) webapp sub-directory.

14. Copy the .war Web application files in [dspace-1.2-source]/build to the webapps sub-directory of your servlet container (e.g. Tomcat). e.g.:

```
cp [dspace-1.2-source]/build/*.war [tomcat]/webapps
```

15. Restart Tomcat.

16. To get image thumbnails generated and full-text extracted for indexing automatically, you need to set up a 'cron' job, for example one like this:

```
# Run the media filter at 02:00 every day
0 2 * * * [dspace]/bin/filter-media
```

You might also wish to run it now to generate thumbnails and index full text for the content already in your system.

17. Note 1: This update process has effectively 'touched' all of your items. Although the dates in the Dublin Core metadata won't have changed (accession date and so forth), the 'last modified' date in the database for each will have been changed.

This means the e-mail subscription tool may be confused, thinking that all items in the archive have been deposited that day, and could thus send a rather long email to lots of subscribers. So, it is recommended that you turn off the e-mail subscription feature for the next day, by commenting out the relevant line in DSpace's cron job, and then re-activating it the next day.

Say you performed the update on 08-June-2004 (UTC), and your e-mail subscription cron job runs at 4am (UTC). When the subscription tool runs at 4am on 09-June-2004, it will find that everything in the system has a modification date in 08-June-2004, and accordingly send out huge emails. So, immediately after the update, you would edit DSpace's 'crontab' and comment out the /dspace/bin/subs-daily line. Then, after 4am on 09-June-2004 you'd 'un-comment' it out, so that things proceed normally.

Of course this means, any real new deposits on 08-June-2004 won't get e-mailed, however if you're updating the system it's likely to be down for some time so this shouldn't be a big problem.

18. Note 2: After consultation with the OAI community, various OAI-PMH changes have occurred:

- The OAI-PMH identifiers have changed (they're now of the form oai:hostname:handle as opposed to just Handles)
- The set structure has changed, due to the new sub-communities feature.
- The default base URL has changed
- As noted in note 1, every item has been 'touched' and will need re-harvesting.

The above means that, if already registered and harvested, you will need to re-register your repository, effectively as a 'new' OAI-PMH data provider. You should also consider posting an announcement to the OAI implementers e-mail list so that harvesters know to update their systems.

Also note that your site may, over the next few days, take quite a big hit from OAI-PMH harvesters. The resumption token support should alleviate this a little, but you might want to temporarily whack up the database connection pool parameters in [dspace]/config/dspace.cfg. See the dspace.cfg distributed with the source code to see what these parameters are and how to use them. (You need to stop and restart Tomcat after changing them.)

I realize this is not ideal; for discussion as to the reasons behind this please see relevant posts to the OAI community: post one, post two, as well as this post to the dspace-tech mailing list.

If you really can't live with updating the base URL like this, you can fairly easily have things proceed more-or-less as they are, by doing the following:



- Change the value of `OAI_ID_PREFIX` at the top of the `org.dspace.app.oai.DSpaceOAICatalog` class to `hdl`:
- Change the servlet mapping for the `OAIHandler` servlet back to `/` (from `/request`)
- Rebuild and deploy `dspace-oai.war`

However, note that in this case, all the records will be re-harvested by harvesters anyway, so you still need to brace for the associated DB activity; also note that the set spec changes may not be picked up by some harvesters. It's recommended you read the above-linked mailing list posts to understand why the change was made.

## 4.8 Updating From 1.1 to 1.1.1

Fortunately the changes in 1.1.1 are only code changes so the update is fairly simple.

In the notes below `[dspace]` refers to the install directory for your existing DSpace installation, and `[dspace-1.1.1-source]` to the source directory for DSpace 1.1.1. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Take down Tomcat.
2. It would be a good idea to update any of the third-party tools used by DSpace at this point (e.g. PostgreSQL), following the instructions provided with the relevant tools.
3. In `[dspace-1.1.1-source]` run:

```
ant -Dconfig=[dspace]/config/dspace.cfg update
```

4. If you have locally modified JSPs of the following JSPs in your `[dspace]/jsp/local` directory, you might like to merge the changes in the new 1.1.1 versions into your locally modified ones. You can use the `diff` command to compare the 1.1 and 1.1.1 versions to do this. The changes are quite minor.

```
collection-home.jsp
admin/authorize-collection-edit.jsp
admin/authorize-community-edit.jsp
admin/authorize-item-edit.jsp
admin/eperson-edit.jsp
```

5. Restart Tomcat.

## 4.9 Updating From 1.0.1 to 1.1

To upgrade from DSpace 1.0.1 to 1.1, follow the steps below. Your `dspace.cfg` does not need to be changed. In the notes below `[dspace]` refers to the install directory for your existing DSpace installation, and `[dspace-1.1-source]` to the source directory for DSpace 1.1. Whenever you see these path references, be sure to replace them with the actual path names on your local system.

1. Take down Tomcat (or whichever servlet container you're using).
2. We recommend that you upgrade to the latest version of PostgreSQL (7.3.2). Included are some notes to help you do this. Note you will also have to upgrade Ant to version 1.5 if you do this.

3. Make the necessary changes to the DSpace database. These include a couple of minor schema changes, and some new indices which should improve performance. Also, the names of a couple of database views have been changed since the old names were so long they were causing problems. First run `psql` to access your database (e.g. `psql -U dspace -W` and then enter the password), and enter these SQL commands:

```
ALTER TABLE bitstream ADD store_number INTEGER;
UPDATE bitstream SET store_number = 0;
ALTER TABLE item ADD last_modified TIMESTAMP;
CREATE INDEX last_modified_idx ON Item(last_modified);
CREATE INDEX eperson_email_idx ON EPerson(email);
CREATE INDEX item2bundle_item_idx ON Item2Bundle(item_id);
CREATE INDEX bundle2bitstream_bundle_idx ON Bundle2Bitstream(bundle_id);
CREATE INDEX dcvalue_item_idx ON DCValue(item_id);
CREATE INDEX collection2item_collection_idx ON Collection2Item(collection_id);
CREATE INDEX resourcepolicy_type_id_idx ON ResourcePolicy (resource_type_id,resource_id);
CREATE INDEX epersongroup2eperson_group_idx ON EPersonGroup2EPerson(eperson_group_id);
CREATE INDEX handle_handle_idx ON Handle(handle);
CREATE INDEX sort_author_idx ON ItemsByAuthor(sort_author);
CREATE INDEX sort_title_idx ON ItemsByTitle(sort_title);
CREATE INDEX date_issued_idx ON ItemsByDate(date_issued);
DROP VIEW CollectionItemsByDateAccessioned;
DROP VIEW CommunityItemsByDateAccessioned;
CREATE VIEW CommunityItemsByDateAccession \
  as SELECT Community2Item.community_id, ItemsByDateAccessioned.* \
  FROM ItemsByDateAccessioned, Community2Item \
  WHERE ItemsByDateAccessioned.item_id = Community2Item.item_id;
CREATE VIEW CollectionItemsByDateAccession \
  AS SELECT collection2item.collection_id, \
  itemsbydateaccessioned.items_by_date_accessioned_id, \
  itemsbydateaccessioned.item_id, \
  itemsbydateaccessioned.date_accessioned \
  FROM itemsbydateaccessioned, collection2item \
  WHERE (itemsbydateaccessioned.item_id = collection2item.item_id);
```

4. Fix your JSPs for Unicode. If you've modified the site 'skin' (`jsp/local/layout/header-default.jsp`) you'll need to add the Unicode header, i.e.:

```
<meta http-equiv="Content-Type" content="text/html; charset=UTF-8">
```

to the `<HEAD>` element. If you have any locally-edited JSPs, you need to add this page directive to the top of all of them:

```
<%@ page contentType="text/html; charset=UTF-8" %>
```

(If you haven't modified any JSPs, you don't have to do anything.)

5. Copy the required Java libraries that we couldn't include in the bundle to the source tree. For example:

```
cd [dspace]/lib
cp *.policy activation.jar servlet.jar mail.jar [dspace-1.1-source]/lib
```

6. Compile up the new DSpace code, replacing [dSPACE]/config/dSPACE.cfg with the path to your current, LIVE configuration. (The second line, touch 'find .', is a precaution, which ensures that the new code has a current timestamp and will overwrite the old code. Note that those are back quotes.)

```
cd [dSPACE-1.1-source]
touch `find .`
ant
ant -Dconfig=[dSPACE]/config/dSPACE.cfg update
```

7. Update the database tables using the upgrader tool, which sets up the new >last\_modified date in the item table:

Run

```
[dSPACE]/bin/dsrun org.dSPACE.administer.Upgrade101To11
```

8. Run the collection default authorisation policy tool:

```
[dSPACE]/bin/dsrun org.dSPACE.authorize.FixDefaultPolicies
```

9. Fix the OAICat properties file.  
Edit [dSPACE]/config/templates/oaicat.properties.  
Change the line that says

```
Identify.deletedRecord=yes
```

To:

```
Identify.deletedRecord=persistent
```

This is needed to fix the OAI-PMH 'Identity' verb response. Then run [dSPACE]/bin/install-configs.

10. Re-run the indexing to index abstracts and fill out the renamed database views:

```
[dSPACE]/bin/index-all
```

11. Restart Tomcat. Tomcat should be run with the following environment variable set, to ensure that Unicode is handled properly. Also, the default JVM memory heap sizes are rather small. Adjust -Xmx512M (512Mb maximum heap size) and -Xms64M (64Mb Java thread stack size) to suit your hardware.

```
JAVA_OPTS="-Xmx512M -Xms64M -Dfile.encoding=UTF-8"
```

## Chapter 5

# Configuration and Customization

There are a number of ways in which DSpace can be configured and/or customized:

- Altering the configuration files in [dspace]/config
- Creating modified versions of the JSPs; these can be placed separately from and override the default installed JSPs, so that future updates of the code won't overwrite your changes
- Implementing a custom 'plug-in' class – for example, an 'authenticator' class, so that user authentication in the Web UI can be adapted and integrated with any existing mechanisms your organization might use, or a 'media filter' to generate thumbnails or extract full text from a new file format
- Editing the source code

Of these methods, only the last is likely to cause any headaches; if you update the DSpace source code directly, particularly core class files in `org.dspace.*` or `org.dspace.storage.*`, it may make applying future updates difficult. Before doing this, it is strongly recommended that you [e-mail the DSpace developer community](#) to find out the best way to proceed, and the best way to implement your change in a way that can be [contributed back](#) to DSpace for everyone's benefit (and saving you from having sole responsibility to maintain the code).

### 5.1 The `dspace.cfg` Configuration Properties File

The primary way of configuring DSpace is to edit the `dspace.cfg`. You'll definitely have to do this before you can operate DSpace properly. `dspace.cfg` contains basic information about a DSpace installation, including system path information, network host information, and other things like site name. The default `dspace.cfg` is a good source of information, and contains comments for all properties. It's a basic Java properties file, where lines are either comments, starting with a '#', blank lines, or property/value pairs of the form:

```
property.name = property value
```

The property value may contain references to other configuration properties, in the form `$property.name`. This follows the ant convention of allowing references in property files. A property may not refer to itself. Examples:

```
property.name = word1 ${other.property.name} more words
property2.name = ${dspace.dir}/rest/of/path
```

Due to time constraints, this document does not contain an exhaustive list of properties; they are all listed in the supplied `dspace.cfg`. Below are some particularly relevant properties with notes for their use:

**dspace.dir** Root directory of DSpace installation. Omit the trailing '/'. Note that if you change this, there are several other parameters you will probably want to change to match, e.g. `assetstore.dir`.

**Example:** `/dspace`

**dspace.url** Main URL at which DSpace Web UI webapp is deployed. Include any port number, but do not include a trailing '/'

**Example:**

`http://dspace.myu.edu`

`http://dspace.myu.edu:8080`

**dspace.hostname** Fully qualified hostname; do not include port number

**Example:** `dspace.myu.edu`

**dspace.name** Short and sweet site name, used throughout Web UI, e-mails and elsewhere (such as OAI protocol) **Example:** `DSpace at My University`

**config.template.foo** When `install-configs` is run, the file `[dspace]/config/templates/foo` file will be filled out with values from `dspace.cfg` and copied to the value of this property, in this example `/opt/othertool/cfg/foo`.

See [here](#) for more information. **Example:** `/opt/othertool/cfg/foo`

**plugin.sequence.org.dspace.eperson.AuthenticationMethod** Comma-separated list of classes implementing the `org.dspace.eperson.AuthenticationMethod` interface, which make up the **authentication stack**. Authentication methods are called on in the order listed.

**Example:**

`org.dspace.eperson.X509Authentication, org.dspace.eperson.PasswordAuthentication`

**authentication.x509.keystore.path** Path to Java keystore containing Client CA's certificate for client X.509 certificates (Optional; only needed if X.509 user authentication is used.) **Example:** `/tomcat/conf/keystore`

**authentication.x509.keystore.password** Password to Java keystore configured above in `authentication.x509.keystore.path`

**Example:** `changeit`

**handle.prefix** The Handle prefix for your site, see the [Handle section](#) **Example:** `1721.1234`

**assetstore.dir** The location in the file system for asset (bitstream) store number zero. This should be a directory for the sole use of DSpace. **Example:** `/bigdisk/store`

**assetstore.dir.n** The location in the file system of asset (bitstream) store number n. When adding additional stores, start with 1 (`assetstore.dir.1`) and count upwards. Always leave asset store zero (`assetstore.dir`). For more details, see the [Bitstream Storage](#) section. **Example:** `/anotherdisk/store1`

**assetstore.incoming** The asset store number to use for storing new bitstreams. For example, if `assetstore.dir.1` is `/anotherdisk/store1`, and `assetstore.incoming` is 1, new bitstreams will be stored under `/anotherdisk/store1`. A value of 0 (zero) corresponds to `assetstore.dir`. For more details, see the [hyperref\[sec:BitstreamStore\]Bitstream Storage](#) section. **Example:** `1`

**srb.xxx, srb.xxx.n** The sets of SRB access parameters (see `dspace.cfg`) if one or more SRB accounts are used. The `srb.xxx` set would correspond to asset (bitstream) store number zero. The `srb.xxx.n` set would correspond to asset (bitstream) store number n. For more details, see the [hyperref\[sec:BitstreamStore\]Bitstream Storage](#) section. **Example:** `/zone/home/user.domain`

**webui.submit.enable-cc** Enable the Creative Commons license step in the submission process. Submitters are given an opportunity to select a Creative Commons license to accompany the Item. Creative Commons licenses govern the use of the content. For more details, see the [Creative Commons](#) website. **Example:**

Property values can include other, previously defined values, by enclosing the property name in `$...`. For example, if your `dspace.cfg` contains: -

```
dspace.dir = /dspace
dspace.history = ${dspace.dir}/history
```

Then the value of the `dspace.history` property is expanded to be `/dspace/history`. This method is especially useful for handling commonly used file paths.

Whenever you edit `dspace.cfg`, you should then run `[dspace]/bin/install-configs` so that any changes you may have made are reflected in the configuration files of other applications, for example Apache. You may then need to restart those applications, depending on what you changed.

## 5.2 Wording of E-mail Messages

Sometimes DSpace automatically sends e-mail messages to users, for example to inform them of a new workflow task, or as a subscription e-mail alert. The wording of emails can be changed by editing the relevant file in `[dspace]/config/emails`. Each file is commented. Be careful to keep the right number 'placeholders' (e.g.2). **Note:** You should replace the contact-information "`dspace-help@myu.edu or call us at xxx-555-xxxx`" with your own contact details in:

- `config/emails/change_password`
- `config/emails/register`

## 5.3 The Metadata and Bitstream Format Registries

The `[dspace]/config/registries` directory contains two XML files. These are used to load the initial contents of the [Dublin Core Metadata registry](#) and [Bitstream Format registry](#). After the initial loading (performed by `ant fresh_install` above), the registries reside in the database; the XML files are not updated.

In order to change the registries, you may adjust the XML files before the first installation of DSpace. On an already running instance it is recommended to change the registries via DSpace admin UI. The changes made via admin UI are not reflected in the XML files.

### 5.3.1 Metadata Format Registries

The default metadata schema is Dublin Core, so DSpace is distributed with a default Dublin Core Metadata Registry. Currently, the system requires that every item have a Dublin Core record. Via the DSpace admin UI you may define new Metadata Schemas, edit existing schemas and move elements between schemas. There is a set of Dublin Core Elements, which is used by the system and should not be removed or moved to another schema, see Appendix: [Default Dublin Core Metadata registry](#). **Note:** altering a Metadata Registry has no effect on corresponding parts, e.g. item submission interface, item display, item import and vice versa. Every metadata element used in submission interface or item import must be registered before using it.

**Note** also that deleting a metadata element will delete all its corresponding values.

### 5.3.2 Bitstream Format Registry

The bitstream formats recognized by the system and levels of support are similarly stored in the bitstream format registry. This can also be edited at install-time via `[dspace]/config/registries/bitstream-formats.xml` or by the administration Web UI. The contents of the bitstream format registry are entirely up to you, though the system requires that the following two formats are present:

- Unknown
- License

**Note:** Deleting a format will cause any existing bitstreams of this format to be reverted to the unknown bitstream format.

## 5.4 The Default Submission License

For each submitted item, a license must be granted. The license will be stored along with the item in the bundle LICENSE in order to keep the information under which terms an items has been published.

You may define a license for each collection separately, when creating/editing a collection. If no collection specific license is defined, the default license is used.

The default license can be found in `[dspace]/config/default.license` and can be edited via the `dspace-admin` interface.

DSpace comes with a demo license, which you must adopt to your institutional needs and the legal regulations of your country.

If in doubt, contact the law department of your institution.

### 5.4.1 Possible Points in a License

**Note:** this is no legal advice, just some starting thoughts for creating you own license.

- Non-exclusive or exclusive right to
  - capture and store
  - distribute
    - \* worldwide
    - \* restricted (e.g. institutional wide)
  - translate
  - transform to other formats or mediums without changing the content
- Make sure no rights (copyright or any other) are violated by this publication
- In case the type of submission (e.g. thesis) needs approval, make sure it is the final and approved version.
- Distinguish between the document itself and the metadata
- Point out that the license granted and the information who granted it will be stored.

## 5.5 Activating Additional OAI-PMH Crosswalks

### 5.5.1 METS

The old DSpace 1.3 METS Crosswalk can be activated as follows:

1. Uncomment the METS Crosswalk entry from the `config/templates/oaicat.properties` file
2. Run the `bin/install-configs` script
3. Restart Tomcat
4. Verify the Crosswalk is activated by accessing a URL such as  
`http://myspace/dspace-oai/request?verb=ListRecords&metadataPrefix=mets`

### 5.5.2 Crosswalk Plugins

OAI-PMH crosswalks based on Crosswalk Plugins, including the new (DSpace 1.4) METS crosswalk, are activated as follows:

1. Ensure the crosswalk plugin has a lower-case name (possibly in addition to its upper-case name) in the plugin configuration.
2. Add a line to the file `config/templates/oaicat.properties` of the form:  
`Crosswalks.plugin_name=org.dspace.app.oai.PluginCrosswalk`  
substituting the plugin's name, e.g. `"mets"` or `"qdc"` for `plugin_name`.
3. Run the `bin/install-configs` script
4. Restart your servlet container, e.g. Tomcat, for the change to take effect.

The DSpace 1.4 source includes the following crosswalk plugins available for use with OAI-PMH:

**mets** The manifest document from a DSpace METS SIP.

**mods** MODS metadata, produced by the table-driven **MODS dissemination** crosswalk.

**qdc** Qualified Dublin Core, produced by the **configurable QDC crosswalk**. Note that this QDC does not include all of the DSpace "dublin core" metadata fields, since the XML standard for QDC is defined for a different set of elements and qualifiers.

### 5.5.3 DIDL

By activating the DIDL provider, DSpace items are represented as MPEG-21 DIDL objects. These DIDL objects are XML documents that wrap both the Dublin Core metadata that describes the DSpace item and its actual bitstreams. A bitstream is provided inline in the DIDL object in a base64 encoded manner, and/or by means of a pointer to the bitstream. The data provider exposes DIDL objects via the `metadataPrefix didl`.

The crosswalk does not deal with special characters and purposely skips dissemination of the `license.txt` file awaiting a better understanding on how to map DSpace rights information to MPEG21-DIDL.

The DIDL Crosswalk can be activated as follows:

- Uncomment the `oai.didl.maxresponse` item in `dspace.cfg`
- Uncomment the DIDL Crosswalk entry from the `config/templates/oaicat.properties` file



- Run the bin/install-configs script
- Restart Tomcat
- Verify the Crosswalk is activated by accessing a URL such as `http://myspace/dspace-oai/request?verb=ListRecords&metadataPrefix=didl`

## 5.6 Configuration Files for Other Applications

To ease the hassle of keeping configuration files for other applications involved in running a DSpace site, for example Apache, in sync, the DSpace system can automatically update them for you when the main DSpace configuration is changed. This feature of the DSpace system is entirely optional, but we found it useful.

The way this is done is by placing the configuration files for those applications in [dspace]/config/templates, and inserting special values in the configuration file that will be filled out with appropriate DSpace configuration properties. Then, tell DSpace where to put filled-out, 'live' version of the configuration by adding an appropriate property to dspace.cfg, and run [dspace]/bin/install-configs.

Take the apache13.conf file as an example. This contains plenty of Apache-specific stuff, but where it uses a value that should be kept in sync across DSpace and associated applications, a 'placeholder' value is written. For example, the host name:

```
ServerName @@dspace.hostname@@
```

The text @@dspace.hostname@@ will be filled out with the value of the dspace.hostname property in dspace.cfg. Then we decide where we want the 'live' version, that is, the version actually read in by Apache when it starts up, will go.

Let's say we want the live version to be located at /opt/apache/conf/dspace-httpd.conf. To do this, we add the following property to dspace.cfg so DSpace knows where to put it:

```
config.template.apache13.conf = /opt/apache/conf/dspace-httpd.conf
```

Now, we run [dspace]/bin/install-configs. This reads in [dspace]/config/templates/apache13.conf, and places a copy at /opt/apache/conf/dspace-httpd.conf with the placeholders filled out.

So, in /opt/apache/conf/dspace-httpd.conf, there will be a line like:

```
ServerName dspace.myu.edu
```

The advantage of this approach is that if a property like the hostname changes, you can just change it in dspace.cfg and run install-configs, and all of your tools' configuration files will be updated.

However, take care to make all your edits to the versions in [dspace]/config/templates! It's a wise idea to put a big reminder at the top of each file, since someone might unwittingly edit a 'live' configuration file which would later be overwritten.

## 5.7 Customizing the Web User Interface

The Web UI is implemented using Java Servlets which handle the business logic, and JavaServer Pages (JSPs) which produce the HTML pages sent to an end-user. Since the JSPs are much closer to HTML than Java code, altering the look and feel of DSpace is relatively easy.

