

MÍMIR

A Multi-paradigm Information Management Index and Repository

version 3.3.1-snapshot

October 26, 2011

Contents

1	Introduction	5
2	Installing and Managing Mimir	7
2.1	Mimir Architecture	7
2.2	Building and running the mimir-demo web application	8
2.2.1	Prerequisites	8
2.2.2	Building	9
2.2.3	Configuring	9
2.2.4	Running	11
2.3	Indexes in Mimir	12
2.3.1	Types of index	12
2.3.2	Creating a local index	13
2.3.3	Working with remote and federated indexes	14
2.3.4	Deleting indexes	18
2.4	“Deleting” Documents from a Mimir Index	19
3	Indexing Documents with Mimir	21
3.1	Configuring the indexer	21
3.2	Adding Documents to an Index	25
3.3	The default representation scheme	25
4	Searching Mimir Indexes	28
4.1	The Mimir query language	28
4.1.1	String Queries	29

4.1.2	AND Operator: “&”	30
4.1.3	OR Operator: “ ”	30
4.1.4	IN and OVER Operators	30
4.1.5	Repeats Operator: “+”	31
4.1.6	Sequence Queries and Gaps	31
4.1.7	Escaping Reserved Words	31
4.2	Search interfaces – how to submit queries to Mimir	32
4.2.1	Mimir Search Web Service	32
4.2.2	GUS example user interface	40
4.2.3	Embedding Mimir in a Grails application	42
5	Standard Mimir Plugins	44
5.1	The <code>db-h2</code> plugin	44
5.2	The <code>sesame</code> plugin	45
5.2.1	Searching in the knowledge base	45
5.3	The <code>measurements</code> plugin	46
5.3.1	Configuring the Measurements SAH	46
5.4	The <code>sparql</code> plugin	48
5.4.1	Creating a SPARQL helper	48
5.4.2	Format of SPARQL queries	50
6	Extending and Customising Mimir	52
6.1	Creating new Semantic Annotation Helpers	52
6.1.1	The <code>SemanticAnnotationHelper</code> interface	53
6.1.2	Abstract base classes	54
6.1.3	Packaging new helper types for use with Mimir	55
6.2	Embedding Mimir in a custom Grails application	56
6.2.1	Installing the Mimir plugin	56
6.2.2	URL mappings	56
6.2.3	Security considerations	57

A	Change Log	58
A.1	Version 3.3.0 (October 2011)	58
A.2	Version 3.2.0 (May 2011)	58

Chapter 1

Introduction

Mimir is a multi-paradigm information management index and repository which can be used to index and search over text, annotations, semantic schemas (ontologies), and semantic meta-data (instance data). It allows queries that arbitrarily mix full-text, structural, linguistic and semantic queries and that can scale to gigabytes of text.

A typical semantic annotation project deals with large quantities of data of different kinds. Mimir provides a framework for implementing indexing and search functionality across all these data types, listed below in the order of increasing information density:

Text

All documents have a textual content. Support for full text search represents the most basic indexing functionality and it is required in most (if not all) cases. Even when semantic annotation is used to abstract away from the actual textual data, the original content still needs to be accessible so that it can be used to provide textual query fragments in the case of more complex conceptual queries.

Mimir uses inverted indexes¹ for indexing the document content (including additional linguistic information, such as part-of-speech or morphological roots), and for associating instance of annotations with the position in the input text where they occur. The inverted index implementation used by Mimir is based on MG4J².

Annotations

The first step in abstracting away from the plain text content is the production of *annotations*. Annotations are meta-data associated to text snippets in the documents. Mimir's view of annotations is based on that of GATE, with each annotation described by

- the document it belongs to;
- the start and end offset of the referred text snippet;

¹*Inverted Indexes* are data structures traditionally used in Information Retrieval to support indexing of text.

²<http://mg4j.dsi.unimi.it/>

- the annotation type;
- an arbitrary set of <feature,value> pairs.

An annotation index supports a more generic search paradigm. Depending on the type of annotations available, the user can search across different dimensions. If, for example, the documents are annotated with occurrences of **Person**, **Location**, **Organization** entities, then searches like {**Person**}, CEO of {**Organization**}, based in {**Location**} become possible. Storage of annotation data in Mimir indexes is handled by plugins, Mimir ships with two storage plugins by default, one storing annotation data in a relational database and the other in a Knowledge Base to support richer semantic querying.

ANNIC (ANNotations In Context)³ is a tool predating Mimir that supports the indexing of annotations, and that has been used to inform the design of Mimir.

Knowledge Base Data

Knowledge Base (KB) Data consists of an ontology populated with instances. The ontology represents the data schema and comprises a hierarchy of class types and a hierarchy of properties that are applicable between instances of classes. The instance data represents facts that are known to the systems and is typically at least partially derived from semantic annotation over documents. KB data is used to reach a higher level of abstraction over the information in the documents which enables conceptual queries such as “find all mentions of **Persons** who are employed by any organization based in Yorkshire”.

A KB that is pre-populated with appropriate world knowledge can perform other generalisations that are natural to humans users, such as being able to identify Vienna as a valid answer to queries relating to Austria, Europe or the Western Hemisphere.

As mentioned above, Mimir can make use of a Knowledge Base to store information relating to annotations. The links between annotations, the textual data, and the knowledge base information are created by the inclusion into the text indexes of a set specially-created URIs that are associated with annotation data. Furthermore, URIs of entities from the Knowledge Base can be stored as annotation features.

Knowledge bases are typically represented as a collection of triples that are kept in highly-specialised and optimised triple stores, using standards such as RDF or one the versions of OWL⁴. The implementation used by Mimir is based on ORDI and OWLIM⁵.

³See <http://gate.ac.uk/userguide/chap:annic>.

⁴See <http://www.w3.org/RDF/> and <http://www.w3.org/TR/owl-features/>.

⁵See <http://www.ontotext.com/ordi/> and <http://www.ontotext.com/owlim/>.

Chapter 2

Installing and Managing Mimir

2.1 Mimir Architecture

Mimir is divided into a number of related modules.

mimir-core The core Java library to create a Mimir index on disk, add GATE documents to the index, and then query the index once it has been built. Also provides some abstract helper classes for the annotation storage layer, but not the actual storage implementations (which are provided by separate plugins, leveraging the CREOLE plugin framework of GATE Embedded).

plugins/db-h2 The default annotation storage implementation. This stores annotation data using H2, an in-process embedded SQL database.

plugins/sesame An alternative annotation storage implementation that stores its annotation data in a triple store using the Sesame API (<http://www.openrdf.org/>). Annotations with an “inst” feature are treated as links into the knowledge base, supporting richer semantic queries.

plugins/sparql A helper that is layered on top of a generic storage implementation to provide semantic querying against a separate knowledge base, accessible at a SPARQL endpoint.

plugins/measurements A special-purpose helper for Measurement annotations created by the GATE `Tagger.Measurements` plugin. Queries are normalized into SI units so can retrieve annotations that express the same measurement in different terms (e.g. an annotation for “90 seconds” would match a query for “1 to 2 minutes”).

mimir-client The client side of the Mimir remote protocol, to support distributed indexing and querying.

mimir-web A Grails plugin (<http://grails.org>) providing both the user interface to create and query indexes over the web, and also the server side of the

remote protocol to expose several distributed Mímir indexes as a single *federated* index for clients. This is provided as a plugin rather than an application to make it more easily customizable.

mimir-demo The simplest possible Grails application using the above plugin. This application is not actually included in the source tree; it is instead generated by the build process. See Section 2.2 for details on how to do that. Once created, it is ready to use as-is or as a base for customization to add security, look and feel customization, etc.

2.2 Building and running the mimir-demo web application

The **mimir-demo** module is intended to be the simplest possible demonstration of how to use the Mímir grails plugin (**mimir-web**) in a Grails application. This application is suitable for experimentation and local use, but as it includes no security or authentication support it may not be appropriate for production deployment. These concerns are addressed further in section 6.2.

2.2.1 Prerequisites

To build the Mímir demo app you will need:

- A Java 6 JDK. Mímir has been tested with the Sun/Oracle and OpenJDK JVMs on Linux and the Apple JVM on Mac OS X.
- Apache Ant 1.8.1 or later.
- Grails (<http://grails.org>). The Mímir plugin was developed with version 1.3.7 but later 1.3.x versions should also work (you will need to run **grails upgrade** in the **mimir-web** directory if you use a version other than 1.3.7). You need to set the `JAVA_HOME` environment variable to point to your JDK, the `GRAILS_HOME` environment variable to point to your Grails installation and put `$GRAILS_HOME/bin` on your `PATH`.
- Google Web Toolkit (<http://code.google.com/webtoolkit>). The GUS web interface is implemented in GWT, using the Grails gwt plugin, but a standalone installation of GWT is required to build the app. Version 2.2.0 or later is recommended, and you need to set the `GWT_HOME` environment variable to point to your GWT installation.

While not strictly a pre-requisite, Mímir performs much better on 64-bit systems than on 32-bit ones, partly due to simply being able to assign more memory to the process, but also because the larger address space allows MG4J to memory-map many of the files that make up the index.

To run a local instance of Mímir you can use the standard **grails prod run-war** command, but to deploy a production instance you will need a separate servlet container such as Tomcat.

2.2.2 Building

There is a top-level Ant build.xml file that should build all the modules in the correct order. To do that simply change to the top level directory containing the Mimir source code, and run the `ant` command. To perform the same build process manually, you need to change to the following directories and run the following `ant` commands in this order:

1. `mimir-core: ant publish`
2. `mimir-client: ant publish`
3. `plugins`: run `ant` in each subdirectory of the `plugins` directory in turn (order is not important here, the plugins do not depend on one another).
4. `mimir-web - grails compile` followed by `grails compile-gwt-modules`.

Once the plugin is built you need to generate the source folder for the demo application. You can do that by running the `ant mimir-demo` command in the top level folder of the Mimir source tree. This will create a new directory (named `mimir-demo`) inside the top level directory, which contains a simple Grails-based web application, that uses the Mimir Grails plugin.

The next step is to configure the `mimir-demo` web application, and is described in the following Section.

2.2.3 Configuring

When the Mimir Grails plugin is installed into a Grails application, it creates a base configuration file at `grails-app/conf/MimirConfig.groovy`. This file contains a number of settings that affect the running of the Mimir components.

```
1 gateInit {
2     gateHome = "WEB-INF/gate-home"
3     userConfigFile = "WEB-INF/gate-home/user.xml"
4 }
```

Since Mimir is based on GATE, the plugin initializes the GATE environment at startup. These parameters control the initialization process. In most cases you can leave the values at their defaults, which use a deliberately cut-down set of GATE configuration files installed into `web-app/WEB-INF` by the Mimir Grails plugin. The available parameters are `gateHome`, `pluginsHome`, `siteConfigFile`, `userConfigFile` and `builtinCreoleDir`, which correspond to the standard settings on the Gate class, and their values can be either absolute URLs (such as `file:/opt/gate`) or paths which are taken relative to the web application (i.e. the `web-app` directory of the Grails application).

```
1 plugins {
2     h2 = "../plugins/db-h2"
3     myCustomPlugin = "file:/data/mimir/plugins/myCustomPlugin"
4 }
```

This section specifies the Mimir plugins that should be loaded, and determines the kinds of annotation helpers you will be able to use in your indexes. You generally need at least one of the standard `db-h2` and/or `sesame` plugins to be able to do anything useful with Mimir, and you may want the `measurements` plugin as well if you will be searching on Measurement annotations and/or the `sparql` plugin if you have an external knowledge base. Section 3.1 has more information about the standard annotation helpers, and section 6.1 discusses how to implement your own custom ones.

Mimir uses the GATE plugin mechanism, so Mimir plugins are actually very simple CREOLE plugins¹, used to add a set of `jar` files to the current classpath.

Plugins can be specified either as absolute URLs or as paths relative to the Grails application base directory. Absolute URLs will be loaded as such both in `run-app` and in WAR deployments, but plugins specified as relative paths are treated slightly differently. They will be loaded directly from the specified paths in `run-app`, but when building a WAR file the referenced plugins will be packaged inside the WAR file and loaded from there at runtime.

```
1 queryTokeniserGapp = "WEB-INF/gate-home/default-query-tokeniser.xgapp"
```

Whereas GATE's usual data model deals with annotations in terms of their *character* offsets from the start of the document, Mimir deals in terms of *tokens*. Queries for plain text strings in Mimir must be tokenised before they can be matched against the index, and the tokenisation applied to the queries must match that applied to the documents that have been indexed. The Mimir Grails plugin uses a saved GATE application state (gapp file) to perform query tokenisation, the location of which is specified here. Again, the location can be an absolute URL or a path relative to the web-app directory, and the default value refers to a simple app installed by the Mimir Grails plugin that contains a single ANNIE tokeniser with its default settings.

If your tokenisation requirements are more complex, you can provide your own saved application, or alternatively you can use your application's `resources.xml` or `resources.groovy` to override the definition of the Spring bean named "query-Tokeniser" – this bean must define a GATE LanguageAnalyser that will produce annotations of type `Token` in the default annotation set.

```
1 indexBaseDirectory = "mimir-indexes"
```

This is the only parameter that you will definitely want to change. It refers to the directory in which newly created Mimir indexes will be stored. Unlike the other parameters discussed above, this parameter is a native file path and a relative path will be resolved against the current working directory of the Mimir process, *not* against the web-app directory. Typically this parameter would be an absolute path, such as `/data/mimir/indexes`.

You can modify these settings directly in `MimirConfig.groovy` or you can specify them in your application's normal `Config.groovy` (or an external file loaded using

¹See <http://gate.ac.uk/userguide/chap:creole-model>.

`grails.config.locations`), under the `gate.mimir` group. In other words, this `MimirConfig.groovy`:

```
1 indexBaseDirectory = "/data/mimir/indexes"
```

is equivalent to this in `Config.groovy`:

```
1 gate {
2     mimir {
3         indexBaseDirectory = "/data/mimir/indexes"
4     }
5 }
```

or this in an external `.properties` file:

```
1 gate.mimir.indexBaseDirectory=/data/mimir/indexes
```

Note that because of the special handling at build time of plugins referenced as relative paths (see above), if you want to load additional plugins into a WAR-packaged Mimir using runtime settings in an external configuration file, then the plugins must be specified using absolute URLs, i.e. `gate.mimir.plugins.custom = "file:/opt/mimir/plugins/custom"`. Relative plugin paths are ignored at runtime by Mimir when running from a WAR deployment. However, since Mimir plugins are simply standard GATE CREOLE plugins and the Mimir Grails plugin initializes GATE Embedded using Spring you can load extra plugins relative to your web app by using Spring configuration in `WEB-INF/spring/resources.xml` (see <http://gate.ac.uk/userguide/sec:api:spring> for details):

```
1 <beans xmlns="http://www.springframework.org/schema/beans"
2     xmlns:gate="http://gate.ac.uk/ns/spring"
3     xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4     xsi:schemaLocation="
5         http://www.springframework.org/schema/beans
6         http://www.springframework.org/schema/beans/spring-beans.xsd
7         http://gate.ac.uk/ns/spring
8         http://gate.ac.uk/ns/spring.xsd">
9     <gate:extra-plugin>WEB-INF/custom-plugin</gate:extra-plugin>
10 </beans>
```

2.2.4 Running

The easiest way to run the Mimir demo web app is to use the normal Grails commands `run-app` or `run-war`. For performance, `grails prod run-war` is preferable. For anything more than the smallest toy index it is advisable to increase the memory available to Mimir by using the `JAVA_OPTS` environment variable. For example (using bash or a similar POSIX shell):

```
$ JAVA_OPTS='-Xmx4G' grails prod run-war
```

To shut down a web app started using `run-app` or `run-war`, simply create an empty file in the `mimir-demo` directory named `“.kill-run-app”`. Grails watches for this file and will shut down gracefully when it detects that the file has been created.

For production deployments, a better option is to build a WAR file using `grails prod war` and deploy that to a standalone servlet container such as Apache Tomcat. If you are using Ubuntu or Debian GNU/Linux, it is better to download the standard Tomcat ZIP package from Apache and use that rather than installing the Tomcat available through `apt-get` as the latter is configured by default with a security manager that interferes with Mimir.

When deployed to a servlet container the demo application reads configuration at runtime from two locations using the Grails standard “externalized configuration” mechanism:

- `WEB-INF/classes/mimir-app-config.groovy` inside the web application.
- `mimir-config.groovy` in the working directory of the Java process.

Any values in these files override values specified in `MimirConfig.groovy` or the main application `Config.groovy`. For production deployments, you should be sure to specify the public URL of your Mimir server in one of these configuration files. For example:

```
1 gate.mimir.indexBaseDirectory = "/data/mimir/indexes"
2 grails.serverURL = "http://example.com/mimir"
3 // or just http://example.com if you have deployed Mimir
4 // as the ROOT web application
```

2.3 Indexes in Mimir

2.3.1 Types of index

A single instance of Mimir can host several indexes. Mimir supports *local* indexes, stored in the file system of the Mimir server, and *remote* indexes, which are a view of an index hosted in another Mimir instance (possibly on a different machine). Several indexes (of any type) can be combined into a *federated* index, which presents the group of indexes as a single virtual index. All the indexing and searching functionality of Mimir applies equally to all three index types.