To make it even easier, DSpace allows you to 'override' the JSPs included in the source distribution with modified versions, that are stored in a separate place, so when it comes to updating your site with a new DSpace release, your modified versions will not be overwritten. It should be possible to dramatically change the look of DSpace to suit your organization by just changing the CSS style file and the site 'skin' or 'layout' JSPs in jsp/layout;

if possible, it is recommended you limit local customizations to these files to make future upgrades easier. You can also easily edit the text that appears on each JSP page by editing the dictionary file. However, note that unless you change the entry in all of the different language message files, users of other languages will still see the default text for their language. See [internationalization](#). Note that the data (attributes) passed from an underlying Servlet to the JSP may change between versions, so you may have to modify your customized JSP to deal with the new data. Thus, if possible, it is recommended you limit your changes to the 'layout' JSPs and the stylesheet. The JSPs are stored in [dspace-source]/jsp. Place your edited version of a JSP in the [dspace-source]/jsp/local directory, with the same path as the original. If they exist, these will be used in preference to the distributed versions in [dspace-source]/jsp. For example:

**DSPACE default** [dspace-source]/jsp/community-list.jsp

**Locally-modified version** [dspace-source]/jsp/local/community-list.jsp

Heavy use is made of a style sheet, in [dspace-source]/jsp/styles.css.jsp. If you make edits, call the local version [dspace-source]/jsp/local/styles.css.jsp, and it will be used automatically in preference to the default, as described above.

Fonts and colors can be easily changed using the stylesheet. The stylesheet is a JSP so that the user's browser version can be detected and the stylesheet tweaked accordingly.

The 'layout' of each page, that is, the top and bottom banners and the navigation bar, are determined by the JSPs [dspace-source]/jsp/layout/header-\*.jsp and [dspace-source]/jsp/layout/footer-\*.jsp. You can provide modified versions of these (in [dspace-source]/jsp/local/layout, or define more styles and apply them to pages by using the "style" attribute of the dspace:layout tag.

After you've customized your JSPs, **you must rebuild the DSpace Web application**. If you haven't already built and installed it, follow the install directions. Otherwise, follow the steps below:

1. Rebuild the dspace.war file by running the following command from your [dspace-source] directory:

```
ant -Dconfig=[dspace]/config/dspace.cfg build_wars
```

2. Shut down Tomcat, and delete the existing [tomcat]/webapps/dspace directory.
3. Copy the new .war file to the Tomcat webapps directory:

```
cp [dspace-source]/build/dspace.war [tomcat]/webapps
```

When you restart the web server you should see your customized JSPs.

## 5.8 Customizing the Simple Item Display Metadata

In dspace.cfg create a webui.itemdisplay.default item specifying a comma-separated list of metadata fields to display. For example,

```
webui.itemdisplay.default = dc.title, \
                           dc.date.issued(date),
                           dc.identifier.uri(link), \
                           dc.description.*
```

will set the default simple item display to show the Dublin Core title, issue date (in date format), handle URI (rendered as a link) and all DC description metadata. If no `webui.itemdisplay.default` is present, the default defers to a preset list hard-coded in DSpace. If an item has no value for a particular field, it won't be displayed. Add entries to the `Messages.properties` file as required. The name of the field for display will be drawn from the current UI dictionary, using the key `metadata.<metadata field>`. For example `metadata.dc.title = Title`. The metadata in `dspace.cfg` can be specified in the form `<schema prefix>.<element>[.<qualifier>|.]*[[<(date)|(link)>]], ...`. For example,

- `dc.title` refers to unqualified Dublin Core (DC) 'title' element
- `dc.title.alternative` refers to the qualified DC element 'title.alternative'
- `dc.title.*` refers to all DC 'title' elements (i.e. any or no qualifier)
- `dc.identifier.uri(link)` refers to DC 'identifier.uri' and render as a link
- `dc.date.issued(date)` refers to DC 'date.issued' and render as a date

### 5.8.1 Customizing the Simple Item Display Metadata for Individual Collections

Create `dspace.cfg` entries for each of the "styles" of item display in the same way the default layout is configured. For example,

```
webui.itemdisplay.thesis-style = dc.contributor.*, dc.identifier.uri, dc.description.abstract
```

Then associate this simple item display with collections using an additional `dspace.cfg` item. For example,

```
webui.itemdisplay.thesis-style.collections = 123456789/1, 123456789/56
```

Don't forget to add any required message keys to the `Messages.properties` file. In the above example, the `Messages.properties` file would need to contain something like:

```
metadata.dc.contributor.* = Authors
metadata.dc.identifier.uri = Citation
metadata.dc.description.abstract = Abstract
```

## 5.9 Custom Authentication Code

Since many institutions and organizations have existing authentication systems, DSpace has been designed to allow these to be easily integrated into an existing authentication infrastructure. It keeps a series, or "stack", of authentication methods, so each one can be tried in turn. This makes it easy to add new authentication methods or rearrange the order without changing any existing code. You can also share authentication code with other sites.

The configuration property `plugin.sequence.org.dspace.eperson.AuthenticationMethod` defines the authentication stack. It is a comma-separated list of class names. Each of these classes implements a different authentication method, or way of determining the identity of the user. They are invoked in the order specified until one succeeds.

An authentication method is a class that implements the interface `org.dspace.eperson.AuthenticationMethod`. It authenticates a user by evaluating the credentials (e.g. username and password) he or she presents and checking that they are valid.

The basic authentication procedure in the DSpace Web UI is this:

1. A request is received from an end-user's browser that, if fulfilled, would lead to an action requiring authorization taking place.
2. If the end-user is already authenticated:
  - (a) If the end-user is allowed to perform the action, the action proceeds
  - (b) If the end-user is NOT allowed to perform the action, an authorization error is displayed.
  - (c) If the end-user is NOT authenticated, i.e. is accessing DSpace anonymously:
3. The parameters etc. of the request are stored
4. The Web UI's `startAuthentication` method is invoked.
5. First it tries all the authentication methods which do implicit authentication (i.e. they work with just the information already in the Web request, such as an X.509 client certificate). If one of these succeeds, it proceeds from Step 2 above.
6. If none of the implicit methods succeed, the UI responds by putting up a "login" page to collect credentials for one of the explicit authentication methods in the stack. The servlet processing that page then gives the proffered credentials to each authentication method in turn until one succeeds, at which point it retries the original operation from Step 2 above.

Please see the source files `AuthenticationManager.java` and `AuthenticationMethod.java` for more details about this mechanism.

### 5.9.1 Authentication by Password

The default method `org.dspace.eperson.PasswordAuthentication` has the following properties:

- Use of inbuilt e-mail address/password-based log-in. This is achieved by forwarding a request that is attempting an action requiring authorization to the password log-in servlet, `/password-login`. The password log-in servlet (`org.dspace.app.webui.servlet.PasswordServlet`) contains code that will resume the original request if authentication is successful, as per step 3. described above.
- Users can register themselves (i.e. add themselves as e-people without needing approval from the administrators), and can set their own passwords when they do this
- Users are not members of any special (dynamic) e-person groups

### 5.9.2 LDAP Authentication

As of version 1.3, the authentication method `org.dspace.eperson.LDAPAuthentication` is also supplied to support **LDAP authentication**.

### 5.9.3 X.509 Certificate Authentication

The X.509 authentication method uses an X.509 certificate sent by the client to establish his/her identity. It requires the client to have a personal Web certificate installed on their browser (or other client software) which is issued by a Certifying Authority (CA) recognized by the web server.

1. See the **HTTPS installation** instructions to configure your Web server. If you are using HTTPS with Tomcat, note that the `<Connector>` tag must include the attribute `clientAuth="true"` so the server requests a personal Web certificate from the client.

2. Add the `org.dspace.eperson.X509Authentication` plugin first to the list of stackable authentication methods in the value of the configuration key `plugin.sequence.org.dspace.eperson.AuthenticationMethod` E.g.:

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

3. You must also configure DSpace with the same CA certificates as the web server, so it can accept and interpret the clients' certificates. It can share the same keystore file as the web server, or a separate one, or a CA certificate in a file by itself. Configure it by one of these methods, either the Java keystore

```
authentication.x509.keystore.path = path to Java keystore file
authentication.x509.keystore.password = password to access the keystore
```

...or the separate CA certificate file (in PEM or DER format):

```
authentication.x509.ca.cert = path to certificate file for CA whose client certs to accept.
```

4. Choose whether to enable auto-registration: If you want users who authenticate successfully to be automatically registered as new E-Persons if they are not already, set the `authentication.x509.autoregister` configuration property to true. This lets you automatically accept all users with valid personal certificates. The default is false.

### 5.9.4 Example of a Custom Authentication Method

Also included in the source is an implementation of an authentication method used at MIT, `edu.mit.dspace.MITSpecialGroup`. This does not actually authenticate a user, it only adds the current user to a special (dynamic) group called 'MIT Users' (which must be present in the system!). This allows us to create authorization policies for MIT users without having to manually maintain membership of the MIT users group.

By keeping this code in a separate method, we can customize the authentication process for MIT by simply adding it to the stack in the DSpace configuration. None of the code has to be touched.

You can create your own custom authentication method and add it to the stack. Use the most similar existing method as a model, e.g. `org.dspace.eperson.PasswordAuthentication` for an "explicit" method (with credentials entered interactively) or `org.dspace.eperson.X509Authentication` for an implicit method.

## 5.10 Configuring LDAP Authentication

You can enable LDAP authentication by adding its method to the stack in the DSpace configuration, e.g.

```
plugin.sequence.org.dspace.eperson.AuthenticationMethod = org.dspace.eperson.LDAPAuthentication
```

If LDAP is enabled in the `dspace.cfg` file, then new users will be able to register by entering their username and password without being sent the registration token. If users do not have a username and password, then they can still register and login with just their email address the same way they do now.

If you want to give any special privileges to LDAP users, create a stackable authentication method to automatically put people who have a netid into a special group. You might also want to give certain email addresses special privileges. Refer to the [Custom Authentication Code](#) section above for more information about how to do this.

Here is an explanation of what each of the different configuration parameters are for:

**ldap.enable** This setting will enable or disable LDAP authentication in DSpace. With the setting off, users will be required to register and login with their email address. With this setting on, users will be able to login and register with their LDAP user ids and passwords.

**webui.ldap.autoregister** This will turn LDAP autoregistration on or off. With this on, a new EPerson object will be created for any user who successfully authenticates against the LDAP server when they first login. With this setting off, the user must first register to get an EPerson object by entering their ldap username and password and filling out the forms.

**ldap.provider\_url = ldap://ldap.myu.edu/o=myu.edu** This is the url to your institution's ldap server. You may or may not need the /o=myu.edu part at the end. Your server may also require the ldaps:// protocol.

**ldap.id\_field = uid** This is the unique identifier field in the LDAP directory where the username is stored.

**ldap.object\_context = ou=people,o=myu.edu** This is the object context used when authenticating the user. It is appended to the ldap.id\_field and username. For example uid=username,ou=people,o=myu.edu. You will need to modify this to match your ldap configuration.

**ldap.search\_context = ou=people** This is the search context used when looking up a user's ldap object to retrieve their data for autoregistering. With ldap.autoregister turned on, when a user authenticates without an EPerson object we search the ldap directory to get their name and email address so that we can create one for them. So after we have authenticated against uid=username,ou=people,o=byu.edu we now search in ou=people for filtering on [uid=username]. Often the ldap.search\_context is the same as the ldap.object\_context parameter. But again this depends on your ldap server configuration.

**ldap.email\_field = mail** This is the ldap object field where the user's email address is stored. "mail" is the default and the most common for ldap servers. If the mail field is not found the username will be used as the email address when creating the eperson object.

**ldap.surname\_field = sn** This is the ldap object field where the user's last name is stored. "sn" is the default and is the most common for ldap servers. If the field is not found the field will be left blank in the new eperson object.

**ldap.givenname\_field = givenName** This is the ldap object field where the user's given names are stored. I'm not sure how common the givenName field is in different ldap instances. If the field is not found the field will be left blank in the new eperson object.

**ldap.phone\_field = telephoneNumber** This is the field where the user's phone number is stored in the ldap directory. If the field is not found the field will be left blank in the new eperson object.

## 5.11 Configuring Lucene Search Indexes

Search Indexes can be configured via the dspace.cfg file. This allows institutions to choose which DSpace metadata fields are indexed by Lucene. Note that as of DSpace 1.4, the schema for the element is also required. For example, the following entries appear in a default DSpace installation:

```
search.index.1 = author:dc.contributor.*
search.index.2 = author:dc.creator.*
search.index.3 = title:dc.title.*
search.index.4 = keyword:dc.subject.*
search.index.5 = abstract:dc.description.abstract
```

```

search.index.6 = author:dc.description.statementofresponsibility
search.index.7 = series:dc.relation.ispartofseries
search.index.8 = abstract:dc.description.tableofcontents
search.index.9 = mime:dc.format.mimetype
search.index.10 = sponsor:dc.description.sponsorship
search.index.11 = id:dc.identifier.*

```

The form of each entry is `search.index.<id> = <search <schema>field>:<metadata field>` where:

- `<id>` is an incremental number to distinguish each search index entry
- `<search field>` is an identifier for the search field this index will correspond to
- `<metadata field>` is the DSpace metadata field to be indexed

So in the example above, `search.indexes1, 2` and `6` are configured as the author search field. The author index is created by Lucene indexing all contributor, creator and `description.statementofresponsibility` metadata fields. After changing the configuration, run `index-all` to recreate the indexes.

**NOTE:** While the indexes are created, this only affects the search results and has no effect on the search components of the user interface. To add new search capability (e.g. to add a new search category to the Advanced Search) requires local customisation to the user interface.

## 5.12 Configuring System Statistical Reports

Statistics for the system can be made available at <http://www.mydspaceinstance.edu/statistics>. To use the system statistics you will have to initialise them as per the installation documentation, but before you do so you need to perform the customisations discussed here in order to ensure that the reports are generated correctly.

### 5.12.1 Configuration File

Configuration for the statistics system are in `[dspace]/config/dstat.cfg` and the file should guide you to correctly filling in the details required. For the most part you will not need to change this file.

### 5.12.2 Customising Shell Scripts

To customise the supplied perl scripts to do monthly and general report generation it is necessary to modify the scripts themselves slightly. This is because these scripts were developed to speed up the process of using DStat at Edinburgh University Library and were not particularly intended for external use. They appear here for the convenience of others and in order to bridge the gap between the report generation and the inclusion of those reports into the DSpace UI, which is currently a clunky process.

In order to get these scripts to work for you, open each of the following in turn:

```

stat-general
stat-initial
stat-monthly
stat-report-general
stat-report-initial
stat-report-monthly

```

scripts ending with `-general` do the work for building reports spanning the entire history of the archive; scripts ending `-initial` are to initialise the reports by doing monthly reports from some start date up to the present;

scripts ending -monthly generate a single monthly report for the current month. These scripts are just designed to make life easier, and are not particularly clever or elegant.

In each file you will find a section:

```
# Details used
#####

... some perl ...

#####
```

the perl between the lines of hashes defines the variables which will be used to do all of the processing in the report. The following explains what the variables mean and what they should be set to for each of the scripts

**stat-initial:**

```
$out_prefix: prefix to place in front of each output file.
$out_suffix: suffix for output file. A date will be inserted between the prefix and suffix
$start_year: year to start back-analysing monthly logs from
$start_month: month to start back-analysing monthly logs from
$dsvrun: path to your dsvrun script, usually [dSPACE]/bin/dsvrun
$out_directory: directory into which to place analysis files, for example [dSPACE]/bin/log/
```

**stat-monthly:**

```
$out_prefix: prefix to place in front of each output file.
$out_suffix: suffix for output file. A date will be inserted between the prefix and suffix
$dsvrun: path to your dsvrun script, usually [dSPACE]/bin/dsvrun
$out_directory: directory into which to place analysis files, for example [dSPACE]/bin/log/
```

**stat-general:**

```
$out_prefix: prefix to place in front of each output file.
$out_suffix: suffix for output file. Today's date will be inserted between the prefix and suffix
$dsvrun: path to your dsvrun script, usually [dSPACE]/bin/dsvrun
$out_directory: directory into which to place analysis files, for example [dSPACE]/bin/log/
```

**stat-report-initial:**

```
$in_prefix: the prefix of the files generated by stat-initial
$in_suffix: the suffix of the files generated by stat-initial
$out_prefix: the report file prefix. Should be "report-" in order to work with DSpace UI
$out_suffix: the report file suffix. Should be ".html" in order to work with DSpace UI
$start_year: the start year used in stat-initial
$start_month: the start month used in stat-initial
$dsvrun: path to your dsvrun script, usually [dSPACE]/bin/dsvrun
$in_directory: directory where analysis files were placed in stat-initial
$out_directory: the live reports directory: [dSPACE]/reports/
```

**stat-report-monthly:**

```
$in_prefix: the prefix of the files generated by stat-monthly
$in_suffix: the suffix of the files generated by stat-monthly
$out_prefix: the report file prefix. Should be "report-" in order to work with DSpace UI
```



\$out\_suffix: the report file suffix. Should be ".html" in order to work with DSpace UI  
 \$dsrun: path to your dsrun script, usually [dSPACE]/bin/dsrun  
 \$in\_directory: directory where analysis files were placed in stat-monthly  
 \$out\_directory: the live reports directory: [dSPACE]/reports/

#### stat-report-general:

\$in\_prefix: the prefix of the files generated by stat-general  
 \$in\_suffix: the suffix of the files generated by stat-general  
 \$out\_prefix: the report file prefix. Should be "report-general-" in order to work with DSpace UI  
 \$out\_suffix: the report file suffix. Should be ".html" in order to work with DSpace UI  
 \$dsrun: path to your dsrun script, usually [dSPACE]/bin/dsrun  
 \$in\_directory: directory where analysis files were placed in stat-general  
 \$out\_directory: the live reports directory: [dSPACE]/reports/

If you want additional customisations, you will need to modify the lines which build the command to be executed and change the parameters passed to the java processes which actually carry out the analysis. For more information on these processes either build the javadocs or run:

```
[dSPACE]/bin/dsrun ac.ed.dSPACE.stats.LogAnalyser -help
[dSPACE]/bin/dsrun ac.ed.dSPACE.stats.ReportGenerator -help
```

## 5.13 MediaFilters

Media Filters are classes used to generate derivative or alternative versions of master bitstreams. For example, the PDF Media Filter will extract textual content from PDF bitstreams, the JPEG Media Filter can create thumbnails from image bitstreams.

Media Filters are configured as a Sequence Plugin, with each filter also having a separate config item indicating which formats it can process. The default configuration is shown below.

```
#### Media Filter plugins (through PluginManager) ####

plugin.sequence.org.dSPACE.app.mediafilter.MediaFilter = \
  org.dSPACE.app.mediafilter.PDFFilter, org.dSPACE.app.mediafilter.HTMLFilter, \
  org.dSPACE.app.mediafilter.WordFilter, org.dSPACE.app.mediafilter.JPEGFilter
# to enable branded preview: remove last line above, and uncomment 2 lines below
# org.dSPACE.app.mediafilter.WordFilter, org.dSPACE.app.mediafilter.JPEGFilter, \
# org.dSPACE.app.mediafilter.BrandedPreviewJPEGFilter

filter.org.dSPACE.app.mediafilter.PDFFilter.inputFormats = Adobe PDF
filter.org.dSPACE.app.mediafilter.HTMLFilter.inputFormats = HTML, Text
filter.org.dSPACE.app.mediafilter.WordFilter.inputFormats = Microsoft Word
filter.org.dSPACE.app.mediafilter.JPEGFilter.inputFormats = GIF, JPEG, image/png
filter.org.dSPACE.app.mediafilter.BrandedPreviewJPEGFilter.inputFormats = GIF, JPEG, image/png
```

To add a new Media Filter, add the new filter class to the plugin.sequence.org.dSPACE.app.mediafilter.MediaFilter config item and add a corresponding filter.<class path>.inputFormats config item. Note the input formats must match the short\_description field in the bitstreamformatregistry table.

## 5.14 Displaying Image Item Preview

For particular types of objects (currently various image formats), a preview feature can be activated. This provides a branded version of the DSpace bitstream for display on the Item Display page.

To activate this feature and display a preview image on the item page (all properties mentioned below are found in `dspace.cfg`):

1. Uncomment the lines defining the list of configured mediafilters at `plugin.sequence.org.dspace.app.mediafilter.MediaFilter` to include `org.dspace.app.mediafilter.BrandedPreviewJPEGFilter`.
2. Set the maximum pixel dimensions for the preview image by altering the `webui.preview.maxwidth` and `webui.preview.maxheight` (default is 600) config items.
3. Set the `webui.preview.brand` to the text you want to brand the image with. The brand will appear as white text on a black background strip across the base of the image. For example you might set the text to the owning organisation. The handle is also displayed as part of the branding.
4. Set the `webui.preview.brand.abbrev`. This is an abbreviated form of the `webui.preview.brand` text and will be shown where the brand text is longer than the image width (e.g. for narrow images).
5. Set the `webui.preview.brand.height` to the height in pixels of the appended branding (default 20).
6. Set the `webui.preview.brand.font` to the font for the brand text (default is SansSerif).
7. Set the `webui.preview.brand.fontpoint` to the font size for the brand text (default is 12).
8. Run `[dspace]/bin/filter-media` to generate the preview images. These will be stored in a `BRANDED_PREVIEW` bundle so won't interfere with existing bundles.
9. Set `webui.preview.enabled = true`.
10. Set `webui.preview.dc` if a metadata value is to appear just below the preview. For example, to highlight item rights you might set `webui.preview.dc = rights` or to show the description as part of the preview set `webui.preview.dc = description.abstract`. This is disabled by default.
11. Restart Tomcat.

## 5.15 Displaying Image Thumbnails

### 5.15.1 Browse and Search Results Page Thumbnails

Image thumbnails can be enabled on the Browse and Search Results pages by setting the appropriate configuration values. To enable the display of thumbnails the following items must be set in the `dspace.cfg` file:

```
webui.browse.thumbnail.show = true
```

If set to false or this configuration item is missing then thumbnails will not be shown. Additionally, appropriate media filters must be configured and the media filter configured to run periodically (for example, via a 'cron' job)

The size of the browse/search thumbnails can also be configured to a smaller size than that generated by the mediafilter. To do this set the following configuration items:

```
webui.browse.thumbnail.maxheight = <maxheight in pixels>
webui.browse.thumbnail.maxwidth = <maxwidth in pixels>
```

If these configuration items are not set, `thumbnail.maxheight` and `thumbnail.maxwidth` are used. Setting these values greater than or equal to the size of the thumbnail generated by the mediafilter (i.e. `thumbnail.maxheight` and `thumbnail.maxwidth`) will have no effect.

**Note:**

- where the primary bitstream is HTML, no thumbnail is shown;
- where the primary bitstream has a thumbnail, its thumbnail is shown;
- where the primary bitstream is not set, the first thumbnail found by DSpace will be shown;
- where the user does not have read access to the thumbnail bitstream, no thumbnail is shown;
- currently, for a thumbnail to display, a JPEG thumbnail under the current implementation rules must exist (i.e. primary bitstream name with ".jpg" suffix).

### 5.15.2 Configuring Thumbnail Link Behaviour

The target of a thumbnail in the Browse and Search Results Page can be configured by setting the following configuration item:

```
webui.browse.thumbnail.linkbehaviour = <target page type>
```

Currently the values `item` and `bitstream` are allowed. If this configuration item is not set, or set incorrectly, the default is `item`.