Each Mimir index has a *state*, and the operations that can be performed on the index depend on which state it is currently in. When first created, a local index will be in the *indexing* state, meaning it is waiting for documents to be added to the index. When all the documents have been added to the index an administrator will close the index, putting it into the *closing* state. For large indexes the closing process can take several hours, and when it is complete the index will enter the *searching* state, at which point it is available for querying. The other possible state for a local index is *failed*, indicating a problem with the index. Typically a failed index will need to be deleted by the administrator. Thus it is apparent that an

index cannot be used simultaneously for searching and indexing. An existing set of index files can be imported into a running Mímir instance as a local index, which will then immediately be in the *searching* state.

Remote indexes inherit their state from the remote server, and federated indexes inherit their state by combining the states of their component indexes. A federated index may occasionally appear in the *working* state if its component indexes are not all in the same state (for example if some of them have started closing but others are still in *indexing* mode), but the working state will usually resolve to a normal state once the component indexes have synchronized.

Note that once a local index has moved from *indexing* to *closing* to *searching* it is not possible to add more documents to the same index. The suggested way to add to an index is to create a new index to hold the new documents, fill it, close it, and then create a federated index consisting of the original index plus the new one (or if the original index was itself federated, add the new index to the existing federation).

A typical setup for a large-scale indexing task would be to have a number of identical “slave” servers running Mímir, each with a single local index. A single “master” Mímir instance could then have one remote index definition pointing to each of the slaves, and a single federated index combining the remote indexes. This federated index would be the point of entry into the system and would share out indexing jobs (round-robin among the slaves) or search requests (to all the slaves in parallel) as appropriate.

2.3.2 Creating a local index

Indexes in Mímir are managed through the web interface. The front page of a newly-installed Mímir is shown in Figure 2.1. The *index templates* mentioned at the bottom, are used to define the properties of new indexes, and are described in more detail in Chapter 3. The Mímir Grails plugin provides a single example template based on ANNIE annotation types.

To create an empty local index ready to receive documents for indexing, select the *create a new local index* link. This will present a form (Figure 2.2) asking for the name of the new index and the template from which it should be created. The “Document URIs are external links” option affects the way documents are presented in the search interface. Every document in Mímir is identified by a URI, and if you intend to use document URIs that are actually resolvable URLs (for example if your documents came from a web crawl) then you should select this option to add a link to the original document to the search results. If the document URIs will not be resolvable URLs then leave the option un-selected. The index will be assigned a unique identifier and a new directory will be created under the `indexBaseDirectory` you configured earlier to hold the index data. The newly-created index will be in the *indexing* state (see Figure 2.3), ready to receive documents for indexing. For details of how to submit documents to the index, see Chapter 3.

This *index information* page can be accessed at any time by clicking the link for the relevant index name from the Mímir front page (Figure 2.4). Once all the documents to be indexed have been submitted the index can be closed using the

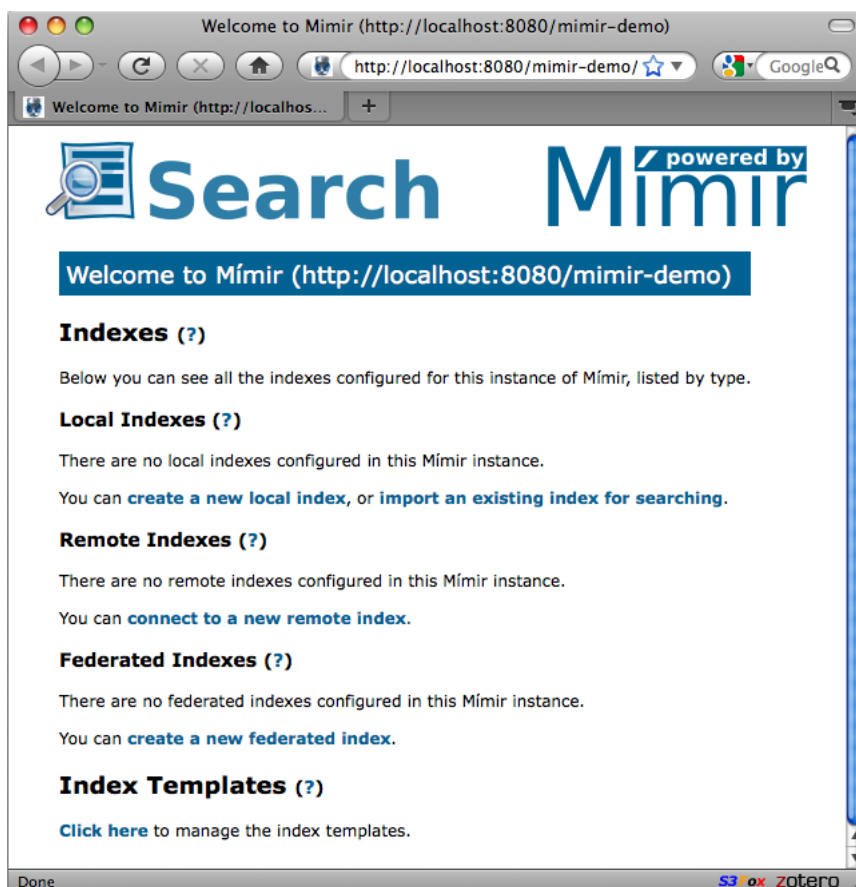


Figure 2.1: The default front page of a new Mimir

Close link on the index information page. This will change the state of the index to *closing* as described above and begin the closing process. The information page will show a live-updating progress bar (Figure 2.5) giving some indication of the time remaining until the index has completely closed.

When the closing process is complete (the progress bar reaches 100%) the index will switch into *searching* mode, and the index can then be searched using the tools described in Chapter 4.

2.3.3 Working with remote and federated indexes

The architecture of Mimir is designed to make working with remote and federated indexes as transparent as possible. The setup process will obviously vary for the different index types, but once created the process of submitting documents for indexing or of performing queries is exactly the same for all indexes.

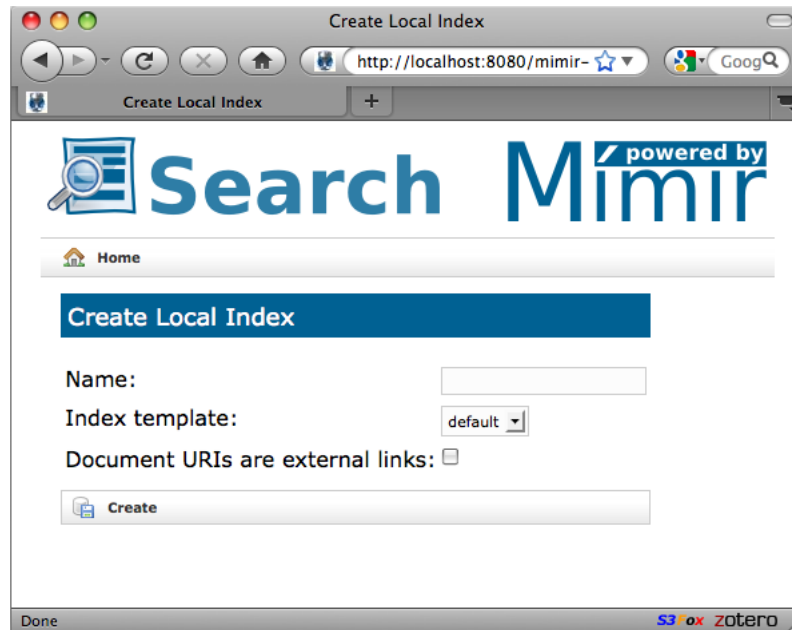


Figure 2.2: Creating a new local index



Figure 2.3: Results of creating a new local index

Local Indexes (?)

The following local indexes are configured:

1. [index-1](#)

You can [create a new local index](#), or [import an existing index for searching](#).

Figure 2.4: List of local indexes on the Mimir front page

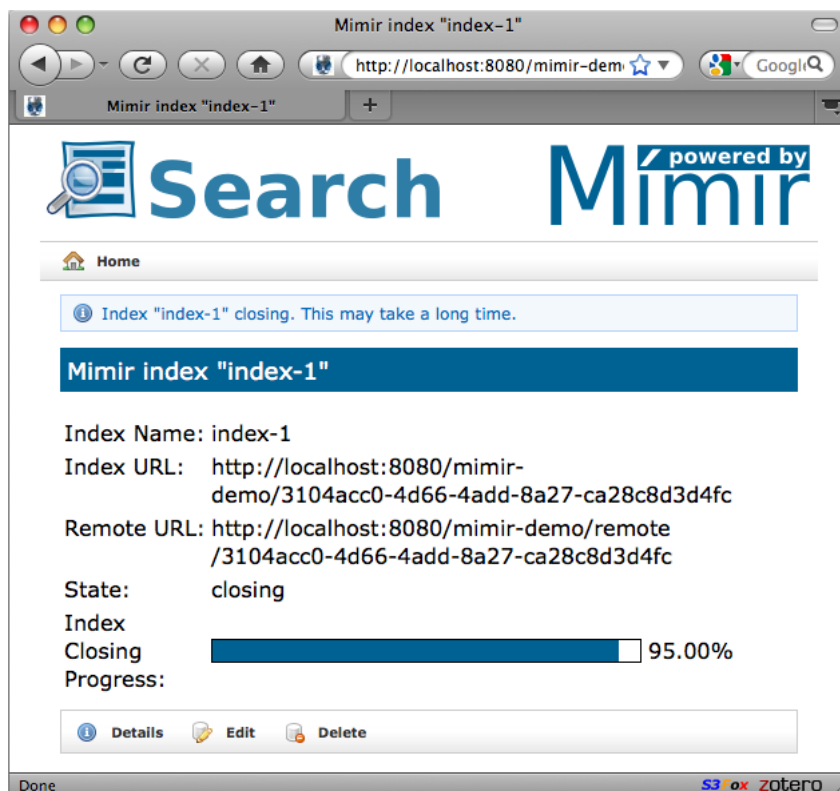


Figure 2.5: Closing a local index

Remote indexes

A *remote* index is a mechanism whereby one Mimir instance can transparently index documents in, or send queries to, an index that is located in a different Mimir instance, typically running on separate hardware. To connect one *master* Mimir instance to an index running in another *slave* instance, first visit the index information page for the relevant index on the slave and make a note of its *remote URL* (typically a URL of the form `http://server:port/mimir/remote/{UUID}`). Now on the front page of the master instance, select the *connect to a new remote index* link. This will present a form (Figure 2.6) asking for a name for the remote index (which need not be the same as the name of the index on the slave), and a *remote URL* which is the one you made a note of from the slave above. You should

never create a remote index pointing to another index in the same Mimir instance. Such a configuration is not supported and will lead to errors!

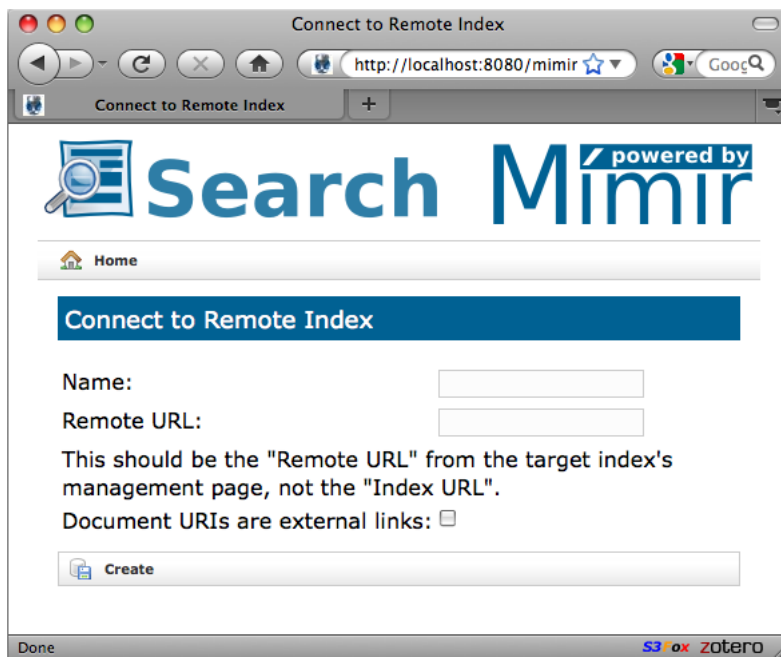


Figure 2.6: Connecting to a remote index

The remote index defined on the master server will synchronize its state with that of the underlying index on the slave, and once created will be usable exactly like a local index. However remote indexes are rarely used directly, as in most cases it is more efficient to operate on the slave instance itself. The main benefit of remote indexes comes when they are used as part of a *federated index*.

Federated indexes

A *federated index* is a device to bundle several indexes (which can themselves be local, remote or federated) together so they can be used as a single index. Documents for indexing are shared out between the component *sub-indexes*, and searches are performed by all sub-indexes in parallel. Thus a federation of five indexes each containing 200,000 documents will typically run queries faster than a single index containing 1 million documents. To create a federated index, go to the Mimir front page and select the *create a new federated index* link. This will present a form (Figure 2.7) asking for a name for the federated index. The form also includes a multiple-selection list to specify the sub-indexes to be included in the federated index. Select the appropriate entries from this list using the usual multiple list selection mechanism (ctrl-click on Windows or Linux, cmd-click on Mac OS X) and press the *Create* button to create the index. Once created the federated index will be usable exactly like a local or remote index.

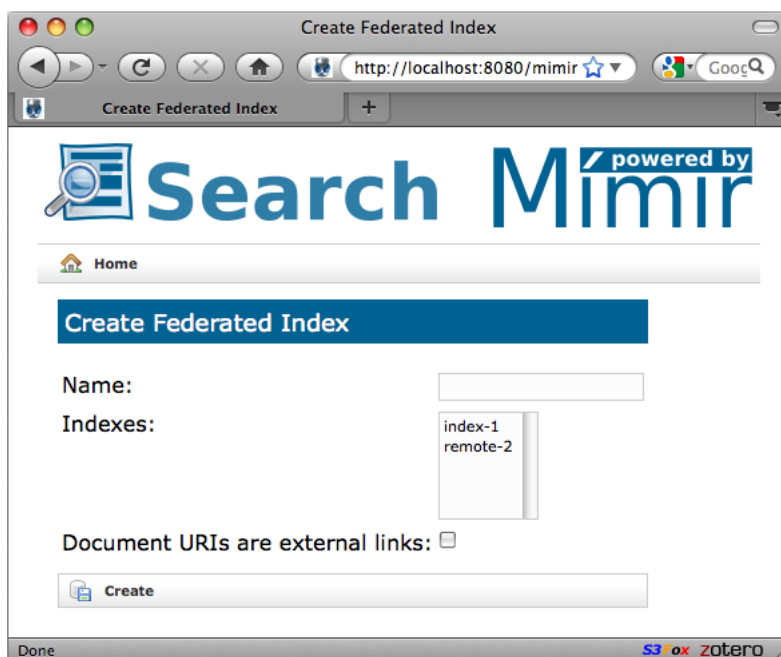


Figure 2.7: Creating a new federated index

2.3.4 Deleting indexes

If an index registered with Mimir is no longer required it can be deleted by selecting the *Delete* button from the index information page (accessible by clicking on the name of the relevant index on the Mimir front page). For remote and federated indexes this simply deletes the “registration” of the index with Mimir, which can be easily re-created as above. For local indexes it also offers the option to delete the underlying index files from disk. If a local index in *searching* state is deleted without deleting the disk files then the index can be re-created later using the *import an existing index for searching* option from the Mimir front page. However, if a local index in the *indexing* or *closing* state is deleted without having been properly closed then the index files will be unusable and will need to be deleted manually.

Mimir will not allow the deletion of an index which is currently part of a federated index in the same Mimir instance. To delete such an index, it must first be removed from the federated index. This guarantee only applies to indexes within a *single* Mimir instance — Mimir does not prevent the deletion of an index on a slave instance which is being used as a remote index by a master instance (it prevents the deletion of the remote index definition in the master but not the slave index it points to). However to do so would put the remote index on the master (and hence any federated index that it is part of) into the *failed* state, preventing further use until the problem is resolved.

2.4 “Deleting” Documents from a Mimir Index

While Mimir indexes are not directly modifiable once they have been created, there are situations in which it is necessary to remove documents that should not have been indexed in the first place, or documents that may be considered libellous, etc. To support this, Mimir provides a mechanism to mark individual documents in the index as “deleted”, and any documents so marked will be excluded from future queries. It is not possible to completely delete the data from the index files on disk, short of completely re-building the index from scratch, but documents marked as deleted are not accessible through any of the public Mimir APIs or user interfaces.

To mark a document as deleted (or to remove an existing deletion marker, making the document available for queries again), use the “Manage deleted documents” link from the index’s administration page. This will present the screen shown in figure 2.8, with a text box into which you can type one or more (space-separated) document IDs, and choose whether to mark them as deleted or as “not deleted” (i.e. to remove any existing deletion markers for those document IDs).