### 5.15.3 Item Display Page Thumbnails

Thumbnails may also be enabled or disabled on the Item Display page by setting the following configuration item in `dspace.cfg`:

```
webui.item.thumbnail.show = true
```

If set to `false` or this configuration item is missing then thumbnails will not be shown.

## 5.16 Displaying Community and Collection Item Counts

To show the item count against communities and collections set the `webui.strengths.show` configuration item in the `dspace.cfg` file as follows:

```
webui.strengths.show = true
```

If this config item is missing or is set to any value other than `true` the item counts will not be shown.

## 5.17 Lucene Analyzer

The Lucene analyzer used in searching and indexing can be configured by setting the `search.analyzer` configuration item in `dspace.cfg` to the class of the desired analyzer. If this item is not present/commented out, the default Lucene analyzer `org.dspace.search.DSAnalyzer` is used.

As well as those analyzers included in the Lucene distribution (see `lucene.jar`), a Chinese analyzer from the Lucene sandbox is included in `lucene-sandbox.jar`. This analyzer is yet to be included in the core Lucene distribution but can be configured by setting `search.analyzer = org.apache.lucene.analysis.cn.ChineseAnalyzer` in `dspace.cfg`.

## 5.18 On-line Help About File Formats

Because the file format support policy is determined by each individual institution, the on-line help on this subject is intentionally left blank. The help file will, however, retrieve a list of formats and the support levels associated with them in your database and display this information to the user. We highly recommend that you edit the "Format Support Policy" section of the file [dspace-source]/jsp/help/formats.jsp.

## 5.19 View Item Licence

Setting

```
webui.licence_bundle = true
```

in dspace.cfg will result in a hyperlink being rendered on the Item View page that points to the item's licence.

## 5.20 Configuring RSS Syndication

In addition to the email subscription service, recent submissions to communities and collections are also available via RSS feeds. The feed link or links appear on the community and collection home pages, in the sidebar underneath the Recent Submissions citations. The content of the feed includes the community/collection handle, name, short description, and icon; for each item it includes: the handle, title, author, date issued, description, and abstract if present. Feed behavior can be configured by setting the following properties in dspace.cfg:

```
webui.feed.enable = true
```

Set the value of this property to true to enable RSS feeds. If false, feeds will not be generated, and the feed links will not appear.

```
webui.feed.items = 4
```

The value of this property governs the number of DSpace items appearing in each feed, which are the most recent n submissions.

```
webui.feed.cache.size = 100
```

To improve performance, generated feeds may be cached in memory. If caching is desired, set the value of this property to a positive number, which represents the total number of feeds kept in the cache at one time, for all communities and collections. A value of 0 disables caching, and the feed is generated on demand for each request. If the cache grows beyond the specified value, the least-requested feeds are removed to maintain its size. To assist you in optimization, whenever cache maintenance is performed a profile of the cache is recorded in the DSpace system log, showing its size, and the total, average, minimum, and maximum number of 'hits' against cached feeds.

```
webui.cache.age = 48
```

To ensure that the feeds in the cache remain current, requested feeds are checked after they reach a certain age. This property specifies that age in hours. A value of 0 will force a check with each request, which guarantees currency, but with a performance penalty. If caching is disabled, this property is ignored.

```
webui.feed.formats = rss_1.0,rss_2.0
```

The RSS feature supports several different syndication formats. If you wish only a single format to be offered, set the value of this property to that format, selecting only from the list of canonical names provided. If you wish multiple formats, set the value to a comma-separated list of canonical names. Each format will appear as a distinct icon and link in the community and collection home page, as well as an 'autodiscovery' link in the page header.

```
webui.feed.localresolve = false
```

By default, the RSS feed will return global handle server-based URLs to items, collections and communities (e.g. <http://hdl.handle.net/123456789/1>). This means if you have not registered your DSpace installation with

the CNRI Handle Server (e.g. development or testing instance) the URLs returned by the feed will return an error if accessed. Setting `webui.feed.localresolve = true` will result in the RSS feed returning localised URLs (e.g. `http://myserver.myorg/handle/123456789/1`). If `webui.feed.localresolve` is set to false or not present the default global handle URL form is used.

```
webui.feed.item.title = dc.title
webui.feed.item.date = dc.date.issued
```

Specify which metadata field you want to be displayed as an item's title and date in the RSS feed. You can only specify a single metadata field for each of these properties.

```
webui.feed.item.description = dc.title, dc.contributor.author, \
dc.contributor.editor, dc.description.abstract, \
dc.description
```

Specify which metadata fields should be displayed in an item's description field in the RSS feed. You can specify as many fields as you wish here. The fields will be listed in the RSS feed in the order they appear, and will only be listed if they have values. Similar to the item display UI, the name of the field displayed in the feed will be drawn from the current UI dictionary, using the key: `"metadata.<field>"` (e.g. `"metadata.dc.title"`, `"metadata.dc.contributor.author"`, `"metadata.dc.description.*"`)

## 5.21 Configuring Item Recommendations

The 'suggest an item' service is a convenient way for DSpace users to notify others via email about content they discover in the repository. If enabled, a link appears on the simple item display page leading to a form where the user can enter an email address, and additional optional information. Upon completion of the form, the system transmits a pre-formatted message about the item to that address. Item recommendation can be configured by setting the following properties in `dspace.cfg`:

```
webui.suggest.enable = true
```

Set the value of this property to true to expose the link to the recommendation form. If false, the link will not display.

```
webui.suggest.sender = A DSpace user
```

The form includes an optional field for the name of the person making the recommendation. The value of this property will be used if the sender is not specified. Note that if the user has been authenticated to DSpace, the e-person record will be used to set the default value of this field, and will also be used to set the 'reply-to' email header to the e-person's address.

```
webui.suggest.recipient = colleague
```

In addition to the required destination email address, the form includes an optional field for the name of person to whom the recommendation is to be sent. The value of this property will be used if the recipient name is not specified.

The wording and layout of the email message is governed by the template file `[dspace]/config/emails/suggest`. You may edit this file to alter the language, layout, type, or amount of information contained in the email, but note that the values passed to the template are limited to those enumerated at the top of the template file and represented in the body of the message with 'placeholders' 0, 1, etc. A general description of DSpace is included as a footer to the email, which may be replaced with site-specific content.

## 5.22 Configuring Controlled Vocabularies

DSpace now supports controlled vocabularies to confine the set of keywords that users can use while describing items.

The need for a limited set of keywords is important since it eliminates the ambiguity of a free description system,

consequently simplifying the task of finding specific items of information.

The controlled vocabulary add-on allows the user to choose from a defined set of keywords organised in an tree (taxonomy) and then use these keywords to describe items while they are being submitted.

We have also developed a small search engine that displays the classification tree (or taxonomy) allowing the user to select the branches that best describe the information that he/she seeks.

The taxonomies are described in XML following this (very simple) structure:

```
<node id="acmccs98" label="ACMCCS98">
  <isComposedBy>
    <node id="A." label="General Literature">
      <isComposedBy>
        <node id="A.0" label="GENERAL"/>
        <node id="A.1" label="INTRODUCTORY AND SURVEY"/>
        ...
      </isComposedBy>
    </node>
    ...
  </isComposedBy>
</node>
```

You are free to use any application you want to create your controlled vocabularies. A simple text editor should be enough for small projects. Bigger projects will require more complex tools. You may use Protegé to create your taxonomies, save them as OWL and then use a XML Stylesheet (XSLT) to transform your documents to the appropriate format. Future enhancements to this add-on should make it compatible with standard schemas such as OWL or RDF.

In order to make DSpace compatible with WAI 2.0, the add-on is turned off by default (the add-on relies strongly on Javascript to function). It can be activated by setting the following property in `dspace.cfg`:

```
webui.controlledvocabulary.enable = true
```

New vocabularies should be placed in `[dspace]/config/controlled-vocabularies/` and must be according to the structure described. A validation XML Schema can be downloaded [here](#).

Vocabularies need to be associated with the correspondant DC metadata fields. Edit the file `[dspace]/config/input-forms.xml` and place a "vocabulary" tag under the "field" element that you want to control. Set value of the "vocabulary" element to the name of the file that contains the vocabulary, leaving out the extension (the add-on will only load files with extension "\*.xml"). For example:

```
<field>
  <dc-schema>dc</dc-schema>
  <dc-element>subject</dc-element>
  <dc-qualifier></dc-qualifier>
  <!-- An input-type of twobox MUST be marked as repeatable -->
  <repeatable>true</repeatable>
  <label>Subject Keywords</label>
  <input-type>twobox</input-type>
  <hint> Enter appropriate subject keywords or phrases below. </hint>
  <required></required>
  <vocabulary>nsi</vocabulary>
</field>
```

The following vocabularies are currently available by default:

**nsi** The Norwegian Science Index

**srsc** Swedish Research Subject Categories

## 5.23 Configuring the checksum checker

There are three aspects of the Checksum Checker's operation that can be configured:

1. the execution mode
2. the logging output
3. the policy for removing old checksum results from the database

### 5.23.1 Checker Execution Mode

Execution mode can be configured using command line options. Information on the options can be found at any time by running `[dspace]/bin/checker -help`. The different modes are described below; see the "Which to use" section that follows for details on the various pros and cons.

Unless a particular bitstream or handle is specified, the Checksum Checker will always check bitstreams in order of the least recently checked bitstream. (Note that this means that the most recently ingested bitstreams will be the last ones checked by the Checksum Checker.)

#### Limited Count Mode

To check a specific number of bitstreams, use the `-c` option followed by an integer number of bitstreams to check:

```
bin/checker -c 10
```

Limited count mode is particularly useful for checking that the checker is executing properly. The Checksum Checker's default execution mode is to check a single bitstream, as if the `-c 1` option had been given.

#### Limited Duration Mode

To run the Checker for a specific period of time, use the `-d` option with a time argument:

```
bin/checker -d 10m
```

```
bin/checker -d 2h
```

Valid options for specifying duration are `s` for seconds, `m` for minutes, `h` for hours, `d` for days, `w` for weeks, and `y` for years (OK, so we're optimists).

The checker will keep starting new bitstream checks for the specified duration, so actual execution duration will be slightly longer than the specified duration. Bear this in mind when scheduling checks.

#### Check Specific Bitstreams

To check one or more particular bitstreams by ID, use the `-b` option followed by one or more bitstream IDs:

```
bin/checker -b 1 2 3 4
```

This mode is useful when analyzing problems reported in the logs and when verifying that a resolution has been successful.

#### Check Specific Handles

Use the `-a` option followed by a handle:

```
bin/checker -a 123456/123
```

This will check all the bitstreams inside an item, collection or community.

### Continuous Looping

There are two looping modes:

```
bin/checker -l      # Loops once through the repository
bin/checker -L      # Loops continuously through the repository
```

The `-l` option can be used if your repository is relatively small and your backup strategy requires it to be completely validated at a particular point. The `-L` option might be useful if you have a large repository, and you don't mind (or can avoid) the IO load caused by the checker.

### Which to Use

The Checksum Checker was designed with the idea that most sys admins will run it from the cron. For small repositories we recommend using the `-l` option in the cron. For larger repositories that cannot be completely checked in a couple of hours, we recommend the `-d` option in the cron.

### 5.23.2 Checker Reporting

Checksum Checker uses `log4j` to report its results. By default it will report to a log called `[dspace]/log/checker.log`, and it will report only on bitstreams for which the newly calculated checksum does not match the stored checksum. To report on all bitstreams checked regardless of outcome, use the `-v` (verbose) command line option:

```
bin/checker -l -v    #Loop through the repository once and
report in detail about every bitstream checked.
```

To change the location of the log, or to modify the prefix used on each line of output, edit the `[dspace]/config/templates/log4j.properties` file and run `[dspace]/bin/install\configs`.

### 5.23.3 Checker Results Pruning

The Checksum Checker will store the result of every check in the `checksum_history` table. By default, successful checksum matches that are eight weeks old or older will be deleted when the `-p` command line option is used (unsuccessful ones will be retained indefinitely). The amount of time for which results are retained in the `checksum_history` table can be modified by one of two methods:

1. editing the retention policies in `[dspace]/config/dspace.cfg` OR
2. passing in a properties file containing retention policies when using the `-p` option.

Pruning is controlled by a number of properties, each of which describes a checksum result code, and the length of time for which results with that code should be retained. The format is `checker.retention.[RESULT CODE]=[duration]`. For example:

```
checker.retention.CHECKSUM_MATCH=8w
```

indicates that successful checksum matches will be retained for eight weeks. Supported units of time are

```
s Seconds
m Minutes
h Hours
d Days
w Weeks
y Years
```



(Note that these units are also used for describing durations for the `-d` limited duration mode.) There is a special property, `checker.retention.default`, that is used to assign a default retention period. To execute the pruning you must use the `-p` command line option (with or without a properties file). Checksum Checker will prune the history table before beginning new checks. We recommend that you use this option regularly, as the `checksum_history` table can grow very large without it.

## 5.24 Configuring Packager Plugins

Package ingester plugins are configured as named or self-named plugins for the interface `org.dspace.content.packager.PackageIngester`.

Package disseminator plugins are configured as named or self-named plugins for the interface `org.dspace.content.packager.PackageDisseminator`.

You can add names for the existing plugins, and add new plugins, by altering these configuration properties. See the Plugin Manager architecture for more information about plugins.

## 5.25 Configuring Crosswalk Plugins

Ingestion crosswalk plugins are configured as named or self-named plugins for the interface `org.dspace.content.crosswalk.IngestionCrosswalk`.

Dissemination crosswalk plugins are configured as named or self-named plugins for the interface `org.dspace.content.crosswalk.DisseminationCrosswalk`.

You can add names for existing crosswalks, add new plugin classes, and add new configurations for the configurable crosswalks as noted below.

### 5.25.1 Configurable MODS dissemination crosswalk

The MODS crosswalk is a self-named plugin. To configure an instance of the MODS crosswalk, add a property to the DSpace configuration starting with `"crosswalk.mods.properties."`; the final word of the property name becomes the plugin's name. For example, a property name `crosswalk.mods.properties.MODS` defines a crosswalk plugin named `"MODS"`. The value of this property is a path to a separate properties file containing the configuration for this crosswalk. The pathname is relative to the DSpace configuration directory, i.e. the config subdirectory of the DSpace install directory. So, a line like:

```
crosswalk.mods.properties.MODS = crosswalks/mods.properties
```

defines a crosswalk named MODS whose configuration comes from the file `[dspace]/config/crosswalks/mods.properties`.

The MODS crosswalk properties file is a list of properties describing how DSpace metadata elements are to be turned into elements of the MODS XML output document. The property name is a concatenation of the metadata schema, element name, and optionally the qualifier. For example, the `contributor.author` element in the native Dublin Core schema would be: `dc.contributor.author`. The value of the property is a line containing two segments separated by the vertical bar (`"|"`): The first part is an XML fragment which is copied into the output document. The second is an XPath expression describing where in that fragment to put the value of the metadata element. For example, in this property:

```
dc.contributor.author = <mods:name>\
<mods:role><mods:roleTerm type="text">author</mods:roleTerm>\
</mods:role><mods:namePart>%s</mods:namePart>\
</mods:name> | mods:namePart/text()
```

Some of the examples include the string "%s" in the prototype XML where the text value is to be inserted, but don't pay any attention to it, it is an artifact that the crosswalk ignores. For example, given an author named Jack Florey, the crosswalk will insert

```
<mods:name>
  <mods:role>
    <mods:roleTerm type="text">author</mods:roleTerm>
  </mods:role>
  <mods:namePart>Jack Florey</mods:namePart>
</mods:name>
```

into the output document. Read the example configuration file for more details.

### 5.25.2 Configurable Qualified Dublin Core (QDC) dissemination crosswalk

The QDC crosswalk is a self-named plugin. To configure an instance of the QDC crosswalk, add a property to the DSpace configuration starting with "crosswalk.qdc.properties."; the final word of the property name becomes the plugin's name. For example, a property name `crosswalk.qdc.properties.QDC` defines a crosswalk plugin named "QDC".

The value of this property is a path to a separate properties file containing the configuration for this crosswalk. The pathname is relative to the DSpace configuration directory, i.e. the `config` subdirectory of the DSpace install directory. So, a line like:

```
crosswalk.qdc.properties.QDC = crosswalks/qdc.properties
```

defines a crosswalk named QDC whose configuration comes from the file `[dspace]/config/crosswalks/qdc.properties`. You'll also need to configure the namespaces and schema location strings for the XML output generated by this crosswalk. The namespaces property names are of the format:

```
crosswalk.qdc.namespace.prefix = uri
```

where `prefix` is the namespace prefix and `uri` is the namespace URI.

For example, this shows how a crosswalk named "QDC" would be configured:

```
crosswalk.qdc.properties.QDC = crosswalks/QDC.properties
crosswalk.qdc.namespace.QDC.dc = http://purl.org/dc/elements/1.1/
crosswalk.qdc.namespace.QDC.dcterms = http://purl.org/dc/terms/
crosswalk.qdc.schemaLocation.QDC = \
  http://purl.org/dc/terms/ http://dublincore.org/schemas/xmls/qdc/2003/04/02/qualifieddc.xsd
```

The QDC crosswalk properties file is a list of properties describing how DSpace metadata elements are to be turned into elements of the Qualified DC XML output document. The property name is a concatenation of the metadata schema, element name, and optionally the qualifier. For example, the `contributor.author` element in the native Dublin Core schema would be: `dc.contributor.author`. The value of the property is an XML fragment, the element whose value will be set to the value of the metadata field in the property key.

For example, in this property:

```
dc.coverage.temporal = <dcterms:temporal />
```

the generated XML in the output document would look like, e.g.:

```
<dcterms:temporal>Fall, 2005</dcterms:temporal>
```

### 5.25.3 XSLT-based crosswalks

The XSLT crosswalks use XSL stylesheet transformation (XSLT) to transform an XML-based external metadata format to or from DSpace's internal metadata. XSLT crosswalks are much more powerful and flexible than the configurable MODS and QDC crosswalks, but they demand some esoteric knowledge (XSL stylesheets). Given that, you can create all the crosswalks you need just by adding stylesheets and configuration lines, without touching any of the Java code.

A submission crosswalk is described by a configuration key starting with "crosswalk.submission.", like

```
crosswalk.submission.PluginName.stylesheet = path
```

The PluginName is, of course, the plugin's name. The path value is the path to the file containing the crosswalk stylesheet (relative to `dspace.dir/config`).

Here is an example that configures a crosswalk named "LOM" using a stylesheet in `[dspace]/config/crosswalks/d-lom.xsl`:

```
crosswalk.submission.stylesheet.LOM = crosswalks/d-lom.xsl
```

A dissemination crosswalk is described by a configuration key starting with "crosswalk.dissemination.", like

```
crosswalk.dissemination.PluginName.stylesheet = path
```

The PluginName is, of course, the plugin's name. The path value is the path to the file containing the crosswalk stylesheet (relative to `dspace.dir/config`).

You can make two different plugin names point to the same crosswalk, by adding two configuration entries with the same path, e.g.

```
crosswalk.submission.MyFormat.stylesheet = crosswalks/myformat.xslt
crosswalk.submission.almost_DC.stylesheet = crosswalks/myformat.xslt
```

The dissemination crosswalk must also be configured with an XML Namespace (including prefix and URI) and an XML Schema for its output format. This is configured on additional properties in the DSpace Configuration, i.e.:

```
crosswalk.dissemination.PluginName.namespace.Prefix = namespace-URI
crosswalk.dissemination.PluginName.schemaLocation = schemaLocation value
```

For example:

```
crosswalk.dissemination.qdc.namespace.dc = http://purl.org/dc/elements/1.1/
crosswalk.dissemination.qdc.namespace.dcterms = http://purl.org/dc/terms/
crosswalk.dissemination.qdc.schemaLocation = \
http://purl.org/dc/elements/1.1/ http://dublincore.org/schemas/xmls/qdc/2003/04/02/qualifieddc.xsd
```

### 5.25.4 DSpace Intermediate Metadata (DIM) format

XSLT crosswalk plugins translate between the external metadata format and an XML format called DSpace Intermediate Metadata, which exists only for the purpose of XSLT crosswalks. It is never to be exported from DSpace, since it is not an acknowledged metadata format, it is simply an expression of the way DSpace stores its metadata fields internally. All the elements in a DIM document are in the namespace `http://www.dspace.org/xmlns/dspace/dim`.

The root element is named `dim`. It has zero or more children, all field elements. It may have an attribute `dspaceType`, which identifies the type of object ("ITEM", "COLLECTION", or "COMMUNITY") this metadata describes. This attribute is only guaranteed to be set for dissemination crosswalks.

Each field element may have the following attributes:

- mdschema (Required) The metadata schema, e.g. 'dc'.
- element (Required) Element name, such as 'contributor'.
- qualifier Qualifier name, such as 'author'.
- lang Language code describing language of this entry.

The value of field is the value of that metadata field. Fields with the same qualifiers may be repeated. Here is an example of the DIM format:

```
<dim:dim xmlns:dim="http://www.dspace.org/xmlns/dspace/dim" dspaceType="ITEM">
  <dim:field mdschema="dc" element="title" lang="en_US">
    The Endochronic Properties of Resublimated Thiotimonline
  </dim:field>
  <dim:field mdschema="dc" element="contributor" qualifier="author">
    Isaac Asimov
  </dim:field>
  <dim:field mdschema="dc" element="language" qualifier="iso">
    eng
  </dim:field>
  <dim:field mdschema="dc" element="subject" qualifier="other" lang="en_US">
    time-travel scifi hoax
  </dim:field>
  <dim:field element="publisher">
    Boston University Department of Biochemistry
  </dim:field>
</dim:dim>
```

### 5.25.5 Testing XSLT Crosswalks

The XSLT crosswalks will automatically reload an XSL stylesheet that has been modified, so you can edit and test stylesheets without restarting DSpace.

You can test a dissemination crosswalk by hooking it up to an OAI-PMH crosswalk and using an OAI request to get the metadata for a known item.

Testing the submission crosswalk is more difficult, so we have supplied a command-line utility to help. It calls the crosswalk plugin to translate an XML document you submit, and displays the resulting intermediate XML (DIM). Invoke it with:

```
[dspace]/bin/dsrun org.dspace.content.crosswalk.XSLTIngestionCrosswalk [-l] plugin input-file
```

..where plugin is the name of the crosswalk plugin to test (e.g. "LOM"), and input-file is a file containing an XML document of metadata in the appropriate format.

Add the -l option to to pass the ingestion crosswalk a list of elements instead of a whole document, as if the List form of the ingest() method had been called. This is needed to test ingesters for formats like DC that get called with lists of elements instead of a root element.

# Chapter 6

## Directories and Files

### 6.1 Source Directory Layout

A complete DSpace installation consists of three separate directory trees:

**The source directory:** This is where (surprise!) the source code lives. Note that the config files here are used only during the initial install process. After the install, config files should be changed in the install directory. It is referred to in this document as [dspace-source].

**The install directory:** This directory is populated during the install process and also by DSpace as it runs. It contains config files, command-line tools (and the libraries necessary to run them), and usually—although not necessarily—the contents of the DSpace archive (depending on how DSpace is configured). After the initial build and install, changes to config files should be made in this directory. It is referred to in this document as [dspace].

**The web deployment directory:** This directory is generated by the web server the first time it finds a dspace.war file in its webapps directory. It contains the unpacked contents of dspace.war, i.e. the JSPs and java classes and libraries necessary to run DSpace. Files in this directory should never be edited directly; if you wish to modify your DSpace installation, you should edit files in the source directory and then rebuild. The contents of this directory aren't listed here since its creation is completely automatic. It is usually referred to in this document as [tomcat]/webapps/dspace.

- dspace-source]
  - build.xml - build file for Ant
  - CHANGES - Detailed list of code changes between versions
  - KNOWN\_BUGS - Known bugs in the current version
  - LICENSE - DSpace source code license
  - README - Obligatory basic information file
  - bin/ - Some shell and Perl scripts for running DSpace command-line tasks
  - config/ - configuration files
    - \* controlled-vocabularies/ - Fixed, limited vocabularies used in metadata entry
    - \* crosswalks/ - Metadata crosswalks - property files or XSL stylesheets
    - \* dspace.cfg - main DSpace configuration file
    - \* dc2mods.cfg - Mappings from Dublin Core metadata to MODS for the METS export
    - \* default.license - default license that users must grant when submitting items

- \* dstat.cfg, dstat.map - statistical report configuration
- \* input-forms.xml - Submission UI metadata field configuration
- \* news-side.html - Text of the front-page news in the sidebar
- \* news-top.html - Text of the front-page news in the top box
- \* emails/ - Text and layout templates for emails sent out by the system
- \* language-packs/ - contains "dictionary files" – Java properties files that contain user interface text in different languages
- \* registries/ - initial contents of the bitstream format registry and Dublin Core element/qualifier registry. These are only used on initial system setup, after which they are maintained in the database.
- \* templates/ - configuration files for libraries and external applications (e.g. Apache, Tomcat) are kept and edited here. They can refer to properties in the main DSpace configuration - have a look at a couple. When they're updated, a command line tool fills out these files with appropriate values from dspace.cfg, and copies them to their appropriate location (hence "templates".)
- docs/ - DSpace system documentation. The technical documentation for functionality, installation, configuration, etc.
- etc/ - miscellaneous stuff need to install DSpace that isn't really to do with system configuration - e.g. the PostgreSQL database schema, and a couple of configuration files that are used during the build process but not by the live system. Also contains the deployment descriptors (web.xml files) for the Web UI and OAI-PMH support .war files.
  - \* oracle/ - versions of the database schema and updater SQL scripts for Oracle
- jsp/ - The Web UI JSPs. As much as possible, these are simply HTML with little bits of Java - the business code resides in the servlets
- lib/ - Library JARs used by the system
  - \* README - Lists the packages third-party libraries (JARs) and their use
  - \* licenses - Contains the licenses associated with the JARs
- src/ - DSpace system source code. For details on how this is laid out, see the overview page of the Javadoc.

Below is the basic layout of a DSpace installation using the default configuration. These paths can be configured if necessary.

- dspace
  - assetstore/ - asset store files
  - bin/ - shell and Perl scripts
  - config/ - configuration, with sub-directories as above
  - handle-server/ - Handles server files
  - history/ - stored history files (generally RDF/XML)
  - lib/ - JARs, including dspace.jar, containing the DSpace classes
  - log/ - Log files
  - reports/ - Reports generated by statistical report generator
  - search/ - Lucene search index files
  - upload/ - temporary directory used during file uploads etc

## 6.2 Contents of Web Application

DSPACE's Ant build file creates a .war file with the following structure:

- (top level dir)
  - The JSPs
  - WEB-INF/
    - \* web.xml - Created from [dSPACE-source]/etc/dSPACE-web.xml, appropriate filled out with configuration parameters (e.g. location of DSPACE installation directory)
    - \* dSPACE-tags.tld - DSPACE custom tag descriptor
    - \* fmt.tld - JSTL message format tag descriptor, for internationalization
    - \* lib/ - All the third-party JARs needed to run DSPACE
    - \* classes/ - The DSPACE class files, plus all of the Messages\_xx.properties files. The latter are placed here so they are in the CLASSPATH and thus can be picked up as resource bundles.

## 6.3 Log Files

The first source of potential confusion is the log files. Since DSPACE uses a number of third-party tools, problems can occur in a variety of places. Below is a table listing the main log files used in a typical DSPACE setup. The locations given are defaults, and might be different for your system depending on where you installed DSPACE and the third-party tools. The ordering of the list is roughly the recommended order for searching them for the details about a particular problem or error.

Log File	What's In It
[dspace]/log/dspace.log	Main DSpace log file. This is where the DSpace code writes a simple log of events and errors that occur within the DSpace code. You can control the verbosity of this by editing the [dspace]/config/templates/log4j.properties file and then running [dspace]/bin/install-configs.
tomcat]/logs/catalina.out	This is where Tomcat's standard output is written. Many errors that occur within the Tomcat code are logged here. For example, if Tomcat can't find the DSpace code (dspace.jar), it would be logged in catalina.out.
[tomcat]/logs/hostname_log.yyyy-mm-dd.txt	If you're running Tomcat stand-alone (without Apache), it logs some information and errors for specific Web applications to this log file. hostname will be your host name (e.g. dspace.myu.edu) and yyyy-mm-dd will be the date.
[tomcat]/logs/apache_log.yyyy-mm-dd.txt	If you're using Apache, Tomcat logs information about Web applications running through Apache (mod_webapp) in this log file (yyyy-mm-dd being the date.)
[apache]/error_log	Apache logs to this file. If there is a problem with getting mod_webapp working, this is a good place to look for clues. Apache also writes to several other log files, though error_log tends to contain the most useful information for tracking down problems.
[dspace]/log/handle-plug.log	The Handle server runs as a separate process from the DSpace Web UI (which runs under Tomcat's JVM). Due to a limitation of log4j's 'rolling file appenders', the DSpace code running in the Handle server's JVM must use a separate log file. The DSpace code that is run as part of a Handle resolution request writes log information to this file. You can control the verbosity of this by editing [dspace]/config/templates/log4j-handle-plugin.properties.
[dspace]/log/handle-server.log	This is the log file for CNRI's Handle server code. If a problem occurs within the Handle server code, before DSpace's plug-in is invoked, this is where it may be logged.
[dspace]/handle-server/error.log	On the other hand, a problem with CNRI's Handle server code might be logged here.
PostgreSQL log	PostgreSQL also writes a log file. This one doesn't seem to have a default location, you probably had to specify it yourself at some point during installation. In general, this log file rarely contains pertinent information—PostgreSQL is pretty stable, you're more likely to encounter problems with connecting via JDBC, and these problems will be logged in dspace.log.

Table 6.1: DSpace Log File Locations



# Chapter 7

## Architecture

### 7.1 Overview

The DSpace system is organized into three layers, each of which consists of a number of components.

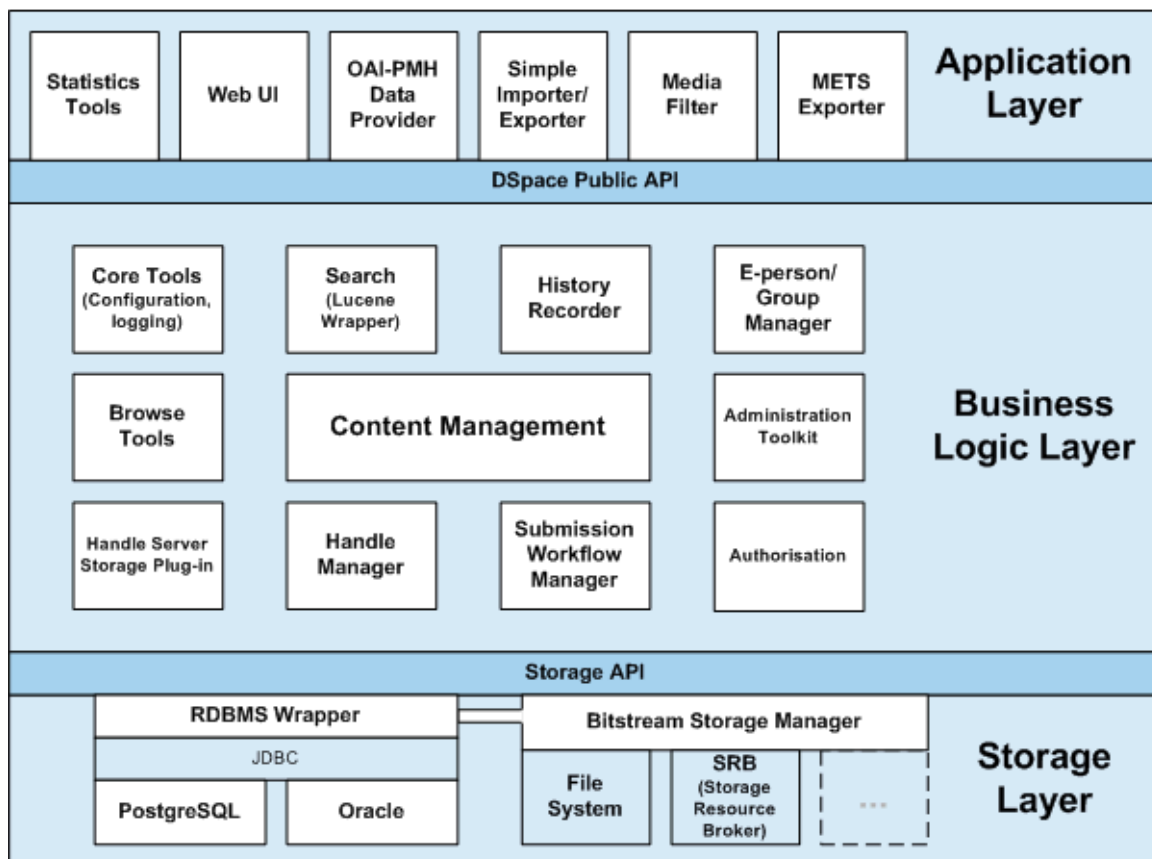


Figure 7.1: DSpace System Architecture

The storage layer is responsible for physical storage of metadata and content. The business logic layer deals with managing the content of the archive, users of the archive (e-people), authorization, and workflow. The application layer contains components that communicate with the world outside of the individual DSpace installation, for example the Web user interface and the [Open Archives Initiative](#) protocol for metadata harvesting service.

Each layer only invokes the layer below it; the application layer may not use the storage layer directly, for example. Each component in the storage and business logic layers has a defined public API. The union of the APIs of those components are referred to as the Storage API (in the case of the storage layer) and the DSpace Public API (in the case of the business logic layer). These APIs are in-process Java classes, objects and methods.

It is important to note that each layer is trusted. Although the logic for authorising actions is in the business logic layer, the system relies on individual applications in the application layer to correctly and securely authenticate e-people. If a 'hostile' or insecure application were allowed to invoke the Public API directly, it could very easily perform actions as any e-person in the system.

The reason for this design choice is that authentication methods will vary widely between different applications, so it makes sense to leave the logic and responsibility for that in these applications.

The source code is organized to cohere very strictly to this three-layer architecture. Also, only methods in a component's public API are given the public access level. This means that the Java compiler helps ensure that the source code conforms to the architecture.

Packages within	Correspond to components in
org.dspace.app	Application layer
org.dspace	Business logic layer (except storage and app)
org.dspace.storage	Storage layer

Table 7.1: Source Code Packages

The storage and business logic layer APIs are extensively documented with Javadoc-style comments. Generate the HTML version of these by entering the source directory and running:

```
ant public_api
```

The package-level documentation of each package usually contains an overview of the package and some example usage. This information is not repeated in this architecture document; this and the Javadoc APIs are intended to be used in parallel.

## 7.2 Storage Layer

### 7.2.1 RDBMS

DSpace uses a relational database to store all information about the organization of content, metadata about the content, information about e-people and authorization, and the state of currently-running workflows. The DSpace system also uses the relational database in order to maintain indices that users can browse. [Graphical visualization of the relational database.](#)

Most of the functionality that DSpace uses can be offered by any standard SQL database that supports transactions. Presently, the browse indices use some features specific to [PostgreSQL](#) and [Oracle](#), so some modification to the code would be needed before DSpace would function fully with an alternative database back-end.

The org.dspace.storage.rdbms package provides access to an SQL database in a somewhat simpler form than using JDBC directly. The main class is DatabaseManager, which executes SQL queries and returns TableRow

or `TableRowIterator` objects. The `InitializeDatabase` class is used to load SQL into the database via JDBC, for example to set up the schema.

All calls to the Database Manager require a **DSpace Context object**. Example use of the database manager API is given in the `org.dspace.storage.rdbms` package Javadoc.

The database schema used by DSpace (for PostgreSQL) is stored in `[dspace-source]/etc/database_schema.sql` in the source distribution. It is stored in the form of SQL that can be fed straight into the DBMS to construct the database. The schema SQL file also directly creates two e-person groups in the database that are required for the system to function properly.

Also in `[dspace-source]/etc` are various SQL files called `database_schema_1x_1y`. These contain the necessary SQL commands to update a live DSpace database from version 1.x to 1.y. Note that this might not be the only part of an upgrade process: see **Updating a DSpace Installation** for details.

The DSpace database code uses an SQL function `getnextid` to assign primary keys to newly created rows. This SQL function must be safe to use if several JVMs are accessing the database at once; for example, the Web UI might be creating new rows in the database at the same time as the batch item importer. The PostgreSQL-specific implementation of the method uses SEQUENCES for each table in order to create new IDs. If an alternative database backend were to be used, the implementation of `getnextid` could be updated to operate with that specific DBMS.

The `etc` directory in the source distribution contains two further SQL files. `clean-database.sql` contains the SQL necessary to completely clean out the database, so use with caution! The Ant target `clean_database` can be used to execute this. `update-sequences.sql` contains SQL to reset the primary key generation sequences to appropriate values. You'd need to do this if, for example, you're restoring a backup database dump which creates rows with specific primary keys already defined. In such a case, the sequences would allocate primary keys that were already used.

Versions of the `*.sql*` files for Oracle are stored in `[dspace-source]/etc/oracle`. These need to be copied over their PostgreSQL counterparts in `[dspace-source]/etc` prior to installation.

## Maintenance and Backup

When using PostgreSQL, it's a good idea to perform regular 'vacuuming' of the database to optimize performance. This is performed by the `vacuumdb` command which can be executed via a 'cron' job, for example by putting this in the system crontab:

```
# clean up the database nightly
40 2 * * * /usr/local/pgsql/bin/vacuumdb --analyze dspace > /dev/null 2>&1
```

The DSpace database can be backed up and restored using usual methods, for example with `pg_dump` and `psql`. However when restoring a database, you will need to perform these additional steps:

- The `fresh_install` target loads up the initial contents of the Dublin Core type and bitstream format registries, as well as two entries in the `epersongroup` table for the system anonymous and administrator groups. Before you restore a raw backup of your database you will need to remove these, since they will already exist in your backup, possibly having been modified. For example, use:

```
DELETE FROM dctyperegistry;
DELETE FROM bitstreamformatregistry;
DELETE FROM epersongroup;
```

- After restoring a backup, you will need to reset the primary key generation sequences so that they do not produce already-used primary keys. Do this by executing the SQL in `[dspace-source]/etc/update-sequences.sql`, for example with:

```
psql -U dspace -f [dspace-source]/etc/update-sequences.sql
```

Future updates of DSpace may involve minor changes to the database schema. Specific instructions on how to update the schema whilst keeping live data will be included. The current schema also contains a few currently unused database columns, to be used for extra functionality in future releases. These unused columns have been added in advance to minimize the effort required to upgrade.

### Configuring the RDBMS Component

The database manager is configured with the following properties in `dspace.cfg`:

**db.url** The JDBC URL to use for accessing the database. This should not point to a connection pool, since DSpace already implements a connection pool.

**db.driver** JDBC driver class name. Since presently, DSpace uses PostgreSQL-specific features, this should be `org.postgresql.Driver`.

**db.username** Username to use when accessing the database.

**db.password** Corresponding password to use when accessing the database.

### 7.2.2 Bitstream Store

DSpace offers two means for storing content. The first is in the file system on the server. The second is using SRB (Storage Resource Broker). Both are achieved using a simple, lightweight API. SRB is purely an option but may be used in lieu of the server's file system or in addition to the file system. Without going into a full description, SRB is a very robust, sophisticated storage manager that offers essentially unlimited storage and straightforward means to replicate (in simple terms, backup) the content on other local or remote storage resources.

The terms "store", "retrieve", "in the system", "storage", and so forth, used below can refer to storage in the file system on the server ("traditional") or in SRB.

The `BitstreamStorageManager` provides low-level access to bitstreams stored in the system. In general, it should not be used directly; instead, use the `Bitstream` object in the content management API since that encapsulated authorization and other metadata to do with a bitstream that are not maintained by the `BitstreamStorageManager`.

The bitstream storage manager provides three methods that store, retrieve and delete bitstreams. Bitstreams are referred to by their 'ID'; that is the primary key `bitstream_id` column of the corresponding row in the database.

As of DSpace version 1.1, there can be multiple bitstream stores. Each of these bitstream stores can be traditional storage or SRB storage. This means that the potential storage of a DSpace system is not bound by the maximum size of a single disk or file system and also that traditional and SRB storage can be combined in one DSpace installation. Both traditional and SRB storage are specified by configuration parameters. Also see `Configuring the Bitstream Store` below.

Stores are numbered, starting with zero, then counting upwards. Each bitstream entry in the database has a store number, used to retrieve the bitstream when required.

At the moment, the store in which new bitstreams are placed is decided using a configuration parameter, and there is no provision for moving bitstreams between stores. Administrative tools for manipulating bitstreams and stores will be provided in future releases. Right now you can move a whole store (e.g. you could move store number 1 from `/localdisk/store` to `/fs/anotherdisk/store` but it would still have to be store number 1 and have the exact same contents.

Bitstreams also have an 38-digit internal ID, different from the primary key ID of the bitstream table row. This is not visible or used outside of the bitstream storage manager. It is used to determine the exact location (relative to the relevant store directory) that the bitstream is stored in traditional or SRB storage. The first three pairs of digits are the directory path that the bitstream is stored under. The bitstream is stored in a file with the internal ID as the filename.

For example, a bitstream with the internal ID 12345678901234567890123456789012345678 is stored in the directory:

```
(assetstore dir)/12/34/56/12345678901234567890123456789012345678
```

The reasons for storing files this way are

- Using a randomly-generated 38-digit number means that the 'number space' is less cluttered than simply using the primary keys, which are allocated sequentially and are thus close together. This means that the bitstreams in the store are distributed around the directory structure, improving access efficiency.
- The internal ID is used as the filename partly to avoid requiring an extra lookup of the filename of the bitstream, and partly because bitstreams may be received from a variety of operating systems. The original name of a bitstream may be an illegal UNIX filename.

When storing a bitstream, the BitstreamStorageManager DOES set the following fields in the corresponding database table row:

- bitstream\_id
- size
- checksum
- checksum\_algorithm
- internal\_id
- deleted
- store\_number

The remaining fields are the responsibility of the Bitstream content management API class.

The bitstream storage manager is fully transaction-safe. In order to implement transaction-safety, the following algorithm is used to store bitstreams:

1. A database connection is created, separately from the currently active connection in the current DSpace context.
2. An unique internal identifier (separate from the database primary key) is generated.
3. The bitstream DB table row is created using this new connection, with the deleted column set to true.
4. The new connection is committed, so the 'deleted' bitstream row is written to the database
5. The bitstream itself is stored in a file in the configured 'asset store directory', with a directory path and filename derived from the internal ID
6. The deleted flag in the bitstream row is set to false. This will occur (or not) as part of the current DSpace Context.

This means that should anything go wrong before, during or after the bitstream storage, only one of the following can be true:

- No bitstream table row was created, and no file was stored
- A bitstream table row with `deleted=true` was created, no file was stored
- A bitstream table row with `deleted=true` was created, and a file was stored

None of these affect the integrity of the data in the database or bitstream store.

Similarly, when a bitstream is deleted for some reason, its `deleted` flag is set to `true` as part of the overall transaction, and the corresponding file in storage is not deleted.

The above techniques mean that the bitstream storage manager is transaction-safe. Over time, the bitstream database table and file store may contain a number of 'deleted' bitstreams. The `cleanup` method of `BitstreamStorageManager` goes through these deleted rows, and actually deletes them along with any corresponding files left in the storage. It only removes 'deleted' bitstreams that are more than one hour old, just in case cleanup is happening in the middle of a storage operation.

This cleanup can be invoked from the command line via the `Cleanup` class, which can in turn be easily executed from a shell on the server machine using `/dspace/bin/cleanup`. You might like to have this run regularly by cron, though since DSpace is read-lots, write-not-so-much it doesn't need to be run very often.

### 7.2.3 Backup

The bitstreams (files) in traditional storage may be backed up very easily by simply 'tarring' or 'zipping' the `assetstore` directory (or whichever directory is configured in `dspace.cfg`). Restoring is as simple as extracting the backed-up compressed file in the appropriate location.

Similar means could be used for SRB, but SRB offers many more options for managing backup.

It is important to note that since the bitstream storage manager holds the bitstreams in storage, and information about them in the database, that a database backup and a backup of the files in the bitstream store must be made at the same time; the bitstream data in the database must correspond to the stored files.

Of course, it isn't really ideal to 'freeze' the system while backing up to ensure that the database and files match up. Since DSpace uses the bitstream data in the database as the authoritative record, it's best to back up the database before the files. This is because it's better to have a bitstream in storage but not the database (effectively non-existent to DSpace) than a bitstream record in the database but not storage, since people would be able to find the bitstream but not actually get the contents.

#### Configuring the Bitstream Store

Both traditional and SRB bitstream stores are configured in `dspace.cfg`.

**Configuring Traditional Storage** Bitstream stores in the file system on the server are configured like this:

```
assetstore.dir = [dspace]/assetstore
```

(Remember that `[dspace]` is a placeholder for the actual name of your DSpace install directory).

The above example specifies a single asset store.

```
assetstore.dir = [dspace]/assetstore_0
assetstore.dir.1 = /mnt/other_filesystem/assetstore_1
```

The above example specifies two asset stores. `assetstore.dir` specifies the asset store number 0 (zero); after that use `assetstore.dir.1`, `assetstore.dir.2` and so on. The particular asset store a bitstream is stored in is held in the database, so don't move bitstreams between asset stores, and don't renumber them.

By default, newly created bitstreams are put in asset store 0 (i.e. the one specified by the `assetstore.dir` property.) This allows backwards compatibility with pre-DSpace 1.1 configurations. To change this, for example when asset store 0 is getting full, add a line to `dspace.cfg` like:

```
assetstore.incoming = 1
```

Then restart DSpace (Tomcat). New bitstreams will be written to the asset store specified by `assetstore.dir.1`, which is `/mnt/other_filesystem/assetstore_1` in the above example.