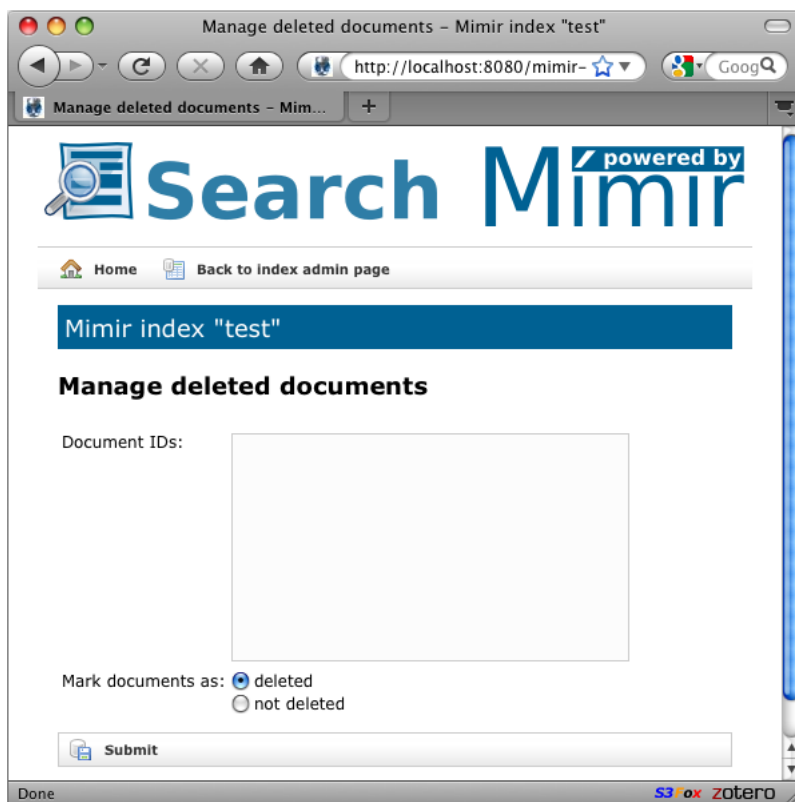


Figure 2.8: Managing deleted documents

Note that the IDs required here are not the URIs that were provided with the documents when they were indexed, but the internal Mimir document IDs which are numbers starting from 0, as returned in the hit lists and “doc-

Stats” by the Mímir query APIs (see section 4.2.1). In the GUS interface (section 4.2.2) the document ID is part of the URL to the “cached” document text, `.../search/gusDocument/<documentID>?queryId=...`

Chapter 3

Indexing Documents with Mimir

Mimir is designed to index semantically annotated documents. It accepts as input GATE documents¹ and produces a set of indexes as a result. The way the text and annotations of the input documents are converted into indexes is controlled through configuration options.

3.1 Configuring the indexer

In the Mimir web interface, the configuration of a new index is represented by an *index template*. This specifies:

- which annotation types and features to index
- which annotation sets contain these annotations
- how to handle the document format and metadata

Index templates can be managed using the “Click here to manage the index templates” link at the bottom of the Mimir front page. An index template is specified in a structured “domain specific language” using Groovy — Listing 3.1 shows an example of the default template provided by the Mimir Grails plugin.

¹<http://gate.ac.uk/userguide/chap:corpora>

```

1  import gate.creole.ANNIEConstants
2  import gate.mimir.index.OriginalMarkupMetadataHelper
3  import gate.mimir.db.DBSemanticAnnotationHelper as DefaultHelper
4
5  tokenASName = "mimir"
6  tokenAnnotationType = ANNIEConstants.TOKEN_ANNOTATION_TYPE
7  tokenFeatures = {
8      string()
9      category()
10     root()
11 }
12
13 semanticASName = "mimir"
14 semanticAnnotations = {
15     index {
16         annotation helper:new DefaultHelper(annType:'Sentence')
17     }
18     index {
19         annotation helper:new DefaultHelper(annType:'Person',
20                                             nominalFeatures:["gender"])
21         annotation helper:new DefaultHelper(annType:'Location',
22                                             nominalFeatures:["locType"])
23         annotation helper:new DefaultHelper(annType:'Organization',
24                                             nominalFeatures:["orgType"])
25         annotation helper:new DefaultHelper(annType:'Date',
26                                             integerFeatures:["normalized"])
27     }
28 }
29 documentRenderer = new OriginalMarkupMetadataHelper()
30 documentMetadataHelpers = [documentRenderer]

```

Listing 3.1: The default index template provided with Mímir

The various sections of the template are as follows:

Imports

```

1  import gate.creole.ANNIEConstants
2  import gate.mimir.index.OriginalMarkupMetadataHelper
3  import gate.mimir.db.DBSemanticAnnotationHelper as DefaultHelper

```

The template can optionally start with import statements to import any Java classes that are used further on in the template.

Token definitions

```

5  tokenASName = "mimir"
6  tokenAnnotationType = ANNIEConstants.TOKEN_ANNOTATION_TYPE
7  tokenFeatures = {
8      string()
9      category()
10     root()
11 }

```

The next section of the template deals with the *tokens* that Mimir will base its index on. Mimir sees every document as a stream of tokens rather than a stream of characters, and all the annotations indexed by Mimir are stored in terms of their starting and ending *tokens*, not character offsets. Thus for Mimir to work correctly it needs to know how to split up the document into tokens and what information to store about each token. For this purpose it uses GATE annotations, and indexes the values of features on the annotations.

The following options can be configured:

tokenASName The name of the annotation set in which the token annotations can be found. This may be left unspecified (or explicitly set to `null`, without quotes) to use the default annotation set, which has no name.

tokenAnnotationType The annotation type that should be used as tokens. This entry is required, and can generally be simply set to the default `ANNIEConstants.TOKEN_ANNOTATION_TYPE` (with a suitable `import` at the top of the template).

tokenFeatures A block of code giving the features from each token annotation that should be indexed.

The *tokenFeatures* block should list the features to be indexed as shown in the example, each feature name followed by parentheses. For advanced users an `MG4J TermProcessor` instance may be provided inside these parentheses. By default, if no term processors are specified, the *first* feature is converted to lowercase and the subsequent features are not modified. Since terms in a query are processed using the same processor as those in the index, this has the effect of making searches on the first feature case-insensitive, and searches on the other features case-sensitive. To stop any processing being done, you should supply a `it.unimi.dsi.mg4j.index.NullTermProcessor` value, by specifying e.g. `string(NullTermProcessor.getInstance())`, after including the relevant `import` statement at the top.

Semantic annotations

```
13  semanticASName = "mimir"
14  semanticAnnotations = {
15    index {
16      annotation helper:new DefaultHelper(annType:'Sentence')
17    }
18  }
```

The next section defines the *semantic annotations* that Mimir will include in the index. Each *index* block defines a set of semantic annotation types that will be indexed and stored together in one index. The choice of how to group annotation types together into indexes can affect the indexing speed, as the annotations within

one index are processed sequentially by a single thread, whereas types in separate indexes can be indexed in parallel.

Each annotation type to be indexed is introduced by “**annotation**”. This is a method call in the Groovy DSL which takes the following named arguments:

helper The *semantic annotation helper* Java object that should be used to index this annotation type.

type The annotation type that is to be indexed. When using the default semantic annotation helper types this can be omitted.²

Semantic Annotation Helpers

Semantic annotations are stored in special indexes that associate URIs with document positions. During indexing, the role of the helper implementations is to store the necessary information about each annotation to be indexed in a persistent form and return one or more URIs that identify it.

One could make a distinction between *generic* semantic annotation helper types, which can be configured to handle any annotation type and features, and *special-purpose* helpers that are designed to handle specific annotation types. Mímir supplies two generic helper implementations in the **db-h2** and **sesame** plugins³ that store annotation information in a relational database and a knowledge base respectively. For the most standard cases, one or other of these default helper implementations should be sufficient. One sample special-purpose helper for **Measurement** annotations (as generated by the GATE **TaggerMeasurements** plugin) is also provided, in the **measurements** plugin. This is intended both to be useful in its own right and to serve as a template for how to implement your own helpers for other complex annotation types. The **sparql** plugin provides a helper that can wrap any other helper and add the ability to query for URI-valued features by making a query to a SPARQL endpoint. These plugins are discussed in detail in chapter 5.

Note for users upgrading from Mímir 3.2.0 and earlier: the previous index template DSL style using the annotation type as the method name and the `nominalFeatures` etc. as parameters is still supported but should be considered deprecated. You should consider porting your index templates to the new style, as support for the old style may be removed in a future release.

Document rendering and metadata

```
25 documentRenderer = new OriginalMarkupMetadataHelper()
26 documentMetadataHelpers = [documentRenderer]
```

The final part of the template concerns how document-level metadata is indexed, and how this can be combined with the document text to render the document

²In particular, if the specified helper has a method “`getAnnotationType()`” then this will be called and the returned value used as the annotation type. All the standard helpers provided with Mímir extend `AbstractSemanticAnnotationHelper` which implements this method.

³A third implementation in the **ordi** plugin is now deprecated. This stores the data in the same underlying OWLIM semantic repository format but accesses it through a different API abstraction.

content at search time, with matches of the query highlighted. These tasks are performed by objects that implement the interfaces `DocumentMetadataHelper` and `DocumentRenderer` respectively (both in the `gate.mimir` package). Mimir provides a single class `gate.mimir.index.OriginalMarkupMetadataHelper` which implements both interfaces, so in most cases the same object can be used for both jobs.

An index template must define one `documentRenderer` and may define any number of `documentMetadataHelpers` (in a square-bracketed list). If the renderer is an `OriginalMarkupMetadataHelper` (or a subclass) then the renderer object must be included in the list of metadata helpers in order to function correctly. Other metadata helpers may be added to the list if required.

3.2 Adding Documents to an Index

Once an index has been created in *indexing* mode, the next stage is to add documents to the index. Mimir provides an HTTP API for this which accepts documents for indexing via HTTP POST requests that include the document in Java serialised format. The easiest way to make use of this API is via GCP (the GATE Cloud Paralleliser batch processing tool) using a `MimirOutputHandler`. This GCP output handler makes use of the `gate.mimir.index.MimirConnector` (in the `mimir-client` module) to actually make the remote call, and you can use the same API in your own code. To add a GATE document to an open index simply call:

```
1 MimirConnector.defaultConnector().sendToMimir(document, uri,  
    indexUri);
```

... with the following parameters:

document a `gate.Document` for indexing.

uri the URI that should be used to identify the document in the Mimir index. May be `null`, in which case Mimir will generate a URI, but in most cases there will be a more meaningful identifier that could be used.

indexUri a `java.net.URL` pointing to the location of the Mimir index. This is the “Index URL” given on the index information page.

The document to be indexed must, of course, contain the token and semantic annotations that the index expects.

It is possible to create your own private instance of `MimirConnector` rather than simply using the default one, but this is not necessary in normal use.

3.3 The default representation scheme

The default generic SAH implementations try to minimise the amount of data stored in their underlying database or semantic repository by creating representation templates that are shared between all occurrences of annotations with the same values for the features. There are two levels of templates, the first defined by the values of

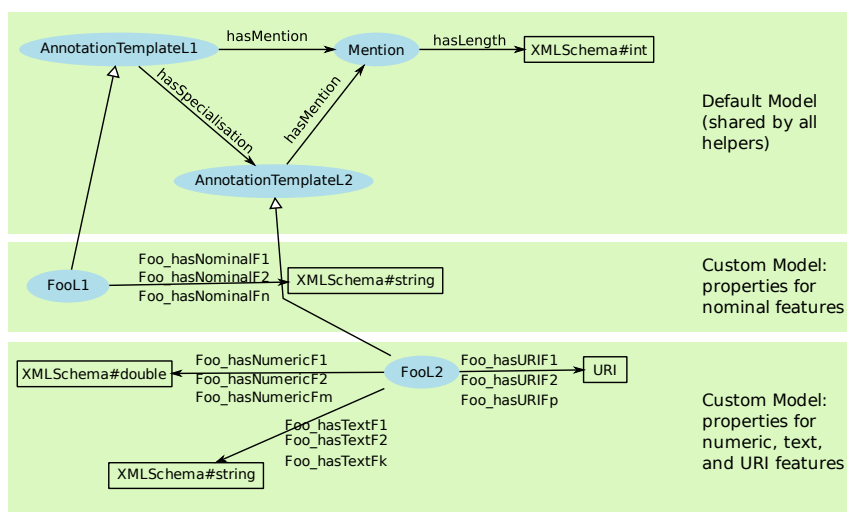


Figure 3.1: Default semantic annotation helper representation schema.

nominal features, and the second that uses the values of all the other features. This is intended to reflect the typical scenario where most annotations are defined by a small set of nominal features, with a few of them having features with arbitrary values. Most annotation types would then only make use of level-1 templates, with a few of them employing both level-1, and level-2 templates.

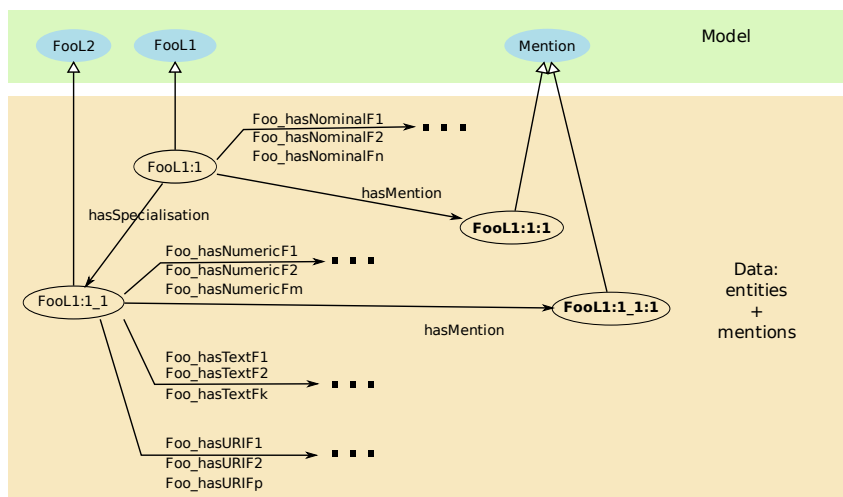


Figure 3.2: Default semantic annotation helper data example.

The representation schema used by the Sesame helper is illustrated in Figure 3.1. Figure 3.2 shows the data created for an example annotation, with the two mention URIs displayed in bold. These URIs will be stored in the mentions index. The DB helper uses a similar strategy, with separate level 1 (nominal) and level 2 (everything else) database tables for each annotation type. Annotation types that only have

nominal features need just a level 1 table.

The execution flow for each annotation includes the following main steps:

- Given the annotation's nominal features, find an appropriate level-1 (L1) template. If none is found, create one.
- For the L1 template, find a mention of appropriate length (the number of tokens covered by the annotation). If none is found, create one. Add the mention URI to the mentions index.
- If the annotation has non-nominal features:
 - Find an appropriate level-2 (L2) template, based on the feature values. If none is found, create one.
 - For the L2 template, find (or create) a mention of appropriate length; add the mention URI to the index.

All annotations sharing the same feature values will share the same database entries or knowledge base entities (the resources with URIs `'FooL1:1'` and `'FooL1:1_1'`), and the same mention objects. This has the advantage of reducing the size of the database, which allows more documents to be indexed, and helps achieve better execution speeds during search. The downside of this is that the indexing process is somewhat slowed down, as pre-existing entities need to be retrieved at every step.

Chapter 4

Searching Mímir Indexes

From a user’s point of view, Mímir is a tool for searching a collection of semantically annotated documents. It provides facilities for searching over different views of the document text, for example one can search the document words, the part-of-speech of those words, or their morphological roots. Beside searching the document text, Mímir also supports searches over the documents’ semantic annotations, where queries are based on annotation types and restrictions over the values of annotation features. These different search paradigms can be combined freely into complex queries, with support for sequences, repetitions, and Boolean operators.

A search session entails the formulation of a query, running the query with the Mímir query engine, and consuming the query results. Mímir queries are expressed in a text-based query language which is described in section 4.1. The way these queries are submitted to Mímir depends on how it is deployed, the various interfaces are discussed in section 4.2.

4.1 The Mímir query language

A Mímir query is either a simple query (i.e. a **String** query, section 4.1.1, or an **Annotation** query, section 4.1.1), or a compound one, which comprises a set of sub-queries linked by operators. If no operator is placed between any two sub-queries, then the **Sequence** operator (see section 4.1.6) is implied. This means that several queries written one after another are interpreted as one sequence query. For example, a query like *‘the brown dog’* is interpreted as a sequence query, having three sub-queries, each of them being a String query. This would match occurrences of the exact phrase *‘the brown dog’* in the indexed documents. Note that this is different from the standard behaviour of search engines, which would simply match documents in which all three query terms occur, in whichever order. That type of search is also supported in Mímir, through the **AND** operator (see section 4.1.2). Parentheses can be used for grouping where the syntax would otherwise be ambiguous.

4.1.1 String Queries

The simplest form of query is a query term. This will match all occurrences of the query term in the indexed documents.

If the Mimir index being interrogated includes multiple token indexes, then the particular index to be searched can be specified by prefixing the query term with the index name and a colon, for example the query `'root:be'`¹ will match all morphological forms of the verb *to be*. If the name of the string index is omitted, then the first configured index is used. By convention (reflected in the default Mimir configuration) the first string index is used to store the terms text, so the default behaviour is to search over the document text, as expected. Double-quoted strings are treated as plain term queries against the first token index in a similar way.

In fact the above is a slight simplification, as bare terms (and double-quoted strings) are actually tokenised before being searched for. This is because Mimir views documents as streams of tokens, not characters, and the query must match the tokenisation that was used to index the documents. For example, the default GATE tokeniser treats “don’t” as two tokens, “do” and “n’t”, so a query for *don’t* as a single token would fail. To get around this, Mimir runs a GATE application over the string of the query to generate Token annotations, and then constructs a query for these tokens in sequence (see section 4.1.6). Named index queries (“root:be”) are not tokenised, so if you want to avoid tokenising a particular query for any reason (e.g. if you suspect there is a mis-tokenised document in your index) you can explicitly name the appropriate index (typically “string”, i.e. `string:don’t`).