**Configuring SRB Storage** The same framework is used to configure SRB storage. That is, the asset store number (0..n) can reference a file system directory as above or it can reference a set of SRB account parameters. But any particular asset store number can reference one or the other but not both. This way traditional and SRB storage can both be used but with different asset store numbers. The same cautions mentioned above apply to SRB asset stores as well: The particular asset store a bitstream is stored in is held in the database, so don't move bitstreams between asset stores, and don't renumber them.

For example, let's say asset store number 1 will refer to SRB. Then there will be a set of SRB account parameters like this:

```
srb.host.1 = mysrbmcatzhost.myu.edu
srb.port.1 = 5544
srb.mcatzzone.1 = mysrbzone
srb.mdasdomainname.1 = mysrbdomain
srb.defaultstorageresource.1 = mydefaultsrbresource
srb.username.1 = mysrbuser
srb.password.1 = mysrbpassword
srb.homedirectory.1 = /mysrbzone/home/mysrbuser.mysrbdomain
srb.parentdir.1 = mysrbdspaceassetstore
```

Several of the terms, such as `mcatzzone`, have meaning only in the SRB context and will be familiar to SRB users. The last, `srb.parentdir.n`, can be used to add (SRB) upper directory structure within an SRB account. This property value could be blank as well.

(If asset store 0 would refer to SRB it would be `srb.host = ...`, `srb.port = ...`, and so on (.0 omitted) to be consistent with the traditional storage configuration above.)

The similar use of `assetstore.incoming` to reference asset store 0 (default) or 1..n (explicit property) means that new bitstreams will be written to traditional or SRB storage determined by whether a file system directory on the server is referenced or a set of SRB account parameters are referenced.

There are comments in `dspace.cfg` that further elaborate the configuration of traditional and SRB storage.

## 7.3 Business Logic Layer

### 7.3.1 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, remember to edit the file in `/dspace/config/templates` and then run `/dspace/bin/install-configs` rather than editing the 'live' version directly! The `ConfigurationManager` class can also be invoked as a command line tool, with two possible uses:

- `/dspace/bin/install-configs` This processes and installs configuration files for other applications, as described in the configuration section.
- `/dspace/bin/dsrun org.dspace.core.ConfigurationManager -property 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. For an example, see `/dspace/bin/start-handle-server`. 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`.

### 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 analysed 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/templates/log4j.properties` and then executing `/dspace/bin/install-configs`.

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=1721.1/1686
```

This breaks down like this: The above format allows the logs to be easily parsed and analysed. The `/dspace/bin/log-reporter` script is a simple tool for analysing 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 methods that are required in a variety of places throughout the code, and thus have no particular 'home' in a subsystem.

Date and time, milliseconds	2002-11-11 08:11:32,903
Level (FATAL, WARN, INFO or DEBUG)	INFO
Java class	org.dspace.app.webui.servlet.DSpaceServlet @
User email or anonymous	anonymous :
Extra log info from context	session_id=BD84E7C194C2CF4BD0EC3A6CAD0142BB :
Action	view_item :
Extra info	handle=1721.1/1686

Table 7.2: Example Logging Entry

### 7.3.2 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.

Items 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 [meta]data 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:

```
Context context = new Context();
Bitstream b = Bitstream.find(context, 1234);
b.setName("newfile.txt");
b.update();
context.complete();
```

**Will change storage.**

```
Context context = new Context();
Bitstream b = Bitstream.find(context, 1234);
b.setName("newfile.txt");
b.update();
context.abort();
```

**Will not change storage (context aborted)**

```
Context context = new Context();
Bitstream b = Bitstream.find(context, 1234);
b.setName("newfile.txt");
context.complete();
```

The new name **will not** be stored since update was not invoked

```
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 BitstreamFormats) to be instantiated.

Instantiating an Item object causes the appropriate Bundle objects (etc.) and hence BitstreamFormats 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 DCValue class is a simple container that represents a single Dublin Core element, optional qualifier, value and language. Note that as of 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
date	Any or unqualified	ISO 8601 in the UTC time zone, with either year, month, day, or second precision. Examples:  2000 2002-10 2002-08-14 1999-01-01T14:35:23Z	DCDate
contributor	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:  Doe, John Smith, John Jr. van Dyke, Dick Massachusetts Institute of Technology	DCPersonName
language	iso	A two letter code taken ISO 639, followed optionally by a two letter country code taken from ISO 3166. Examples:  en fr en_US	DCLanguage
relation	ispartofseries	The series name, following by a semicolon followed by the number in that series. Alternatively, just free text.  MIT-TR; 1234 My Report Series; ABC-1234 NS1234	DCSeriesNumber

Table 7.3: Dublin Core API

### 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 `DCValue` 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);
```

### 7.3.3 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 lifecycle 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 implementations 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 Objects 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 intface)` 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 `PluginConfigurationException` 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 no 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.MediaFilter`
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 have 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 guide](#) for further details.

### 7.3.4 Workflow System

The primary classes are: The workflow system models the states of an Item in a state machine with 5 states

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

Table 7.4: Workflow System - primary classes

(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.

### 7.3.5 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/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.DCValue` class. Only administrators may modify the Dublin Core type registry.

The `org.dspace.administer.RegistryLoader` class contains methods for initialising 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.

### 7.3.6 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.

### 7.3.7 Authorization

The primary classes are: The authorization system is based on the classic 'police state' model of security; no

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

Table 7.5: Authorization - primary classes

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.

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.

### 7.3.8 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.hdlib.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 held as a template in `/dspace/config/templates/log4j-handle-plugin.properties`, written to `/dspace/config/log4j-handle-plugin.properties` by the `install-configs` script. The `/dspace/bin/start-handle-server` script passes in the appropriate command line parameters so that the Handle server uses this configuration.

### 7.3.9 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 is invoked by the `dspace/bin/index-all` 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 subelements. 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.

#### Our 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:

#### 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.

Search Field	Taken from Dublin Core Fields
Authors	contributor.* creator.*
Titles	description.statementsofresponsibility title.*
Keywords	subject.*
Abstracts	description.abstract description.tableofcontents
Series	relation.ispartofseries
MIME types	format.mimetype
Sponsors	description.sponsorship
Identifiers	identifier.*

Table 7.6: Indexed Fields

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.

### 7.3.10 Browse API

The browse API maintains indices 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 SystemAppears 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 cataloguers. For these reasons, DSpace does not yet feature 'authority control' functionality.

**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 flavour 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 indices 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 indices from scratch. In general the content management API invokes the necessary browse API calls to keep the browse indices 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/index-all` shell script) that causes the indices 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, normalising 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.

### 7.3.11 History Recorder

The purpose of the history subsystem is to capture a time-based record of significant changes in DSpace, in a manner suitable for later refactoring or repurposing. Note that the history data is not expected to provide current information about the archive; it simply records what has happened in the past.

The [Harmony project](#) describes a simple and powerful approach for modeling temporal data. The DSpace history framework adopts this model. The Harmony model is used by the serialization mechanism (and ultimately by agents who interpret the serializations); users of the History API need not be aware of it. The content management API handles invocations of the history system. Users of the DSpace public API do not generally need to use the history API.

When anything of archival interest occurs in DSpace, the `saveHistory` method of the `HistoryManager` is invoked. The parameters contains a reference to anything of archival interest. Upon reception of the object, it serializes the state of all archive objects referred to by it, and creates Harmony-style objects and associations to describe the relationships between the objects. (A simple example is given below). Note that each archive object must have a unique identifier to allow linkage between discrete events; this is discussed under "Unique IDs" below. The serializations (including the Harmony objects and associations) are persisted as files in the `/dspace/history` (or other configured) directory. The `history` and `historystate` tables contain simple indices into the serializations in the file system.

## Archival Events

The following events are significant enough to warrant history records:

- Communities
  - create/modify/delete
  - add/remove Collection to/from Community

- Collections
  - create/modify/delete
  - add/remove Item to/from Collection
- Items
  - create/modify/delete
  - assign Handle to Item
  - modify Item contents (Bundles, Bitstreams, metadata fields, etc)
- EPerson
  - create/modify/delete
- Workflow
  - Workflow completed

### Serializations

The serialization of an archival object consists of:

- Its instance fields (ie, non-static, non-transient fields)
- The serializations of associated objects (or references to these serializations).

### Unique Ids

To be able to trace the history of an object, it is essential that the object have a unique identifier. Since not all objects in the system have Handles, the unique identifiers are only weakly tied to the Handle system. Instead, the identifier consists of:

- an identifier for the project
- a site id (using the handle prefix)
- an RDBMS-based id for objects

### Storage

When an archive object is serialized, an object ID and MD5 checksum are recorded. When another object is serialized, the checksum for the serialization is matched against existing checksums for that object. If the checksum already exists, the object is not stored; a reference to the object is used instead. Note that since none of the serializations are deleted, reference counting is unnecessary.

The history data is not initially stored in a queryable form. Two simple RDBMS tables give basic indications of what is stored, and where. The history table is an index of serializations with checksums and dates. The history\_id column corresponds to the file in which a serialization is stored. For example, if the history ID is 123456, it will be stored in the file:

```
/dspace/history/00/12/34/123456
```

The table also contains the date the serialization was written and the MD5 checksum of the serialization. The historystate table is supposed to indicate the most recent serialization of any given object.

**Example**

An item is submitted to a collection via bulk upload. When (and if) the item is eventually added to the collection, the history method is called, with references to the item, its collection, the e-person who performed the bulk upload, and some indication of the fact that it was submitted via a bulk upload.

When called, the HistoryManager does the following: It creates the following new resources (all with unique ids):

- An event
- A state
- An action

It also generates the following relationships

```

event --atTime-->    time
event --hasOutput--> state
Item  --inState-->   state
state --contains-->  Item
action --creates-->  Item
event --hasAction--> action
action --usesTool--> DSpace Upload
action --hasAgent--> User

```

The history component serializes the state of all archival objects involved (in this case, the item, the e-person, and the collection). It creates entries in the history database tables which associate the archival objects with the generated serializations.

**Caveats**

This history system is a largely untested experiment. It also needs further documentation. There have been no serious efforts to determine whether the information written by the history system, either to files or the database tables, is accurate. In particular, the historystate table does not seem to be correctly written.

**7.3.12 Checksum checker**

The architecture of the checker is documented in the package javadocs, run `ant public_api`, and look in `build/public_api/index.html`.

**7.4 Application Layer****7.4.1 Web User Interface**

The DSpace Web UI is the largest and most-used component in the application layer. Built on Java Servlet and JavaServer Page technology, it allows end-users to access DSpace over the Web via their Web browsers. As of Dspace 1.3.2 the UI meets both XHTML 1.0 standards and Web Accessibility Initiative (WAI) level-2 standard. It also features an administration section, consisting of pages intended for use by central administrators. Presently, this part of the Web UI is not particularly sophisticated; users of the administration section need to know what they are doing! Selected parts of this may also be used by collection [FIXME: administrators or editors?]



## Web UI Files

The Web UI-related files are located in a variety of directories in the DSpace source tree. Note that as of DSpace version 1.2, the deployment mechanism has changed; the build process creates easy-to-deploy Web application archives (.war files).

Location	Description
org.dspace.app.webui	Web UI source files
org.dspace.app.webui.filter	Servlet Filters (Servlet 2.3 spec)
org.dspace.app.webui.jsptag	Custom JSP tag class files
org.dspace.app.webui.servlet	Servlets for main Web UI (controllers)
org.dspace.app.webui.util	Miscellaneous classes used by the servlets and filters
[dspace-source]/jsp	The JSP files
[dspace-source]/jsp/local	This is where you can place customized versions of JSPs
[dspace-source]/jsp/WEB-INF/dspace-tags.tld	Custom DSpace JSP tag descriptor
[dspace-source]/etc/dspace-web.xml	The Web application deployment descriptor. Before including in the .war file, the text @@dspace.dir@@ will be replaced with the DSpace installation directory (referred to as [dspace] elsewhere in this system documentation). This allows the Web application to pick up the DSpace configuration and environment.

Table 7.7: Locations of Web UI Source Files

## The Build Process

The DSpace build process constructs a Web application archive, which is placed in [dspace-source]/build/dspace.war. The build\_wars Ant target does the work. The process works as follows:

- All the DSpace source code is compiled.
- [dspace-source]/etc/dspace-web.xml is copied to [dspace-source]/build and the @@dspace.dir@@ token inside it replaced with the DSpace installation directory (dspace.dir property from dspace.cfg)
- The JSPs are all copied to [dspace-source]/build/jsp
- Customized JSPs from [dspace-source]/jsp/local are copied on top of these, thus 'overriding' the default versions
- [dspace-source]/build/dspace.war is built

The contents of dspace.war are:

- (Top level) – the JSPs (customized versions from [dspace-source]/jsp/local will have overwritten the defaults from the DSpace source distribution)
- WEB-INF/classes – the compiled DSpace classes
- WEB-INF/lib – the third party library JAR files from [dspace-source]/lib, minus servlet.jar which will be available as part of Tomcat (or other servlet engine)
- WEB-INF/web.xml – web deployment descriptor, copied from [dspace-source]/build/dspace-web.xml

- WEB-INF/dspace-tags.tld – tag descriptor

Note that this does mean there are multiple copies of the compiled DSpace code and third-party libraries in the system, so care must be taken to ensure that they are all in sync. (The storage overhead is a few megabytes, totally insignificant these days.) In general, when you change any DSpace code or JSP, it's best to do a complete update of both the installation ([dspace]), and to rebuild and redeploy the Web UI and OAI .war files, by running this in [dspace-source]:

```
ant -D[dspace]/config/dspace.cfg update
```

and then following the instructions that command writes to the console.

### Servlets and JSPs

The Web UI is loosely based around the MVC (model, view, controller) model. The content management API corresponds to the model, the Java Servlets are the controllers, and the JSPs are the views. Interactions take the following basic form:

1. An HTTP request is received from a browser
2. The appropriate servlet is invoked, and processes the request by invoking the DSpace business logic layer public API
3. Depending on the outcome of the processing, the servlet invokes the appropriate JSP
4. The JSP is processed and sent to the browser

The reasons for this approach are:

- All of the processing is done before the JSP is invoked, so any error or problem that occurs does not occur halfway through HTML rendering
- The JSPs contain as little code as possible, so they can be customized without having to delve into Java code too much

The `org.dspace.app.webui.servlet.LoadDSpaceConfig` servlet is always loaded first. This is a very simple servlet that checks the `dspace-config` context parameter from the DSpace deployment descriptor, and uses it to locate `dspace.cfg`. It also loads up the Log4j configuration. It's important that this servlet is loaded first, since if another servlet is loaded up, it will cause the system to try and load DSpace and Log4j configurations, neither of which would be found.

All DSpace servlets are subclasses of the `DSpaceServlet` class. The `DSpaceServlet` class handles some basic operations such as creating a DSpace Context object (opening a database connection etc.), authentication and error handling. Instead of overriding the `doGet` and `doPost` methods as one normally would for a servlet, DSpace servlets implement `doDSGet` or `doDSPost` which have an extra context parameter, and allow the servlet to throw various exceptions that can be handled in a standard way.

The DSpace servlet processes the contents of the HTTP request. This might involve retrieving the results of a search with a query term, accessing the current user's eperson record, or updating a submission in progress. According to the results of this processing, the servlet must decide which JSP should be displayed. The servlet then fills out the appropriate attributes in the `HttpRequest` object that represents the HTTP request being processed. This is done by invoking the `setAttribute` method of the `javax.servlet.http.HttpServletRequest` object that is passed into the servlet from Tomcat. The servlet then forwards control of the request to the appropriate JSP using the `JSPManager.showJSP` method.

The `JSPManager.showJSP` method uses the standard Java servlet forwarding mechanism is then used to forward the HTTP request to the JSP. The JSP is processed by Tomcat and the results sent back to the user's browser. There is an exception to this servlet/JSP style: `index.jsp`, the 'home page', receives the HTTP request directly from Tomcat without a servlet being invoked first. This is because in the servlet 2.3 specification, there is no way to map a servlet to handle only requests made to `'/'`; such a mapping results in every request being directed to that servlet. By default, Tomcat forwards requests to `'/'` to `index.jsp`. To try and make things as clean as possible, `index.jsp` contains some simple code that would normally go in a servlet, and then forwards to `home.jsp` using the `JSPManager.showJSP` method. This means localized versions of the 'home page' can be created by placing a customized `home.jsp` in `[dspace-source]/jsp/local`, in the same manner as other JSPs. `[dspace-source]/jsp/dspace-admin/index.jsp`, the administration UI index page, is invoked directly by Tomcat and not through a servlet for similar reasons.

At the top of each JSP file, right after the license and copyright header, is documented the appropriate attributes that a servlet must fill out prior to forwarding to that JSP. No validation is performed; if the servlet does not fill out the necessary attributes, it is likely that an internal server error will occur.

Many JSPs containing forms will include hidden parameters that tell the servlets which form has been filled out. The submission UI servlet (`SubmitServlet` is a prime example of a servlet that deals with the input from many different JSPs. The step hidden parameter is used to inform the servlet which form has been filled out (which step of submission the user has just completed.)

Below is a detailed, scary diagram depicting the flow of control during the whole process of processing and responding to an HTTP request. More information about the authentication mechanism is mostly described in the configuration section.

### Custom JSP Tags

The DSpace JSPs all use some custom tags defined in `/dspace/jsp/WEB-INF/dspace-tags.tld`, and the corresponding Java classes reside in `org.dspace.app.webui.jsptag`. The tags are listed below. The `dspace-tags.tld` file contains detailed comments about how to use the tags, so that information is not repeated here.

**layout** Just about every JSP uses this tag. It produces the standard HTML header and `<BODY>` tag. Thus the content of each JSP is nested inside a `<dspace:layout>` tag. The (XML-style) attributes of this tag are slightly complicated—see `dspace-tags.tld`. The JSPs in the source code bundle also provide plenty of examples.

**sidebar** Can only be used inside a layout tag, and can only be used once per JSP. The content between the start and end sidebar tags is rendered in a column on the right-hand side of the HTML page. The contents can contain further JSP tags and Java 'scriptlets'.

**date** Displays the date represented by an `org.dspace.content.DCDate` object. Just the one representation of date is rendered currently, but this could use the user's browser preferences to display a localized date in the future.

**include** Obsolete, simple tag, similar to `jsp:include`. In versions prior to DSpace 1.2, this tag would use the locally modified version of a JSP if one was installed in `jsp/local`. As of 1.2, the build process now performs this function, however this tag is left in for backwards compatibility.

**item** Displays an item record, including Dublin Core metadata and links to the bitstreams within it. Note that the displaying of the bitstream links is simplistic, and does not take into account any of the bundling structure. This is because DSpace does not have a fully-fledged dissemination architectural piece yet. Displaying an item record is done by a tag rather than a JSP for two reasons: Firstly, it happens in several places (when verifying an item record during submission or workflow review, as well as during standard

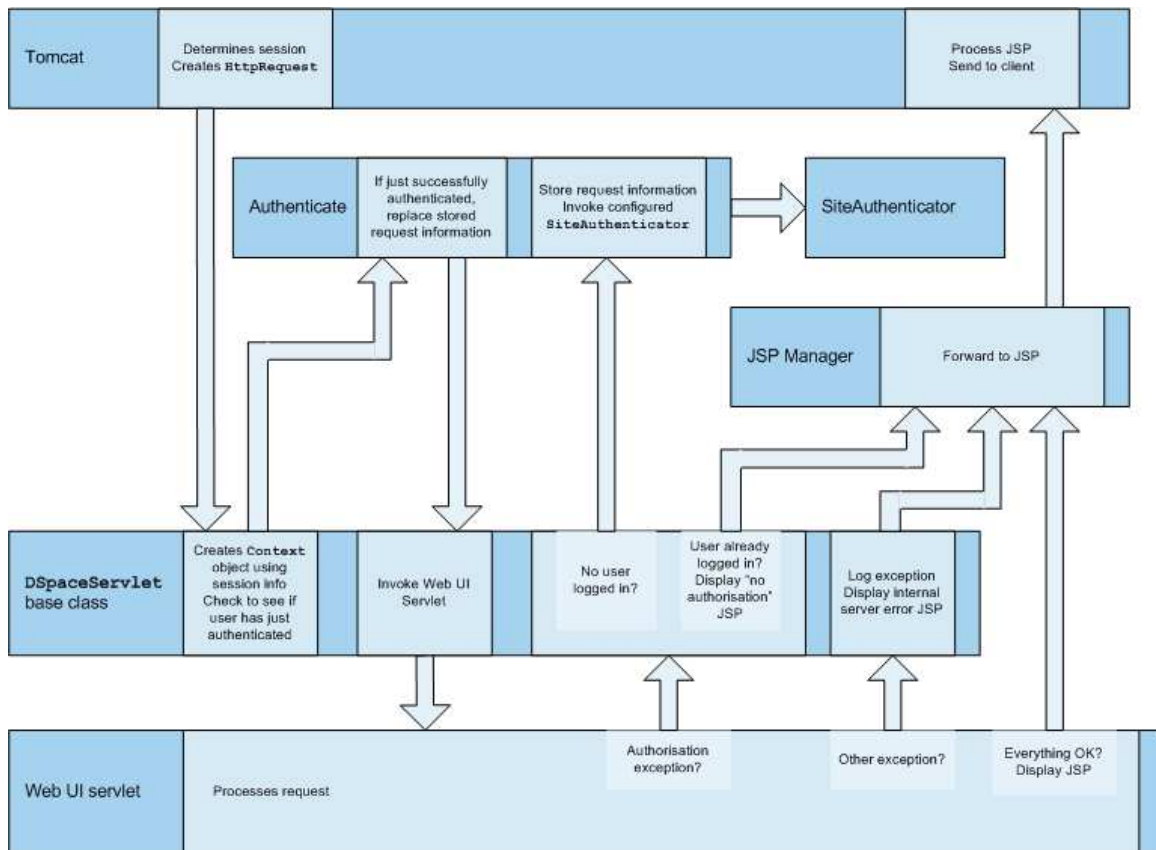


Figure 7.2: Flow of Control During HTTP Request Processing

item accesses), and secondly, displaying the item turns out to be mostly code-work rather than HTML anyway. Of course, the disadvantage of doing it this way is that it is slightly harder to customize exactly what is displayed from an item record; it is necessary to edit the tag code (`org.dspace.app.webui.jsp.tag.ItemTag`). Hopefully a better solution can be found in the future.

**itemlist, collectionlist, communitylist** These tags display ordered sequences of items, collections and communities, showing minimal information but including a link to the page containing full details. These need to be used in HTML tables.

**popup** This tag is used to render a link to a pop-up page (typically a help page.) If Javascript is available, the link will either open or pop to the front any existing DSpace pop-up window. If Javascript is not available, a standard HTML link is displayed that renders the link destination in a window named 'dspace.popup'. In graphical browsers, this usually opens a new window or re-uses an existing window of that name, but if a window is re-used it is not 'raised' which might confuse the user. In text browsers, following this link will simply replace the current page with the destination of the link. This obviously means that Javascript offers the best functionality, but other browsers are still supported.

**selectperson** A tag which produces a widget analogous to HTML `<SELECT>`, that allows a user to select one or multiple e-people from a pop-up list.

**sfxlink** Using an item's Dublin Core metadata DSpace can display an SFX link, if an SFX server is available. This tag does so for a particular item if the `sfx.server.url` property is defined in `dspace.cfg`.

### Internationalisation

The [Java Standard Tag Library v1.0](#) is used to specify messages in the JSPs like this:

OLD:

```
<H1>Search Results</H1>
```

NEW:

```
<H1><fmt:message key="jsp.search.results.title" /></H1>
```

Phrases may have parameters to be passed in, to make the job of translating easier, reduce the number of 'keys' and to allow translators to make the translated text flow more appropriately for the target language.

OLD:

```
<P>Results <%= r.getFirst() %> to <%= r.getLast() %> of <%= r.getTotal() %></P>
```

NEW:

```
<fmt:message key="jsp.search.results.text">
  <fmt:param><%= r.getFirst() %></fmt:param>
  <fmt:param><%= r.getLast() %></fmt:param>
  <fmt:param><%= r.getTotal() %></fmt:param>
</fmt:message>
```

(Note: JSTL 1.0 does not seem to allow JSP `<%= %>` expressions to be passed in as values of attribute in `<fmt:param value="" />`)

The above would appear in the `Messages_xx.properties` file as:

```
jsp.search.results.text = Results {0}-{1} of {2}
```

Introducing number parameters that should be formatted according to the locale used makes no difference in the message key compared to string parameters:

```
jsp.submit.show-uploaded-file.size-in-bytes = {0} bytes
```

In the JSP using this key can be used in the way below:

```
<fmt:message key="jsp.submit.show-uploaded-file.size-in-bytes">
  <fmt:param><fmt:formatNumber><%= bitstream.getSize() %></fmt:formatNumber></fmt:param>
</fmt:message>
```

(Note: JSTL offers a way to include numbers in the message keys as `jsp.foo.key = 0,number bytes`. Setting the parameter as `<fmt:param value=""$variable" />` works when `variable` is a single variable name and doesn't work when trying to use a method's return value instead: `bitstream.getSize()`. Passing the number as string (or using the `<%= %>` expression) also does not work.)

Multiple `Messages.properties` can be created for different languages. See `ResourceBundle.getBundle`. e.g. you can add German and Canadian French translations:

```
Messages_de.properties
Messages_fr_CA.properties
```

The end user's browser settings determine which language is used. The English language file `Messages.properties` (or the default server locale) will be used as a default if there's no language bundle for the end user's preferred language. (Note that the English file is not called `Messages_en.properties` – this is so it is always available as a default, regardless of server configuration.)

The `dspace:layout` tag has been updated to allow dictionary keys to be passed in for the titles. It now has two new parameters: `titlekey` and `parenttitlekey`. So where before you'd do:

```
<dspace:layout title="Here"
               parentlink="/myspace"
               parenttitle="My DSpace">
```

You now do:

```
<dspace:layout titlekey="jsp.page.title"
               parentlink="/myspace"
               parenttitlekey="jsp.myspace">
```

And so the layout tag itself gets the relevant stuff out of the dictionary. `title` and `parenttitle` still work as before for backwards compatibility, and the odd spot where that's preferable.

**Message Key Convention** When translating further pages, please follow the convention for naming message keys to avoid clashes.

For text in JSPs use the complete path + filename of the JSP, then a one-word name for the message. e.g. for the title of `jsp/myspace/main.jsp` use:

```
jsp.myspace.main.title
```

Some common words (e.g. "Help") can be brought out into keys starting `jsp.` for ease of translation, e.g.:

```
jsp.admin = Administer
```

Other common words/phrases are brought out into 'general' parameters if they relate to a set (directory) of JSPs, e.g.

```
jsp.tools.general.delete = Delete
```

Phrases that relate strongly to a topic (eg. `MyDSpace`) but used in many JSPs outside the particular directory are more convenient to be cross-referenced. For example one could use the key below in `jsp/submit/saved.jsp` to provide a link back to the user's `MyDSpace`:

(Cross-referencing of keys in general is not a good idea as it may make maintenance more difficult. But in some cases it has more advantages as the meaning is obvious.)

```
jsp.myspace.general.goto-myspace = Go to My DSpace
```

For text in servlet code, in custom JSP tags or wherever applicable use the fully qualified classname + a one-word name for the message. e.g.

```
org.dspace.app.webui.jsptag.ItemListTag.title = Title
```

**Which Languages are currently supported?** To view translations currently being developed, please refer to the [i18n](#) page of the `DSpace Wiki`.

**HTML Content in Items**

For the most part, the DSpace item display just gives a link that allows an end-user to download a bitstream. However, if a bundle has a primary bitstream whose format is of MIME type text/html, instead a link to the HTML servlet is given.

So if we had an HTML document like this:

```
contents.html
chapter1.html
chapter2.html
chapter3.html
figure1.gif
figure2.jpg
figure3.gif
figure4.jpg
figure5.gif
figure6.gif
```

The Bundle's primary bitstream field would point to the contents.html Bitstream, which we know is HTML (check the format MIME type) and so we know which to serve up first.

The HTML servlet employs a trick to serve up HTML documents without actually modifying the HTML or other files themselves. Say someone is looking at contents.html from the above example, the URL in their browser will look like this:

```
https://dspace.mit.edu/html/1721.1/12345/contents.html
```

If there's an image called figure1.gif in that HTML page, the browser will do HTTP GET on this URL:

```
https://dspace.mit.edu/html/1721.1/12345/figure1.gif
```

The HTML document servlet can work out which item the user is looking at, and then which Bitstream in it is called figure1.gif, and serve up that bitstream. Similar for following links to other HTML pages. Of course all the links and image references have to be relative and not absolute.

This can cope with relative links that refer to a deeper path, e.g.

```
<IMG SRC="images/figure1.gif">
```

Remember that in the Bitstream table in the database we have the 'name' field, which always contains the filename with no path (figure1.gif). We also have the source field, which may contain the full pathname of the file as it appeared on the submitter's hard drive, but this is browser- and OS-dependent, so we can't rely on it. All we can rely on is the filename.

We can still work out what images/figure1.gif is by making the HTML document servlet strip any path that comes in from the URL, e.g.

```
https://dspace.mit.edu/html/1721.1/12345/images/figure1.gif
```

Strip this

BUT all the filenames (regardless of directory names) must be unique. For example, this wouldn't work:

```
contents.html
chapter1.html
chapter2.html
chapter1_images/figure.gif
chapter2_images/figure.gif
```

since the HTML document servlet wouldn't know which bitstream to serve up for:

```
https://dspace.mit.edu/html/1721.1/12345/chapter1_images/figure.gif
https://dspace.mit.edu/html/1721.1/12345/chapter2_images/figure.gif
```

since it would just have figure.gif in the Bitstream table. Thus, the limitations are:

All links must be relative and not refer to parents (e.g. ../images/foo.gif or /images/foo.gif) If links refer to deeper directory levels, all the filenames must be different (as explained above)

### Thesis Blocking

The submission UI has an optional feature that came about as a result of MIT Libraries policy. If the block.theses parameter in dspace.cfg is true, an extra checkbox is included in the first page of the submission UI. This asks the user if the submission is a thesis. If the user checks this box, the submission is halted (deleted) and an error message displayed, explaining that DSpace should not be used to submit theses. This feature can be turned off and on, and the message displayed (/dspace/jsp/submit/no-theses.jsp can be localized as necessary).

## 7.4.2 OAI-PMH Data Provider

The DSpace platform supports the [Open Archives Initiative Protocol for Metadata Harvesting](#) (OAI-PMH) version 2.0 as a data provider. This is accomplished using the [OAICat](#) framework from OCLC.

The DSpace build process builds a Web application archive, [dspace-source]/build/dspace-oai.war), in much the same way as the Web UI build process described above. The only differences are that the JSPs are not included, and [dspace-source]/etc/oai-web.xml is used as the deployment descriptor. This 'webapp' is deployed to receive and respond to OAI-PMH requests via HTTP. Note that typically it should not be deployed on SSL (https: protocol). In a typical configuration, this is deployed at dspace-oai, for example:

```
http://dspace.myu.edu/dspace-oai/request?verb=Identify
```

The 'base URL' of this DSpace deployment would be:

```
http://dspace.myu.edu/dspace-oai/request
```

It is this URL that should be registered with [www.openarchives.org](http://www.openarchives.org). Note that you can easily change the 'request' portion of the URL by editing [dspace-source]/etc/oai-web.xml and rebuilding and deploying dspace-oai.war.

DSpace provides implementations of the OAICat interfaces AbstractCatalog, RecordFactory and Crosswalk that interface with the DSpace content management API and harvesting API (in the search subsystem).

Only the basic oai\_dc unqualified Dublin Core metadata set export is enabled by default; this is particularly easy since all items have qualified Dublin Core metadata. When this metadata is harvested, the qualifiers are simply stripped; for example, description.abstract is exposed as unqualified description. The description.provenance field is hidden, as this contains private information about the submitter and workflow reviewers of the item, including their e-mail addresses. Additionally, to keep in line with OAI community practices, values of contributor.author are exposed as creator values.

Other metadata formats are supported as well, using other Crosswalk implementations; consult the oaicat.properties file described below. To enable a format, simply uncomment the lines beginning with Crosswalks.\*. Multiple formats are allowed, and the current list includes, in addition to unqualified DC: MPEG DIDL, METS, MODS. There is also an incomplete, experimental qualified DC.

Note that the current simple DC implementation (org.dspace.app.oai.OAIDCCrosswalk) does not currently strip out any invalid XML characters that may be lying around in the data. If your database contains a DC value with, for example, some ASCII control codes (form feed etc.) this may cause OAI harvesters problems. This should rarely occur, however. XML entities (such as >) are encoded (e.g. to &gt;);



In addition to the implementations of the OAICat interfaces, there are two configuration files relevant to OAI support:

**oaiicat.properties** This resides as a template in `[dSPACE]/config/templates`, and the live version is written to `[dSPACE]/config`. You probably won't need to edit this; the `install-configs` script fills out the relevant deployment-specific parameters. You might want to change the `earliestDatestamp` field to accurately reflect the oldest datestamp in the system. (Note that this is the value of the `last_modified` column in the `Item` database table.)

**oai-web.xml** This standard Java Servlet 'deployment descriptor' is stored in the source as `[dSPACE-source]/etc/oai-web.xml`, and is written to `/dSPACE/oai/WEB-INF/web.xml`.

## Sets

OAI-PMH allows repositories to expose an hierarchy of sets in which records may be placed. A record can be in zero or more sets.

DSPACE exposes collections as sets. The organization of communities is likely to change over time, and is therefore a less stable basis for selective harvesting.

Each collection has a corresponding OAI set, discoverable by harvesters via the `ListSets` verb. The `setSpec` is the Handle of the collection, with the `'` and `'/'` converted to underscores so that the Handle is a legal `setSpec`, for example:

```
hdl_1721.1_1234
```

Naturally enough, the collection name is also the name of the corresponding set.

## Unique Identifier

Every item in OAI-PMH data repository must have a unique identifier, which must conform to the URI syntax. As of DSPACE 1.2, Handles are not used; this is because in OAI-PMH, the OAI identifier identifies the metadata record associated with the resource. The resource is the DSPACE item, whose resource identifier is the Handle. In practical terms, using the Handle for the OAI identifier may cause problems in the future if DSPACE instances share items with the same Handles; the OAI metadata record identifiers should be different as the different DSPACE instances would need to be harvested separately and may have different metadata for the item.

The OAI identifiers that DSPACE uses are of the form:

```
oai:host name:handle
```

For example:

```
oai:dSPACE.myu.edu:123456789/345
```

If you wish to use a different scheme, this can easily be changed by editing the value of `OAI_ID_PREFIX` at the top of the `org.dSPACE.app.oai.DSPACEOAICatalog` class. (You do not need to change the code if the above scheme works for you; the code picks up the host name and Handles automatically from the DSPACE configuration.)

## Access control

OAI provides no authentication/authorisation details, although these could be implemented using standard HTTP methods. It is assumed that all access will be anonymous for the time being.

A question is, "is all metadata public?" Presently the answer to this is yes; all metadata is exposed via OAI-PMH, even if the item has restricted access policies. The reasoning behind this is that people who do actually

have permission to read a restricted item should still be able to use OAI-based services to discover the content. If in the future, this 'expose all metadata' approach proves unsatisfactory for any reason, it should be possible to expose only publicly readable metadata. The authorisation system has separate permissions for READING and item and READING the content (bitstreams) within it. This means the system can differentiate between an item with public metadata and hidden content, and an item with hidden metadata as well as hidden content. In this case the OAI data repository should only expose items those with anonymous READ access, so it can hide the existence of records to the outside world completely. In this scenario, one should be wary of protected items that are made public after a time. When this happens, the items are "new" from the OAI-PMH perspective.

### Modification Date (OAI Date Stamp)

OAI-PMH harvesters need to know when a record has been created, changed or deleted. DSpace keeps track of a 'last modified' date for each item in the system, and this date is used for the OAI-PMH date stamp. This means that any changes to the metadata (e.g. admins correcting a field, or a withdrawal) will be exposed to harvesters.

### 'About' Information

As part of each record given out to a harvester, there is an optional, repeatable "about" section which can be filled out in any (XML-schema conformant) way. Common uses are for provenance and rights information, and there are schemas in use by OAI communities for this. Presently DSpace does not provide any of this information.

### Deletions

DSpace keeps track of deletions (withdrawals). These are exposed via OAI, which has a specific mechanism for dealing with this. Since DSpace keeps a permanent record of withdrawn items, in the OAI-PMH sense DSpace supports deletions 'persistently'. This is as opposed to 'transient' deletion support, which would mean that deleted records are forgotten after a time.

Once an item has been withdrawn, OAI-PMH harvests of the date range in which the withdrawal occurred will find the 'deleted' record header. Harvests of a date range prior to the withdrawal will not find the record, despite the fact that the record did exist at that time.

As an example of this, consider an item that was created on 2002-05-02 and withdrawn on 2002-10-06. A request to harvest the month 2002-10 will yield the 'record deleted' header. However, a harvest of the month 2002-05 will not yield the original record.

Note that presently, the deletion of 'expunged' items is not exposed through OAI.

### Flow Control (Resumption Tokens)

An OAI data provider can prevent any performance impact caused by harvesting by forcing a harvester to receive data in time-separated chunks. If the data provider receives a request for a lot of data, it can send part of the data with a resumption token. The harvester can then return later with the resumption token and continue.

DSpace supports resumption tokens for 'ListRecords' OAI-PMH requests. ListIdentifiers and ListSets requests do not produce a particularly high load on the system, so resumption tokens are not used for those requests.

Each OAI-PMH ListRecords request will return at most 100 records. This limit is set at the top of `org.dspace.app.oai.DSpaceOAIcatalog.java` (`MAX_RECORDS`). A potential issue here is that if a harvest yields an exact multiple of `MAX_RECORDS`, the last operation will result in a harvest with no records in it. It is unclear from the OAI-PMH specification if this is acceptable.

When a resumption token is issued, the optional `completeListSize` and `cursor` attributes are not included.

OAICat sets the `expirationDate` of the resumption token to one hour after it was issued, though in fact since DSpace resumption tokens contain all the information required to continue a request they do not actually expire. Resumption tokens contain all the state information required to continue a request. The format is:

```
from/until/setSpec/offset
```

`from` and `until` are the ISO 8601 dates passed in as part of the original request, and `setSpec` is also taken from the original request. `offset` is the number of records that have already been sent to the harvester. For example:

```
2003-01-01//hdl_1721_1_1234/300
```

This means the harvest is 'from' 2003-01-01, has no 'until' date, is for collection `hdl:1721.1/1234`, and 300 records have already been sent to the harvester. (Actually, if the original OAI-PMH request doesn't specify a 'from' or 'until', OAICat fills them out automatically to '0000-00-00T00:00:00Z' and '9999-12-31T23:59:59Z' respectively. This means DSpace resumption tokens will always have from and until dates in them.)

### 7.4.3 Package Importer and Exporter

This command-line tool gives you access to the Packager plugins. It can ingest a package to create a new DSpace Item, or disseminate an Item as a package.

To see all the options, invoke it as:

```
[dspace]/bin/dsrun org.dspace.app.packager.Packager --help
```

This mode also displays a list of the names of package ingesters and disseminators that are available.

#### Ingesting

To ingest a package from a file, give the command:

```
[dspace]/bin/dsrun org.dspace.app.packager.Packager -e user \  
-c handle \  
-t packager path
```

Where `user` is the e-mail address of the E-Person under whose authority this runs; `handle` is the Handle of the collection into which the Item is added, `packager` is the plugin name of the package ingester to use, and `path` is the path to the file to ingest (or "-" to read from the standard input).

Here is an example that loads a PDF file with internal metadata as a package:

```
/dspace/bin/dsrun org.dspace.app.packager.Packager -e florey@mit.edu \  
-c 1721.2/13 \  
-t pdf thesis.pdf
```

This example takes the result of retrieving a URL and ingests it:

```
wget -O - http://alum.mit.edu/jarandom/my-thesis.pdf | \  
/dspace/bin/dsrun org.dspace.app.packager.Packager -e florey@mit.edu -c 1721.2/13 -t pdf -
```

#### Disseminating

To disseminate an Item as a package, give the command:

```
[dspace]/bin/dsrun org.dspace.app.packager.Packager -e user \  
-d -i handle  
-t packager path
```

Where user is the e-mail address of the E-Person under whose authority this runs; handle is the Handle of the Item to disseminate; packager is the plugin name of the package disseminator to use; and path is the path to the file to create (or "-" to write to the standard output). This example writes an Item out as a METS package in the file "454.zip":

```
/dspace/bin/dsrun org.dspace.app.packager.Packager -e florey@mit.edu \
                                                    -d -i 1721.2/454 \
                                                    -t METS 454.zip
```

#### 7.4.4 Item Importer and Exporter

DSpace has a set of command line tools for importing and exporting items in batches, using the DSpace simple archive format. The tools are not terribly robust, but are useful and are easily modified. They also give a good demonstration of how to implement your own item importer if desired.

**Warning: templates may be applied**

Due to a bug as of 1.2 beta 2, if you have an Item template in your Collection, then those default values may be added to Items that you import. Be sure to remove the template if this is unwanted behavior.

#### DSpace simple archive format

The basic concept behind the DSpace's simple archive format is to create an archive, which is directory full of items, with a subdirectory per item. Each item directory contains a file for the item's descriptive metadata, and the files that make up the item.

```
archive_directory/
  item_000/
    dublin_core.xml -- qualified Dublin Core metadata
    contents        -- text file containing one line per filename
    file_1.doc      -- files to be added as bitstreams to the item
    file_2.pdf
  item_001/
    dublin_core.xml
    contents
    file_1.png
    ...
```

The dublin\_core.xml file has the following format, where each Dublin Core element has its own entry within a <dcvalue> tagset. There are currently three tag elements available in the <dcvalue> tagset:

- <element> - the Dublin Core element
- <qualifier> - the element's qualifier
- <language> - (optional)ISO language code for element

```
<dublin_core>
  <dcvalue element="title" qualifier="none">A Tale of Two Cities</dcvalue>
  <dcvalue element="date" qualifier="issued">1990</dcvalue></dublin_core>
  <dcvalue element="title" qualifier="alternate" language="fr" ">J'aime les Printemps</dcvalue>
</dublin_core>
```

(Note the optional language tag which notifies the system that the optional title is in French.)

The contents file simply enumerates, one file per line, the bitstream file names. The bitstream name may optionally be followed by the sequence:

```
\tbundle:bundlename
```

where `\t` is the tab character and 'bundlename' is replaced by the name of the bundle to which the bitstream should be added. If no bundle is specified, the bitstream will be added to the 'ORIGINAL' bundle.

### Importing Items

Note: Before running the item importer over items previously exported from a DSpace instance, please first refer to Transferring Items Between DSpace Instances.

The item importer is in `org.dspace.app.itemimport.ItemImport`, and is run with the `dstrun` utility in the `dspace/bin` directory. Running it with `-h` gets the current command-line arguments. Another very important flag is the `-test` flag, which you can use with any command to simulate all of the actions it will perform without actually making any changes to your DSpace instance - very useful for validating your item directories before doing an import. In the importer's arguments you can use either the user's database ID or email address and the `eperson` ID, and the collection's database ID or handle as arguments. Currently with the importer you can add, remove, and replace items in a collection. If you specify more than one collection argument then the items will be imported to multiple collections, and the first collection specified becomes the "owning" collection. If there is an error and the import is aborted, there is a `-resume` flag that you can try to resume the import where you left off after you fix the error.

To add items to a collection with an EPerson as the submitter, type:

```
dstrun org.dspace.app.itemimport.ItemImport --add \  
      --eperson=joe@user.com \  
      --collection=collectionID \  
      --source=items_dir \  
      --mapfile=mapfile
```

(or by using the short form)

```
dstrun org.dspace.app.itemimport.ItemImport -a \  
      -e joe@user.com \  
      -c collectionID \  
      -s items_dir \  
      -m mapfile
```

which would then cycle through the archive directory's items, import them, and then generate a map file which stores the mapping of item directories to item handles. Save this map file! Using the map file you can then 'unimport' with the command:

```
dstrun org.dspace.app.itemimport.ItemImport --delete \  
      --mapfile=mapfile
```

The imported items listed in the map file would then be deleted. If you wish to replace previously imported items, you can give the command:

```
dstrun org.dspace.app.itemimport.ItemImport --replace \  
      --eperson=joe@user.com \  
      --collection=collectID \  
      --source=items_dir \  
      --mapfile=mapfile
```

Replacing items uses the map file to replace the old items and still retain their handles.

The importer usually bypasses any workflow assigned to a collection, but adding the `--workflow` option will route the imported items through the workflow system.

The importer also has a `--test` flag that will simulate the entire import process without actually doing the import. This is extremely useful for verifying your import files before doing the import step.

### Exporting Items

The item exporter can export a single item or a collection of items, and creates a DSpace simple archive for each item to be exported. To export a collection's items you type:

```
dsrun org.dspace.app.itemexport.ItemExport --type=COLLECTION \
      --id=collID \
      --dest=dest_dir \
      --number=seq_num
```

The keyword `COLLECTION` means that you intend to export an entire collection. The ID can either be the database ID or the handle. The exporter will begin numbering the simple archives with the sequence number that you supply. To export a single item use the keyword `ITEM` and give the item ID as an argument:

```
dsrun org.dspace.app.itemexport.ItemExport --type=ITEM \
      --id=itemID \
      --dest=dest_dir \
      --number=seq_num
```

Each exported item will have an additional file in its directory, named 'handle'. This will contain the handle that was assigned to the item, and this file will be read by the importer so that items exported and then imported to another machine will retain the item's original handle.

### 7.4.5 Transferring Items Between DSpace Instances

Where items are to be moved between DSpace instances (for example from a test DSpace into a production DSpace) the item exporter and item importer can be used in conjunction with a script to assist in this process. After running the item exporter each `dublin_core.xml` file will contain metadata that was automatically added by DSpace. These fields are as follows:

- `date.accessioned`
- `date.available`
- `date.issued`
- `description.provenance`
- `format.extent`
- `format.mimetype`
- `identifier.uri`

In order to avoid duplication of this metadata, run