Annotation Queries

If annotation indexes were used during indexing, Mimir allows searching for annotation-based patterns. An annotation is a piece of metadata associated with a text segment, with a **type** and optionally **features**. An annotation query takes the following form: `{Type feature1=value1 feature2=value2 ...}`. The annotation type is required, the feature constraints are optional.

While the example above uses equality for the feature constraints, other operators are also available. Here is the full list:

equality: represented by the sign `=`, matches annotations which have the given value for the specified feature. The equality operator is applicable to features of any type.

comparison operators: represented by one of the following symbols: `<`, `<=`, `>`, `>=`, with the usual meaning. These operators can apply to features of type **nominal**, **numeric**, or **text**.

regular expressions: can be specified using the syntax `REGEX(pattern, flags)`, where the **pattern** represents the regular expression sought, and the **flags** are optional, and can be used to change the way matching is performed. See <http://www.w3.org/TR/xpath-functions/#regex-syntax> for a full specification of the regular expression support. The `REGEX` operator can only be used for **nominal**, and **text** features.

¹This assumes that an index named `root` exists, and was used to store the morphological root of the words.

Examples:

`{Person gender = female}` – searches for annotations of type **Person**, which have a feature named **gender**, with the value *female*.

`{Measurement type = scalar normalisedValue > 0 normalisedValue < 10 normalisedUnit = m}` – searches for scalar measurements, with a unit of *metre*, and a normalised value between 0 and 10². Note that the same feature name can be used in several constraints, in which case only annotations where the feature meets *all* the constraints will be matched by the query. For disjunctive queries, use the OR operator described below.

In order to be able to search on a particular feature, that feature must have been specified in the index template when the index was created – Mímir indexes only the features it has been told to index. There may be additional “synthetic” features available at query time depending on the semantic annotation helper that the index uses for the given annotation type, for example the SPARQL helper allows queries on the “feature” named “sparql”, the measurements helper allows queries for “spec” etc.

4.1.2 AND Operator: “&”

The ‘AND’ (also ‘&’) operator can be used to specify queries that should match document segments that include at least one hit from each of the sub-queries. The results returned will always be the shortest document segments that satisfy the query.

4.1.3 OR Operator: “|”

OR queries are used to search hits that match one of a set of alternative query expressions. This is indicated by using the ‘OR’ (also ‘|’) operator between the sub-queries. A query of the form **Query1 | Query2** will return hits that match either sub-query **Query1** or sub-query **Query2**.

4.1.4 IN and OVER Operators

The operators **IN** and **OVER** are used to search for hits of a query that contain, or are contained in the hits of another query. For example:

Query1 IN Query2 will match all the hits of **Query1** that are contained in a hit of **Query2**.

Query1 OVER Query2 will match all hits of **Query1** that contain (are overlapping) a hit of **Query2**.

²The extended support for Measurement annotations is discussed further in section 5.3.

4.1.5 Repeats Operator: “+”

The `+` operator can be used to match text segments that comprise a sequence of hits from the same sub-query. The length of the sequence is specified through a number (representing the **maximum** number of repetitions) or through two numeric values (representing the **minimum** and **maximum** number of repetitions). For example:

`to+3` will match one, two, or three repeated occurrences of the word *to*. The returned hits will be of the form “*to*”, “*to to*”, or “*to to to*”).

`{Person}+2..5` will match sequences of 2, 3, 4, or 5 adjacent **Person** annotations.

`({Location locType = city} | {Location locType = country})+3` will match any sequence of up to three **Location** annotations where each one refers to either a city or a country.

Note that there is no support for a repetition count of zero (an optional match) – you will need to reformulate the query to cover the versions with and without the optional match separately and combine them with an OR, for example `(term1 term2+2 term3) | (term1 term3)`. Similarly there is no support for unbounded wildcards (*n* times or more).

4.1.6 Sequence Queries and Gaps

As sequence is the default operator in Mimir, there is no graphical sign for it: simply writing a set of queries one after another will cause a search for sequences of hits, one from each sub-query. For example, the query “**the energy level**” is actually a sequence query where the first sub-query searches for the word “*the*”, the second for “*energy*”, and the last for “*level*”.

It is sometimes useful to include gaps in a sequence query, that is to allow arbitrary text fragments (of specified length) to occur in-between the hits from some of the sub-queries. This can be done by using the gap markers “[*n*]”, or “[*m*..*n*]”. These will match a sequence of length *n*, or with a length of between *m* and *n* of arbitrary tokens.

For example the query “**the [2] root:time**” will match phrases like “*the best of times*” or “*the worst of times*”, whereas the query “**the [2..10] root:time**” would also match “*the best use of one’s time*” (where the gap consists of six tokens – five words and an apostrophe).

4.1.7 Escaping Reserved Words

Some words are part of the query language definition so they cannot be used directly as query terms. If that is desired, then these constructs must be escaped as shown in the following table:

Reserved input	Escaped form
{, }	\{, \}
(,)	\(, \)
[,]	\[, \]
:	\:
+	\+
	\
&	\&
?	\?
\	\\
.	\.
"	\"
=	\=
IN	"IN"
OVER	"OVER"
OR	"OR"
AND	"AND"

Escaping reserved constructs in the Mímir query language

4.2 Search interfaces – how to submit queries to Mímir

The Mímir Grails plugin supplies two search interfaces by default, with the infrastructure to implement other interfaces as required. An XML-based service interface allows other applications to submit queries to the indexes hosted by a Mímir web application by POSTing requests over HTTP (described in section 4.2.1). There is also an example user-facing search interface called *GUS*, intended primarily for testing and demonstration purposes (described in section 4.2.2). Both of these interfaces interact with the underlying indexes through the **SearchService** Grails service provided by the plugin. When embedding the Mímir Grails plugin in another Grails application this service is the primary means for application code to interact with Mímir, and is described in section 4.2.3.

4.2.1 Mímir Search Web Service

The Mímir web application exposes the search functionality as a web service that can be accessed through a simple HTTP interface. All requests are performed by calling an action with a set of parameters; the results of a call are encoded in XML and returned as the response to the request. All the example URLs in this section assume the **mimir-demo** application with its default URL mappings, running on **localhost** port 8080.

The Mímir web service can be accessed at a URL like:

http://localhost:8080/mimir-demo/{index ID}/search/{action}, where the **action** value is the name of one of the supported actions, described below. The actual URL (with the correct index ID included) can be obtained from

the *index information page* presented by the Mimir web application. Parameters may be supplied as query parameters with a GET request or in normal `application/x-www-form-urlencoded` form in a POST request. Alternatively, they may be supplied as XML (if the request content type is `text/xml` or `application/xml`) of the form:

```
<request xmlns="http://gate.ac.uk/ns/mimir">
  <firstParam>value</firstParam>
  <secondParam>value</secondParam>
</request>
```

The first request to the service will return a session cookie, which must be passed back with all subsequent requests.

When accessing the service URL with no value provided for **action**, a help page will be returned presenting the documentation associated with the XML web service.

The following actions are available:

help

Function	Obtain service documentation.
Parameters	none
Returns	A help message describing how to use the service.

postQuery

Function	Starts a new query. The query will execute asynchronously in a background thread, and will initially run until it has found a maximum of 1 million hits or has been running for 30 seconds, whichever comes first. If the query runner hits one of these limits before the search is complete it can be restarted by calling the getMoreHits action.
Parameters	queryString : the text of the query, using the Mimir query language.
Returns	An XML message with the ID of the new query, or an error message if there were any problems while parsing the query. Example request: <code>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/postQuery?queryString=%22the%22</code> Example response: <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <queryId>a28656e2-18f4-4b58-b9d3-9a9378eb14d0</queryId> </data> </message></pre>

hitCount

Function	Obtains the number of hits collected so far. If a query is not complete, more hits may be available at later time. If a query has stopped being active before completing, it can be restarted by calling <code>getMoreHits</code> .
Parameters	queryId : the ID for the query, as returned by the <code>postQuery</code> action.
Returns	<p>An XML message encapsulating a numeric value, or an error message if there were any problems.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/hitCount?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <value>6257</value> </data> </message></pre> <p>Example error response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>ERROR</state> <error>Query ID a28656e2-18f4-4b58-b9d3-9a9378eb14d1 not known! </error> </message></pre>

docCount

Function	Obtains the number distinct documents that have been found so far to contain hits.
Parameters	queryId : the ID for the query, as returned by the <code>postQuery</code> action.
Returns	<p>An XML message encapsulating a numeric value, or an error message if there were any problems.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/docCount?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <value>11</value> </data> </message></pre>

docStats

Function	Obtains the statistics for the documents that have been found so far to contain hits. These include the document IDs and the number of hits for each individual document.
----------	---