```
dspace_migrate <exported item directory>
```

prior to running the item importer. This will remove the above metadata items, except for `date.issued` - if the item has been published or publicly distributed before and `identifier.uri` - if it is not the handle, from the `dublin_core.xml` file and remove all handle files. It will then be safe to run the item exporter. Use

```
dspace_migrate --help
```

for instructions on use of the script.

### 7.4.6 Registering (Not Importing) Bitstreams

Registration is an alternate means of incorporating items, their metadata, and their bitstreams into DSpace by taking advantage of the bitstreams already being in storage accessible to DSpace. An example might be that there is a repository for existing digital assets. Rather than using the normal interactive ingest process or the batch import to furnish DSpace the metadata and to upload bitstreams, registration provides DSpace the metadata and the location of the bitstreams. DSpace uses a variation of the import tool to accomplish registration.

#### Accessible Storage

To register an item its bitstreams must reside on storage accessible to DSpace and therefore referenced by an asset store number in `dspace.cfg`. The configuration file `dspace.cfg` establishes one or more asset stores through the use of an integer asset store number. This number relates to a directory in the DSpace host's file system or a set of SRB account parameters. This asset store number is described in The [dspace.cfg Configuration Properties File](#) section and in the `dspace.cfg` file itself. The asset store number(s) used for registered items should generally not be the value of the `assetstore.incoming` property since it is unlikely that that you will want to mix the bitstreams of normally ingested and imported items and registered items.

#### Registering Items Using the Item Importer

DSpace uses the same import tool that is used for batch import except that several variations are employed to support registration. The discussion that follows assumes familiarity with the import tool.

The archive format for registration does not include the actual content files (bitstreams) being registered. The format is however a directory full of items to be registered, with a subdirectory per item. Each item directory contains a file for the item's descriptive metadata (`dublin_core.xml`) and a file listing the item's content files (`contents`), but not the actual content files themselves.

The `dublin_core.xml` file for item registration is exactly the same as for regular item import.

The `contents` file, like that for regular item import, lists the item's content files, one content file per line, but each line has the one of the following formats:

```
-r -s n -f filepath
-r -s n -f filepath\tbundle:bundlename
```

where

- `-r` indicates this is a file to be registered
- `-s n` indicates the asset store number (`n`)
- `-f filepath` indicates the path and name of the content file to be registered (`filepath`)
- `\t` is a tab character
- `bundle:bundlename` is an optional bundle name

The bundle, that is everything after the filepath, is optional and is normally not used. The command line for registration is just like the one for regular import:

```
dsrun org.dspace.app.itemimport.ItemImport --add \
      --eperson=joe@user.com \
      --collection=collectionID \
      --source=items_dir \
      --mapfile=mapfile
```

(or by using the short form)

```
dsrun org.dspace.app.itemimport.ItemImport -a \
      -e joe@user.com \
      -c collectionID \
      -s items_dir \
      -m mapfile
```

The `-workflow` and `-test` flags will function as described in [Importing Items](#).

The `-delete` flag will function as described in [Importing Items](#) but the registered content files will not be removed from storage. See [Deleting Registered Items](#).

The `-replace` flag will function as described in [Importing Items](#) but care should be taken to consider different cases and implications. With old items and new items being registered or ingested normally, there are four combinations or cases to consider. Foremost, an old registered item deleted from DSpace using `-replace` will not be removed from the storage. See [Deleting Registered Items](#). where it resides. A new item added to DSpace using `-replace` will be ingested normally or will be registered depending on whether or not it is marked in the contents files with the `-r`.

### Internal Identification and Retrieval of Registered Items

Once an item has been registered, superficially it is indistinguishable from items ingested interactively or by batch import. But internally there are some differences:

First, the randomly generated internal ID is not used because DSpace does not control the file path and name of the bitstream. Instead, the file path and name are that specified in the contents file.

Second, the `store_number` column of the bitstream database row contains the asset store number specified in the contents file.

Third, the `internal_id` column of the bitstream database row contains a leading flag (`-R`) followed by the registered file path and name. For example, `-Rfilepath` where `filepath` is the file path and name relative to the asset store corresponding to the asset store number. The asset store could be traditional storage in the DSpace server's file system or an SRB account.

Fourth, an MD5 checksum is calculated by reading the registered file if it is in local storage. If the registered file is in remote storage (say, SRB) a checksum is calculated on just the file name! This is an efficiency choice since registering a large number of large files that are in SRB would consume substantial network resources and time. A future option could be to have an SRB proxy process calculate MD5s and store them in SRB's metadata catalog (MCAT) for rapid retrieval. SRB offers such an option but it's not yet in production release.

Registered items and their bitstreams can be retrieved transparently just like normally ingested items.

### Exporting Registered Items

Registered items may be exported as described in [Exporting Items](#). If so, the export directory will contain actual copies of the files being exported but the lines in the contents file will flag the files as registered. This means that if DSpace items are "round tripped" (see [Transferring Items Between DSpace Instances](#)) using the



exporter and importer, the registered files in the export directory will again be registered in DSpace instead of being uploaded and ingested normally.

### METS Export of Registered Items

The METS **Export Tool** can also be used but note the cautions described in that section and note that MD5 values for items in remote storage are actually MD5 values on just the file name.

### Deleting Registered Items

If a registered item is deleted from DSpace, either interactively or by using the `-delete` or `-replace` flags described in Importing Items, the item will disappear from DSpace but its registered content files will remain in place just as they were prior to registration. Bitstreams not registered but added by DSpace as part of registration, such as `license.txt` files, will be deleted.

## 7.4.7 METS Tools

The experimental (incomplete) METS export tool writes DSpace items to a filesystem with the metadata held in a more standard format based on METS.

### The Export Tool

The METS export tool is invoked via the command line like this:

```
[dspace]/bin/dsrun org.dspace.app.mets.METSExport --help
```

The tool can export an individual item, the items within a given collection, or everything in the DSpace instance. To export an individual item, use:

```
[dspace]/bin/dsrun org.dspace.app.mets.METSExport --item [handle]
```

To export the items in collection `hdl:123.456/789`, use:

```
[dspace]/bin/dsrun org.dspace.app.mets.METSExport --collection hdl:123.456/789
```

To export all the items in DSpace, use:

```
[dspace]/bin/dsrun org.dspace.app.mets.METSExport --all
```

With any of the above forms, you can specify the base directory into which the items will be exported, using `-destination [directory]`. If this parameter is omitted, the current directory is used.

### The AIP Format

Each exported item is written to a separate directory, created under the base directory specified in the command-line arguments, or in the current directory if `-destination` is omitted. The name of each directory is the Handle, URL-encoded so that the directory name is 'legal'.

Within each item directory is a `mets.xml` file which contains the METS-encoded metadata for the item. Bitstreams in the item are also stored in the directory. Their filenames are their MD5 checksums, firstly for easy integrity checking, and also to avoid any problems with 'special characters' in the filenames that were legal on the original filing system they came from but are illegal in the server filing system. The `mets.xml` file includes XLink pointers to these bitstream files.

An example AIP might look like this:

- hdl%3A123456789%2F8/
  - mets.xml – METS metadata
  - 184BE84F293342 – bitstream
  - 3F9AD0389CB821
  - 135FB82113C32D

The contents of the METS in the mets.xml file are as follows:

- A dmdSec (descriptive metadata section) containing the item’s metadata in **Metadata Object Description Schema (MODS)** XML. The Dublin Core descriptive metadata is mapped to MODS since there is no official qualified Dublin Core XML schema in existence as of yet, and the Library Application Profile of DC that DSpace uses includes some qualifiers that are not part of the **DCMI Metadata Terms**.
- An amdSec (administrative metadata section), which contains the a rights metadata element, which in turn contains the base64-encoded deposit license (the license the submitter granted as part of the submission process).
- A fileSec containing a list of the bitstreams in the item. Each bundle constitutes a fileGrp. Each bitstream is represented by a file element, which contains an FLocat element with a simple XLink to the bitstream in the same directory as the mets.xml file. The file attributes consist of most of the basic technical metadata for the bitstream. Additionally, for those bitstreams that are thumbnails or text extracted from another bitstream in the item, those ‘derived’ bitstreams have the same GROUPID as the bitstream they were derived from, in order that clients understand that there is a relationship. The OWNERID of each file is the ‘persistent’ **bitstream identifier** assigned by the DSpace instance. The ID and GROUPID attributes consist of the item’s Handle, together with the bitstream’s sequence ID, which underscores used in place of dots and slashes. For example, a bitstream with sequence ID 24, in the item hdl:123.456/789 will have the ID 123\_456\_789\_24. This is because ID and GROUPID attributes must be of type xsd:id.

### Limitations

- No corresponding import tool yet
- No structmap section
- Some technical metadata not written, e.g. the primary bitstream in a bundle, original filenames or descriptions.
- Only the MIME type is stored, not the (finer grained) bitstream format.
- Dublin Core to MODS mapping is very simple, probably needs verification

### 7.4.8 Media Filters

DSpace can apply filters to content/bitstreams, creating new content. Filters are included that extract text for full-text searching, and create thumbnails for items that contain images. The media filters are controlled by the MediaFilterManager which traverses the asset store, invoking the MediaFilter subclasses on bitstreams. The MediaFilter plugin config item plugin.named.org.dspace.app.mediafilter.MediaFilter in dspace.cfg contains a list of bitstream format types and the filters that operate on bitstreams of that type. The media filter system is intended to be run from the command line (or regularly as a cron task):

```
dspace/bin/filter-media
```

Traverse the asset store, applying media filters to bitstreams, skipping bitstreams that have already been filtered.

```
dspace/bin/filter-media -f
```

Apply filters to ALL bitstreams, even if they've already been filtered.

```
dspace/bin/filter-media -v
```

Verbose mode - print all extracted text and other filter details to STDOUT.

```
dspace/bin/filter-media -n
```

Suppress index creation - by default, a new search index is created for full-text searching. This option suppresses index creation if you intend to run index-all elsewhere.

```
dspace/bin/filter-media -i 123456789/2
```

Restrict processing to the community, collection, or item named by the identifier - by default, all bitstreams of all items in the repository are processed. The identifier must be a handle, not a DB key. This option may be combined with any other option.

```
dspace/bin/filter-media -m 1000
```

Suspend operation after the specified maximum number of items have been processed - by default, no limit exists. This option may be combined with any other option.

Adding your own filters is done by creating a sub-class of the MediaFilter class. See the comments in the source file MediaFilter.java for more information. In theory filters could be implemented in any language (C, Perl, etc.) They only need to be invoked by the Java code in the MediaFilter class that you create.

### 7.4.9 Sub-Community Management

DSpace provides an administrative tool - 'CommunityFiliator' - for managing community sub-structure. Normally this structure seldom changes, but prior to the 1.2 release sub-communities were not supported, so this tool could be used to place existing pre-1.2 communities into a hierarchy. It has two operations, either establishing a community to sub-community relationship, or dis-establishing an existing relationship.

The familiar parent/child metaphor can be used to explain how it works. Every community in DSpace can be either a 'parent' community - meaning it has at least one sub-community, or a 'child' community - meaning it is a sub-community of another community, or both or neither. In these terms, an 'orphan' is a community that lacks a parent (although it can be a parent); 'orphans' are referred to as 'top-level' communities in the DSpace user-interface, since there is no parent community 'above' them. The first operation - establishing a parent/child relationship - can take place between any community and an orphan. The second operation - removing a parent/child relationship - will make the child an orphan.

Using the dsrun utility in the dspace/bin directory, the establish operation looks like this:

```
dsrun org.dspace.administer.CommunityFiliator --set --parent=parentID --child=childID
```

(or using the short form)

```
dsrun org.dspace.administer.CommunityFiliator -s -p parentID -c childID
```

where '-s' or '--set' means establish a relationship whereby the community identified by the '-p' parameter becomes the parent of the community identified by the '-c' parameter. Both the 'parentID' and 'childID' values may be handles or database IDs.

The reverse operation looks like this:

```
dsrun org.dspace.administer.CommunityFiliator --remove --parent=parentID --child=childID
```

(or using the short form)

```
dsrun org.dspace.administer.CommunityFiliator -r -p parentID -c childID
```

where '-r' or '--remove' means dis-establish the current relationship in which the community identified by 'parentID' is the parent of the community identified by 'childID'. The outcome will be that the 'childID' community will become an orphan, i.e. a top-level community.

If the required constraints of operation are violated, an error message will appear explaining the problem, and no change will be made. An example in a removal operation, where the stated child community does not have the stated parent community as its parent: "Error, child community not a child of parent community".

It is possible to effect arbitrary changes to the community hierarchy by chaining the basic operations together. For example, to move a child community from one parent to another, simply perform a 'remove' from its current parent (which will leave it an orphan), followed by a 'set' to its new parent.

It is important to understand that when any operation is performed, all the sub-structure of the child community follows it. Thus, if a child has itself children (sub-communities), or collections, they will all move with it to its new 'location' in the community tree.

# Chapter 8

## Version History

### 8.1 Changes in DSpace 1.4.1

#### 8.1.1 General Improvements

- HandleServlet and BitstreamServlet support If-Modified-Since requests
- Improved sanity-checking of XSL-based ingest crosswalks
- Remove thumbnail filename from alt-text
- Include item title in HTML title element
- Improvements to help prevent spammers and sploggers
- Make cleanup() commit outstanding work every 100 iterations
- Better handling where email send failed due to wrong address for new user
- Include robots.txt to limit bots navigating author, date and browse by subject pages
- Add css styles for print media
- RSS made more configurable and provide system-wide RSS feed, also moves text to Messages.properties
- Jar file updates (includes required code changes for DSIndexer and DSQuery and new jars fontbox.jar and serializer.jar)
- Various documentation additions and cleanups
- XHTML compliance improvements

#### 8.1.2 Bug fixes

- 1532389 - Item Templates do not work for non-dc fields
- 1066771 - Metadata edit form dropping DC qualifier
- 1548738 - Multiple Metadata Schema, schema not shown on edit item page
- 1589895 - Not possible to add unqualified Metadata Field

- 1543853 - Statistics do not work in 1.4
- 1541381 - Browse-by-date and browse-by-title not working
- 1556947 - NullPointerException when no user selected to del/edit
- 1554064 - Fix exception handling for ClassCastException in BitstreamServlet
- 1548865 - Browse errors on withdrawn item
- 1554056 - Community/collection handle URL with / redirects to homepage
- 1571490 - UTF-8 encoded characters in licence
- 1571519 - UTF-8 in statistics
- 1544807 - Browse-by-Subject/Author paging mechanism broken
- 1543966 - "Special" groups inside groups bug
- 1480496 - Cannot turn off "ignore authorization" flag!
- 1515148 - Community policies not deleting correctly
- 1556829 - Docs mention old SiteAuthenticator class
- Fix for bitstream authorization timeout
- Fix to make sure cleanup() doesn't fail with NullPointerException
- Fix for removeBitstream() failing to update primary bitstream
- Fix for Advanced Search ignoring conjunctions for arbitrary number of queries
- Fix minor bug in Harvest.java for Oracle users
- Fix missing title for news editor page
- Small Messages.properties modification (change of DSpace copyright text)
- fix PDFBox tmp file issue
- Fix HttpServletRequest encoding issues

## 8.2 Changes in DSpace 1.4

### 8.2.1 General Improvements

- Content verification through periodic checksum checking
- Support for branded preview image
- Add/replace Creative Commons in 'edit item' tool
- Customisable item listing columns and browse indices
- Script for updating handle prefixes (e.g. for moving from development to production)

- Configurable boolean search operator
- Controlled vocabulary patch to provide search on classification terms, and addition of terms during submission.
- Add 'visibility' element to input-forms.xml
- Browse by subject feature
- Log4J enhancement to use XML configuration
- QueryArgs class can support any number of fields in advanced search.
- Community names no longer have to be unique
- Enhanced Windows support
- Support for multiple (flat) metadata schemas
- Suggest an item page
- RSS Feeds
- Performance enhancements
- Stackable authentication methods
- Plug-in manager
- Pluggable SIP/DIP support and metadata crosswalks
- Nested groups of e-people
- Expose METS and MPEG-21 DIDL DIPs via OAI-PMH
- Configurable Lucene search analyzer (e.g. for Chinese metadata)
- Support for SMTP servers requiring authentication

### 8.2.2 Bug fixes

- 1358197 - Edit Item, empty DC fields not removable
- 1363633 - Submission step 1 fails when there are no collections
- 1255264 - Resource policy eperson value was set to wrong column
- 1380494 - Error deleting an item with multiple metadata schema support
- 1443649 - Cannot configure unqualified elements for advanced search index
- 1333687 - Browse-(title|date) fails on withdrawn item
- 1066713 - Two (sub)communities cannot have one name
- 1284055 - Two Communities of same name throws error
- 1035366 - Bitstream size column should be bigint

- 1352257 - Selecting a Group for GroupToGroup while Creating Collection
- 1352226 - Navigation and Sorting in Group List (Select Groups) fails
- 1348276 - Null in collection name causes OAI ListSets to fail
- 1160898 - dspace\_migrate removes Date.Issued from prev published items
- 1261191 - Malformed METS metadata exported

## 8.3 Changes in DSpace 1.3.2

### 8.3.1 General Improvements

- DSpace UI XHTML/WAI compliant
- Configure metadata fields shown on simple item display
- Supervisor/workspace help documentation

### 8.3.2 Bug fixes

- Oracle compatibility fixes
- Item exporter now correctly exports metadata in UTF-8
- fixed to handle 'null' values passed in

## 8.4 Changes in DSpace 1.3.1

### 8.4.1 Bug fixes

- 1252153 - Error on fresh install

## 8.5 Changes in DSpace 1.3

### 8.5.1 General Improvements

- Initial i18n Support for JSPs - Note: the implementation of this feature required changes to almost all JSP pages
- LDAP authentication support
- Log file analysis and report generation
- Configurable item licence viewing
- Supervision order/collaborative workspace administrative tools
- Basic workspace for submissions in progress, with support for supervision
- SRB storage system option
- Updated handle server system



- Database optimisations
- Latest versions of Xerces, Xalan and OAICAT jars
- Various documentation additions and cleanups

### 8.5.2 Bug fixes

- 1161459 - ItemExporter fails with Too many open files
- 1167373 - Email date field not populated
- 1193948 - New item submit problem
- 1188132 - NullPointerException when Adding EPerson
- 1188016 - Cannot Edit an Eperson
- 1219701 - Unable to open unfinished submission
- 1206836 - community strengths not reflecting sub-community
- 1238262 - Submit UI nav/progress buttons no longer show progress
- 1238276 - Double quote problem in some fields in submit UI
- 1238277 - format support level not shown in "uploaded file" page
- 1242548 - Uploading non-existing files
- 1244743 - Bad lookup key for special case of DC Title in ItemTag.java
- 1245223 - Subscription Emailer fails
- 1247508 - Error when browsing item with no content/bitstream collections
- Set the content type in the HTTP header
- Fix issue where EPerson edit would not work due to form indexing (partial fix)
- POST handling in HTMLServlet
- Missing ContentType directives added to some JSPs
- Name dependency on Collection Admin and Submitter groups fixed
- Fixed OAI-PMH XML encoding

## 8.6 Changes in DSpace 1.2.2

### 8.6.1 General Improvements

- Customisable submission forms added
- Configurable number of index terms in Lucene for full-text indexing
- Improved scalability in media filter
- Submit button on collection pages only appears if user has authorisation
- PostgreSQL 8.0 compatibility
- Search scope retention to improve browsing
- Community and collection strengths displayed
- Upgraded OAICat software

### 8.6.2 Bug fixes

- Fix for Oracle too many cursors problem.
- Fix for UTF-8 encoded searches in advanced search.
- Fix for handling \ in bitstream names.
- Fix to prevent delete of "unknown" bitstream format
- Fix for ItemImport creating new handles for replaced items

### 8.6.3 Changes in JSPs

- collection-home.jsp changed
- community-home.jsp changed
- community-list.jsp changed
- home.jsp changed
- dspace-admin/list-formats.jsp changed
- dspace-admin/wizard-questions.jsp changed
- search/results.jsp changed
- submit/cancel.jsp changed
- submit/change-file-description.jsp changed
- submit/choose-file.jsp changed
- submit/complete.jsp changed
- submit/creative-commons.jsp changed

- submit/edit-metadata.jsp new
- submit/get-file-format.jsp changed
- submit/initial-questions.jsp changed
- submit/progressbar.jsp changed
- submit/review.jsp changed
- submit/select-collection.jsp changed
- submit/show-license.jsp changed
- submit/show-uploaded-file.jsp changed
- submit/upload-error.jsp changed
- submit/upload-file-list.jsp changed

## 8.7 Changes in DSpace 1.2.1

### 8.7.1 General Improvements

- Oracle support added
- Thumbnails in item view can now be switched off/on
- Browse and search thumbnail options
- Improved item importer
  - can now import to multiple collections
  - added `-test` flag to simulate an import, without actually making any changes
  - added `-resume` flag to try to resume the import in case the import is aborted
- Configurable fields for the search index
- Script for transferring items between DSpace instances
- Sun library JARs (JavaMail, Java Activation Framework and Servlet) now included in DSpace source code bundle

### 8.7.2 Bug fixes

- A logo to existing collection can now be added. Fixes SF bug #1065933
- The community logo can now be edited. Fixes SF bug #1035692
- MediaFilterManager doesn't 'touch' every item every time. Fixes SF bug #1015296
- Supported formats help page, set the format support level to "known" as default
- Fixed various database connection pool leaks

### 8.7.3 Changed JSPs

- collection-home changed
- community-home changed
- display-item changed
- dspace-admin/confirm-delete-collection moved to tools/ and changed
- dspace-admin/confirm-delete-community moved to tools/ and changed
- dspace-admin/edit-collection moved to tools/ and changed
- dspace-admin/edit-community moved to tools/ and changed
- dspace-admin/index changed
- dspace-admin/upload-logo changed
- dspace-admin/wizard-basicinfo changed
- dspace-admin/wizard-default-item changed
- dspace-admin/wizard-permissions changed
- dspace-admin/wizard-questions changed
- help/formats.html removed
- help/formats changed
- index changed
- layout/navbar-admin changed

## 8.8 Changes in DSpace 1.2

### 8.8.1 General Improvements

- Communities can now contain sub-communities
- Items may be included in more than one collection
- Full text extraction and searching for MS Word, PDF, HTML, text documents
- Thumbnails displayed in item view for items that contain images
- Configurable MediaFilter tool creates both extracted text and thumbnails
- Bitstream IDs are now persistent - generated from item's handle and a sequence number
- Creative Commons licenses can optionally be added to items during web submission process

## 8.9 Administration

- If you are logged in as administrator, you see admin buttons on item, collection, and community pages
- New collection administration wizard
- Can now administer collection's submitters from collection admin tool
- Delegated administration - new 'collection editor' role - edits item metadata, manages submitters list, edits collection metadata, links to items from other collections, and can withdraw items
- Admin UI moved from /admin to /dspace-admin to avoid conflict with Tomcat /admin JSPs
- New EPerson selector popup makes Group editing much easier
- 'News' section is now editable using admin UI (no more mucking with JSPs)

### 8.9.1 Import/Export/OAI

- New tool that exports DSpace content in AIPs that use METS XML for metadata (incomplete)
- OAI - sets are now collections, identified by Handles ('safe' with /, : converted to \_)
- OAI - contributor.author now mapped to oai\_dc:creator

### 8.9.2 Miscellaneous

- Build process streamlined with use of WAR files, symbolic links no longer used, friendlier to later versions of Tomcat
- MIT-specific aspects of UI removed to avoid confusion
- Item metadata now rendered to avoid interpreting as HTML (displays as entered)
- Forms now have no-cache directive to avoid trouble with browser 'back' button
- Bundles now have 'names' for more structure in item's content

### 8.9.3 JSP file changes between 1.1 and 1.2

This list generated with `cvs -Q rdiff -s -r dspace-1_1 dspace` and a sprinkling of perl.

- Changed: `dspace/jsp/collection-home.jsp`
- Changed: `dspace/jsp/community-home.jsp`
- Changed: `dspace/jsp/community-list.jsp`
- Changed: `dspace/jsp/display-item.jsp`
- Changed: `dspace/jsp/index.jsp`
- Changed: `dspace/jsp/home.jsp`
- Changed: `dspace/jsp/styles.css.jsp`
- Moved to `dspace-admin` and changed: `dspace/jsp/admin/authorize-advanced.jsp`

- Moved to dspace-admin and changed: `dspace/jsp/admin/authorize-collection-edit.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/authorize-community-edit.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/authorize-item-edit.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/authorize-main.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/authorize-policy-edit.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/collection-select.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/community-select.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/confirm-delete-collection.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/confirm-delete-community.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/confirm-delete-dctype.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/confirm-delete-eperson.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/confirm-delete-format.jsp`
- Moved to dspace/jsp/tools: `dspace/jsp/admin/confirm-delete-item.jsp`
- Moved to dspace/jsp/tools: `dspace/jsp/admin/confirm-withdraw-item.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/edit-collection.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/edit-community.jsp`
- Moved to dspace/jsp/tools and changed: `dspace/jsp/admin/edit-item-form.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/eperson-browse.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/eperson-confirm-delete.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/eperson-edit.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/eperson-main.jsp`
- Moved to dspace/jsp/tools and changed: `dspace/jsp/admin/get-item-id.jsp`
- Moved to dspace/jsp/tools and changed: `dspace/jsp/admin/group-edit.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/group-eperson-select.jsp`
- Moved to dspace/jsp/tools and changed: `dspace/jsp/admin/group-list.jsp`
- Moved to dspace-admin: `dspace/jsp/admin/index.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/item-select.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/list-communities.jsp`
- Moved to dspace-admin and changed: `dspace/jsp/admin/list-dc-types.jsp`
- Removed: `dspace/jsp/admin/list-epeople.jsp`

- Moved to dspace-admin and changed: dspace/jsp/admin/list-formats.jsp
- Moved to dspace/jsp/tools: dspace/jsp/admin/upload-bitstream.jsp
- Moved to dspace-admin and changed: dspace/jsp/admin/upload-logo.jsp
- Moved to dspace-admin: dspace/jsp/admin/workflow-abort-confirm.jsp
- Moved to dspace-admin and changed: dspace/jsp/admin/workflow-list.jsp
- Changed: dspace/jsp/browse/authors.jsp
- Changed: dspace/jsp/browse/items-by-author.jsp
- Changed: dspace/jsp/browse/items-by-date.jsp
- Changed: dspace/jsp/browse/no-results.jsp
- New: dspace-admin/eperson-deletion-error.jsp
- New: dspace/jsp/dspace-admin/news-edit.jsp
- New: dspace/jsp/dspace-admin/news-main.jsp
- New: dspace/jsp/dspace-admin/wizard-basicinfo.jsp
- New: dspace/jsp/dspace-admin/wizard-default-item.jsp
- New: dspace/jsp/dspace-admin/wizard-permissions.jsp
- New: dspace/jsp/dspace-admin/wizard-questions.jsp
- Changed: dspace/jsp/components/contact-info.jsp
- Changed: dspace/jsp/error/internal.jsp
- New: dspace/jsp/help/formats.jsp
- Changed: dspace/jsp/layout/footer-default.jsp
- Changed: dspace/jsp/layout/header-default.jsp
- Changed: dspace/jsp/layout/navbar-admin.jsp
- Changed: dspace/jsp/layout/navbar-default.jsp
- Changed: dspace/jsp/login/password.jsp
- Changed: dspace/jsp/mydspace/main.jsp
- Changed: dspace/jsp/mydspace/perform-task.jsp
- Changed: dspace/jsp/mydspace/preview-task.jsp
- Changed: dspace/jsp/mydspace/reject-reason.jsp
- Changed: dspace/jsp/mydspace/remove-item.jsp
- Changed: dspace/jsp/register/edit-profile.jsp

- Changed: `dspace/jsp/register/inactive-account.jsp`
- Changed: `dspace/jsp/register/new-password.jsp`
- Changed: `dspace/jsp/register/registration-form.jsp`
- Changed: `dspace/jsp/search/advanced.jsp`
- Changed: `dspace/jsp/search/results.jsp`
- Changed: `dspace/jsp/submit/cancel.jsp`
- New: `dspace/jsp/submit/cc-license.jsp`
- Changed: `dspace/jsp/submit/choose-file.jsp`
- New: `dspace/jsp/submit/creative-commons.css`
- New: `dspace/jsp/submit/creative-commons.jsp`
- Changed: `dspace/jsp/submit/edit-metadata-1.jsp`
- Changed: `dspace/jsp/submit/edit-metadata-2.jsp`
- Changed: `dspace/jsp/submit/get-file-format.jsp`
- Changed: `dspace/jsp/submit/initial-questions.jsp`
- Changed: `dspace/jsp/submit/progressbar.jsp`
- Changed: `dspace/jsp/submit/review.jsp`
- Changed: `dspace/jsp/submit/select-collection.jsp`
- Changed: `dspace/jsp/submit/show-license.jsp`
- Changed: `dspace/jsp/submit/show-uploaded-file.jsp`
- Changed: `dspace/jsp/submit/upload-error.jsp`
- Changed: `dspace/jsp/submit/upload-file-list.jsp`
- Changed: `dspace/jsp/submit/verify-prune.jsp`
- New: `dspace/jsp/tools/edit-item-form.jsp`
- New: `dspace/jsp/tools/eperson-list.jsp`
- New: `dspace/jsp/tools/itemmap-browse.jsp`
- New: `dspace/jsp/tools/itemmap-info.jsp`
- New: `dspace/jsp/tools/itemmap-main.jsp`



## 8.10 Changes in DSpace 1.1.1

### 8.10.1 Bug fixes

- non-administrators can now submit again
- installations now preserve file creation dates, eliminating confusion with upgrades
- authorization editing pages no longer create null entries in database, and no longer handles them poorly (no longer gives blank page instead of displaying policies.)
- registration page Invalid token error page now displayed when an invalid token is received (as opposed to internal server error.) Fixes SF bug #739999
- eperson admin 'recent submission' links fixed for DSpaces deployed somewhere other than at / (e.g. /dspace).
- help pages Link to help pages now includes servlet context (e.g. '/dspace'). Fixes SF bug #738399.

### 8.10.2 Improvements

- bin/dspace-info.pl now checks jsp and asset store files for zero-length files
- make-release-package now works with SourceForge CVS
- eperson editor now doesn't display the spurious text 'null'
- item exporter now uses Jakarta's cli command line arg parser (much cleaner)
- item importer improvements:
  - now uses Jakarta's cli command line arg parser (much cleaner)
  - imported items can now be routed through a workflow
  - more validation and error messages before import
  - can now use email addresses and handles instead of just database IDs
  - can import an item to a collection with the workflow suppressed

## 8.11 Changes in DSpace 1.1

- Fixed various OAI-related bugs; DSpace's OAI support should now be correct. Note that harvesting is now based on the new Item 'last modified' date (as opposed to the Dublin Core date.available date.)
- Fixed Handle support—DSpace now responds to naming authority requests correctly.
- Multiple bitstream stores can now be specified; this allows DSpace storage to span several disks, and so there is no longer a hard limit on storage.
- Search improvements:
  - New fielded searching UI
  - Search results are now paged
  - Abstracts are indexed

- Better use of Lucene API; should stop the number of open file handles getting large
- Submission UI improvements:
  - now insists on a title being specified
  - fixed navigation on file upload page
  - citation & identifier fields for previously published submissions now fixed
- Many Unicode fixes to the database and Web user interface
- Collections can now be deleted
- Bitstream descriptions (if available) displayed on item display page
- Modified a couple of servlets to handle invalid parameters better (i.e. to report a suitable error message instead of an internal server error)
- Item templates now work
- Fixed registration token expiration problem (they no longer expire.)

# Appendix A

## List of Abbreviations

- AIFF** Audio Interchange File Format
- AIP** Archival Information Package
- API** Application programming interface
- CA** Certificate Authority
- CC** Creative Commons
- CNRI** Corporation for National Research Initiatives
- DBMS** Database Management System
- DC** Dublin Core
- DCMI** Dublin Core Metadata Initiative
- DDC** Dewey Decimal Classification
- DIDL** Digital Item Declaration Language
- DIM** DSpace Intermediate Metadata
- FAQ** Frequently asked questions
- GIF** Graphics Interchange Format
- HTML** HyperText Markup Language
- HTTPD** Hyper Text Transfer Protocol Daemon
- ID** Identifier
- IEC** International Electrotechnical Commission
- IR** Institutional Repository
- ISO** International Organization for Standardization
- ISBN** International Standard Book Number

- ISMN** International Standard Music Number
- ISSN** International Standard Serial Number
- JDBC** Java Database Connectivity
- JPEG** Joint Photographic Experts Group
- JSP** JavaServer Pages
- LCC** Library of Congress Classification
- LCSH** Library of Congress Subject Headings
- LDAP** Lightweight Directory Access Protocol
- MARC** Machine-Readable Catalog
- MESH** Medical Subject Headings
- METS** Metadata Encoding and Transmission Standard
- MIME** Multipurpose Internet Mail Extensions
- MIT** Massachusetts Institute of Technology,
- MODS** Metadata Object Description Schema
- MPEG** Moving Picture Experts Group
- NSI** Norwegian Science Index
- OAI** Open Archives Initiative
- OAI-PMH** Open Archive Initiative - Protocol for Metadata Harvesting
- OAIS** Open Archival Information System
- PDF** Portable Document Format
- QDC** Qualified Dublin Core
- RDBMS** Relational Database Management System
- RSS** Really Simple Syndication
- SFX** Self-extracting archive
- SGML** Standard Generalized Markup Language
- SICI** Serial Item and Contribution Identifier
- SIP** Submission Information Package
- SRB** Storage Resource Broker
- SRSC** Swedish Research Subject Categories
- SSL** Secure Sockets Layer

**TCP** Transmission Control Protocol

**TIFF** Tagged Image File Format

**UI** User Interface

**URI** Uniformed Resource Identifier

**URL** Uniform Resource Locator

**UTF** Unicode Transformation Format

**XML** Extensible Markup Language

**XSL** eXtensible Stylesheet Language

**XSLT** Extensible Stylesheet Language Transformations

## Appendix B

# Default Dublin Core Metadata Registry

Element	Qualifier	Scope Note
contributor		A person, organization, or service responsible for the content of the resource. Catch-all for unspecified contributors.
contributor	advisor	Use primarily for thesis advisor.
contributor	author	
contributor	editor	
contributor	illustrator	
contributor	other	
coverage	spatial	Spatial characteristics of content.
coverage	temporal	Temporal characteristics of content.
creator	Do not use; only for harvested meta-data.	
date		Use qualified form if possible.
date	accessioned	Date DSpace takes possession of item.
date	available	Date or date range item became available to the public.
date	copyright	Date of copyright.
date	created	Date of creation or manufacture of intellectual content if different from date.issued.
date	issued	Date of publication or distribution.
date	submitted	Recommend for theses/dissertations.
identifier		Catch-all for unambiguous identifiers not defined by qualified form; use identifier.other for a known identifier common to a local collection instead of unqualified form.
identifier	citation	Human-readable, standard bibliographic citation of non-DSpace format of this item
identifier	govdoc	A government document number
identifier	isbn	International Standard Book Number
identifier	ismn	International Standard Music Number
identifier	issn	International Standard Serial Number
identifier	sici	Serial Item and Contribution Identifier
identifier	uri	Uniform Resource Identifier

Element	Qualifier	Scope Note
contributor		A person, organization, or service responsible for the content of the resource. Catch-all for unspecified contributors.
contributor	advisor	Use primarily for thesis advisor.
identifier	other	A known identifier type common to a local collection.
description		Catch-all for any description not defined by qualifiers.
description	abstract	Abstract or summary.
description	provenance	The history of custody of the item since its creation, including any changes successive custodians made to it.
description	sponsorship	Information about sponsoring agencies, individuals, or contractual arrangements for the item.
description	statementofresponsibility	To preserve statement of responsibility from MARC records.
description	tableofcontents	A table of contents for a given item.
description	uri	Uniform Resource Identifier pointing to description of this item.
format		Catch-all for any format information not defined by qualifiers.
format	extent	Size or duration.
format	medium	Physical medium.
format	mimetype	Registered MIME type identifiers.
language		Catch-all for non-ISO forms of the language of the item, accommodating harvested values.
language	iso	Current ISO standard for language of intellectual content, including country codes (e.g. "en_US").
publisher		Entity responsible for publication, distribution, or imprint.
relation		Catch-all for references to other related items.
relation	isformatof	References additional physical form.
relation	ispartof	References physically or logically containing item.
relation	ispartofseries	Series name and number within that series, if available.
relation	haspart	References physically or logically contained item.
relation	isversionof	References earlier version.
relation	hasversion	References later version.
relation	isbasedon	References source.
relation	isreferencedby	Pointed to by referenced resource.
relation	requires	Referenced resource is required to support function, delivery, or coherence of item.
relation	replaces	References preceding item.
relation	isreplacedby	References succeeding item.
relation	uri	References Uniform Resource Identifier for related item.
rights		Terms governing use and reproduction.
rights	uri	References terms governing use and reproduction.
source		Do not use; only for harvested metadata.
source	uri	Do not use; only for harvested metadata.
subject		Uncontrolled index term.
subject	classification	Catch-all for value from local classification system. Global classification systems will receive specific qualifier
subject	ddc	Dewey Decimal Classification Number

<b>Element</b>	<b>Qualifier</b>	<b>Scope Note</b>
contributor		A person, organization, or service responsible for the content of the resource. Catch-all for unspecified contributors.
contributor	advisor	Use primarily for thesis advisor.
subject	lcc	Library of Congress Classification
subject	LCSH	Library of Congress Subject Headings
subject	mesh	Medical Subject Headings
subject	other	Local controlled vocabulary; global vocabularies will receive specific qualifier.
title		Title statement/title proper.
title	alternative	Varying (or substitute) form of title proper appearing in item, e.g. abbreviation or translation
type		Nature or genre of content.

Table B.1: Default Dublin Core Registry



## Appendix C

# Default Bitstream Format Registry

Mimetype	Short Description	Description	Support Level	Internal	Extensions
application/octet-stream	Unknown	Unknown data format	Unknown	Unknown	
text/plain	License	Item-specific license agreed upon to submission	Known	true	
application/marc	MARC	Machine-Readable Cataloging records	Known	true	
application/mathematica	Mathematica	Mathematica Notebook	Known	true	ma
application/msword	Microsoft Word	Microsoft Word	Known	true	doc
application/pdf	Adobe PDF	Adobe Portable Document Format	Known	true	pdf
application/postscript	Postscript	Postscript Files	Known	true	ai, eps, ps
application/sgml	SGML	SGML application (RFC 1874)	Known	true	sgm, sgml
application/vnd.ms-excel	Microsoft Excel	Microsoft Excel	Known	true	xls
application/vnd.ms-powerpoint	Microsoft Powerpoint	Microsoft Powerpoint	Known	true	ppt
application/vnd.ms-project	Microsoft Project	Microsoft Project	Known	true	mpd, mpp, mpx
application/vnd.visio	Microsoft Visio	Microsoft Visio	Known	true	vsd
application/wordperfect5.1	WordPerfect	WordPerfect 5.1 document	Known	true	wpd
application/x-dvi	TeX dvi	TeX dvi format	Known	true	dvi
application/x-filemaker	FMP3	Filemaker Pro	Known	true	fm
application/x-latex	LateX	LateX document	Known	true	latex
application/x-photoshop	Photoshop	Photoshop	Known	true	pdd, psd
application/x-tex	TeX	Tex/LateX document	Known	true	tex
audio/basic	audio/basic	Basic Audio	Known	true	au, snd
audio/x-aiff	AIFF	Audio Interchange File Format	Known	true	aif, aifc, aiff

Mimetype	Short Description	Description	Support Level	Internal	Extensions
audio/x-mpeg	MPEG Audio	MPEG Audio	Known	true	abs, mpa, mpega
audio/x-pn-realaudio	RealAudio	RealAudio file	Known	true	ra, ram
audio/x-wav	WAV	Broadcast Wave Format	Known	true	wav
image/gif	GIF	Graphics Interchange Format	Known	true	gif
image/jpeg	JPEG	Joint Photographic Experts Group/JPEG File Interchange Format (JFIF)	Known	true	jpeg, jpg
image/png	PNG	Portable Network Graphics	Known	true	png
image/tiff	TIFF	Tag Image File Format	Known	true	tif, tiff
image/x-ms-bmp	BMP	Microsoft Windows bitmap	Known	true	bmp
image/x-photo-cd	Photo CD	Kodak Photo CD image	Known	true	pcd
text/css	CSS	Cascading Style Sheets	Known	true	css
text/html	HTML	Hypertext Markup Language	Known	true	htm, html
text/plain	Text	Plain Text	Known	true	asc, txt
text/richtext	RTF	Rich Text Format	Known	true	rtf
text/xml	XML	Extensible Markup Language	Known	true	xml
video/mpeg	MPEG	Moving Picture Experts Group	Known	true	mpe, mpeg, mpg
video/quicktime	Video Quicktime	Video Quicktime	Known	true	mov, qt

Table C.1: Default Bitstream Format Registry

# Index

## Application Layer

- item importer and exporter, 123
- media filters, 129
- METS tools, 128
- OAI-PMH data provider, 119
- package importer and exporter, 122
- registering bitstreams, 126
- sub-community management, 130
- transferring items between instances, 125
- web user interface, 111

## Architecture

- application layer, 111
- business logic layer, 86
- overview, 80
- storage layer, 81

## Authentication

- custom authentication method, 60
- functional, 15
- LDAP, 59
- password, 59
- X.509, 59

## Authorization

- Bitstream, 15
- Bundle, 15
- Collection, 15
- Community, 15
- functional, 15
- Item, 15

## Bitstream

- Authorization, 15
- persistent identifier, 19

## Bitstream Format

- Support Level, 12

## Browse

- functional, 20

## Bugs, 35

## Build Process, 112

## Bundle

- authorization, 15

## Business Logic Layer

- administration toolkit, 102
- authorization, 103
- browse, 106
- checksum checker, 111
- content management API, 89
- core classes, 86
- e-person/group manager, 102
- handle manager, 104
- history system, 109
- plugin manager, 93
- search, 105
- workflow system, 101

## Collection

- authorization, 15

## Community

- authorization, 15

## Configuration

- OAI-PMH crosswalks, 55

## Configuration, 51

- bitstream format registry, 53
- checksum checker, 70
- configuration properties, 51
- controlled vocabularies, 68
- crosswalk plugins, 55, 72
- DIDL, 55
- e-mail messages, 53
- intermediate metadata format, 74
- item recommendations, 68
- LDAP authentication, 60
- media filters, 64
- metadata registry, 53
- METS crosswalk, 55
- other applications, 56
- packager plugins, 72
- plugins, 97
- RSS, 67
- statistical reports, 62
- submission license, 54

- thumbnail display, 65
- Creative Commons
  - functional, 21
- Crosswalk Plugins
  - functional, 14
- Customization, 51
  - authentication, 58, 60
  - e-mail messages, 53
  - item display, 57
  - item image preview, 65
  - license display, 67
  - Lucene analyzer, 66
  - OAI-PMH crosswalks, 55
  - statistical reports, shell scripts, 62
  - thumbnail link behaviour, 66
  - user interface, 56
- Data Model
  - Diagram, 12
- Directories and Files
  - log files, 78
  - source directory, 76
  - web application, 78
- E-Person
  - functional, 14
- Export
  - functional, 22
- Groups
  - functional, 14
- Handles
  - functional, 18
- History System
  - functional, 22
- HTML Support, 20
- HTTPS, 28
  - Tomcat, 29
- Import
  - functional, 22
  - registration, 22
- Ingest
  - functional, 16
  - visualization, 16
  - workflow, 16
- Installation
  - certificate, 29
  - compile, 27
  - create administrator, 27
  - cron jobs, 28
  - cron jobs, statistics, 28
  - cron jobs, vacuumdb, 28
  - download source code, 26
  - dspace.cfg, 26
  - Handle Server, 32
  - install, 27
  - Oracle, 26
  - Postgres, 26
  - quick installation, 25
  - update, 37
  - web application archives, 27
  - windows, 33
- Installation
  - create user, 26
- Internationalisation, 116
- Item
  - authorization, 15
- Item Export, 123
- Item Import, 123
- JSP Tags, 114
- Known Bugs, 35
- Metadata
  - descriptive, 13
  - Functional, 13
  - structural, 13
- Metadata
  - administrative, 13
- OAI
  - functional, 21
- OAI-PMH Data Provider, 119
  - about, 121
  - access control, 120
  - date stamp, 121
  - deletions, 121
  - resumption tokens, 121
  - sets, 120
  - unique identifier, 120
- OpenURL, 21
- Packager Plugins
  - functional, 13
- Plugin Manager
  - functional, 12
- Prerequisite Software

- Ant, [24](#)
  - Caucho Resin, [25](#)
  - Java SDK, [24](#)
  - Jetty, [25](#)
  - Oracle, [24](#)
  - Postgres, [24](#)
  - Relational Database, [24](#)
  - Tomcat, [25](#)
  - Unix, [24](#)
- Search
- functional, [20](#)
- Servlets and JSP's, [113](#)
- SSL on Apache, [31](#)
- Statistics
- functional, [22](#)
- Storage Layer
- backup, [85](#)
  - bitstream store, [83](#)
  - RDBMS, [81](#)
- Storage Resource Broker
- functional, [19](#)
- Submission
- thesis blocking, [119](#)
- Submission License
- default, [54](#)
- Subscriptions, [22](#)
- Supervision, [18](#)
- Version history, [132](#)
- Web User Interface, [111](#)
- Wiki, [8](#)
- Workflow
- functional, [17](#)
  - step 1, [17](#)
  - step 2, [17](#)
  - step 3, [17](#)