Parameters	<p>queryId: the ID for the query, as returned by the <code>postQuery</code> action.</p> <p>startIndex the first requested document. A value of 0 requests the details for the first document found to contain hits.</p> <p>count the number of documents for which the details are requested.</p>
Returns	<p>An XML message encapsulating a set of <code><document></code> elements, one for each individual document.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/docStats?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0&startIndex=0&count=11</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <document id='0' hitCount='622' /> <document id='1' hitCount='300' /> <document id='2' hitCount='448' /> <document id='3' hitCount='273' /> <document id='4' hitCount='1053' /> <document id='5' hitCount='677' /> <document id='6' hitCount='86' /> <document id='7' hitCount='1399' /> <document id='8' hitCount='356' /> <document id='9' hitCount='841' /> <document id='10' hitCount='202' /> </data> </message></pre> <p>Note that as our query was simply <i>“the”</i>, we have found hits on every document. In a more typical case, not all document IDs will be represented in the results.</p>

hits

Function	Obtains a set of hits. Each hit is defined by a document ID, a position and a length, both of which are defined in terms of tokens, not characters (see Section 3.1 for details).
Parameters	<p>queryId the ID for the query, as returned by the postQuery action.</p> <p>startIndex the first requested hit. A value of 0 requests the details for the first hit found.</p> <p>count the number of hits for which the details are requested.</p>
Returns	<p>An XML message encapsulating a set of <hit> elements, one for each individual hit.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/hits?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0&startIndex=0&count=11</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <hits> <hit documentId='0' position='257' length='1' /> <hit documentId='0' position='266' length='1' /> <hit documentId='0' position='290' length='1' /> <hit documentId='0' position='303' length='1' /> <hit documentId='0' position='309' length='1' /> <hit documentId='0' position='316' length='1' /> <hit documentId='0' position='320' length='1' /> <hit documentId='0' position='332' length='1' /> <hit documentId='0' position='335' length='1' /> <hit documentId='0' position='342' length='1' /> <hit documentId='0' position='348' length='1' /> </hits> </data> </message></pre>

getMoreHits

Function	Requests a query that has stopped collecting hits before completing to restart collecting hits. If the query has already completed, or is already active, this call will simply be ignored (it will not cause an error).
Parameters	queryId the ID for the query, as returned by the postQuery action.
Returns	<p>An XML message reporting success or failure.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/getMoreHits?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> </message></pre>

isActive

Function	Checks if a query is still working on collecting hits.
----------	--

Parameters	queryId: the ID for the query, as returned by the <code>postQuery</code> action.
Returns	<p>An XML message encapsulating a Boolean value, or an error message if there were any problems.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/isActive?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <value>false</value> </data> </message></pre>

isComplete

Function	Checks if a query has finished collecting all the hits. If a query is not complete, more hits may be available at a later time. If a query has stopped being active before completing, it can be restarted by calling <code>getMoreHits</code> .
Parameters	queryId: the ID for the query, as returned by the <code>postQuery</code> action.
Returns	<p>An XML message encapsulating a Boolean value, or an error message if there were any problems.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/isComplete?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <value>true</value> </data> </message></pre>

renderDocument

Function	Renders the document text and hits for a given document, in the context of a given query. The HTML of the document is rendered directly to the response stream of the connection.
Parameters	<p>queryId: the ID for the query, as returned by the <code>postQuery</code> action.</p> <p>documentId the ID for the requested document, as returned by e.g. a call to the <code>docStats</code> action.</p>

Returns	<p>HTML content. The hits are rendered as <code>...</code>.</p> <p>Example request:</p> <p><code>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/renderDocument?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0&documentId=1</code></p> <p>Example response fragment:</p> <pre>... <p num="p0002">Moreover, the present invention further relates to a method of higher purification effective in the higher purification of metal which reduces the oxygen content caused by organic matter.</p> ...</pre>
---------	--

documentMetadata

Function	Returns the title and URI associated with a document. These values were provided at indexing time.
Parameters	<p>queryId: the ID for the query that has returned the document ID being used, as returned by the <code>postQuery</code> action.</p> <p>documentId the ID for the requested document, as returned by e.g. a call to the <code>docStats</code> action.</p>
Returns	<p>An XML message encapsulating the two string values, or an error message if there were any problems.</p> <p>Example request:</p> <p><code>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/documentMetadata?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0&documentId=1</code></p> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <documentTitle>EP-1288339-A9</documentTitle> <documentURI>urn:matrixware.com:alexandria:EP-1288339-A9</ documentURI> </data> </message></pre>

documentText

Function	Action for obtaining (a segment of) the text of a document.
Parameters	<p>queryId: the ID for the query that has returned the document ID being used, as returned by the <code>postQuery</code> action.</p> <p>documentId the ID for the requested document, as returned by e.g. a call to the <code>docStats</code> action.</p> <p>position the position of the first returned token. This parameter is optional; defaults to 0 is not provided, which means the first token of the document.</p> <p>length the number of tokens to be returned. This parameter is optional, if omitted, all the document tokens will be returned.</p>

Returns	<p>An XML message containing the text of all the individual tokens and, if available, the spaces between them.</p> <p>This action could be used, for example, to obtain text snippets around a query hit.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/documentText?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0&documentId=1&position=100&length=10</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> <data> <text position='100'>25</text> <text position='101'>C</text> <space> </space> <text position='102'>1</text> <text position='103'>/</text> <text position='104'>08</text> <space> </space> <text position='105'>C</text> <text position='106'>25</text> <text position='107'>C</text> <space> </space> <text position='108'>1</text> <text position='109'>/</text> </data> </message></pre>
---------	---

close

Function	<p>Closes a query, releasing all resources allocated for supporting it. After a query is closed, no more actions can be performed for it. It is important to close queries, as each running query uses up memory on the server. Queries are also closed automatically after a period of inactivity (upon session expiration, the time for which is defined in the configuration of the web application server – this is why it is important to pass the session cookie you received from <code>postQuery</code> back to the server with subsequent calls).</p>
Parameters	<p>queryId: the ID for the query, as returned by the <code>postQuery</code> action.</p>
Returns	<p>An XML message with a success or failure value.</p> <p>Example request:</p> <pre>http://localhost:8080/mimir-demo/a4300d00-2dd1-4797-8eaa-e65b0c7d879b/search/close?queryId=a28656e2-18f4-4b58-b9d3-9a9378eb14d0</pre> <p>Example response:</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>SUCCESS</state> </message></pre> <p>Example failure (using the ID for an already closed query):</p> <pre><?xml version="1.0"?> <message xmlns='http://gate.ac.uk/ns/mimir'> <state>ERROR</state> <error>Query ID a28656e2-18f4-4b58-b9d3-9a9378eb14d0 not known! </error> </message></pre>



Figure 4.1: Front page of the GUS user interface

4.2.2 GUS example user interface


The GUS search tool is a browser-based search interface intended to serve as a platform for experimentation with a Mimir index, and as a demonstration of the capabilities of the Mimir framework and API. It is written using the Google Web Toolkit, and the source code is included in the Mimir Grails plugin.

In the demo web application with its default URL mappings, the GUS interface for an index in searching mode is available at `http://localhost:8080/mimir-demo/{index ID}/gus/search`. The initial page, shown in figure 4.1, provides a text area into which you can type queries in the Mimir query language. It provides auto-completion for annotation types and features (by asking the index what types it was configured with when it was created).

Clicking the Search button starts a search on the server. The query runs asynchronously, collecting hits in the background until either 30 seconds have passed or 1 million hits have been collected, at which point it stops. To restart the search, click the “>> keep searching” link.

Hits are shown below the search box, as shown in figure 4.2, with the hit text highlighted in bold and with five tokens of left and right context. The document title is a link, in this example to the original document as the index was created with the “Document URIs are external links” option. The “cached” link shows Mimir’s cached copy of the document, with all the hits from that document highlighted in red. For indexes where the document URIs are not external links the document title would link directly to the cached version and there would be no separate “cached” link.

At the bottom of the page is a row of pagination links (figure 4.3). Since, on a large index, there can be many hundreds of thousands or even millions of hits to navigate, GUS provides links where necessary to skip over large numbers of pages in one click.



Search powered by Mimir

Searching index: tna-1.024

Results 4,351 - 4,360 of at least 1,000,000 >> keep searching

About Becta - Register of Board members' interests - Becta (cached)
 Get involved Publications © Becta **2009** About this site Freedom of

Becta's Board (cached)
 Employment Parliamentary Select Committee (**1997**) and has also been

Figure 4.2: GUS search results page

Becta's Board (cached)
 the position of Chair in **May 2009**. In December 2008 he

Becta's Board (cached)
 in May 2009. In **December 2008** he was also appointed by

Becta's Board (cached)
 Safeguarding Children Board. In **January 2008** he was awarded the CBE

Becta's Board (cached)
 Mail Group plc) in **1985** and had Board responsibility,

Becta's Board (cached)
 He became Chief Executive in **1995** and saw the organisation through

Page [< Prev](#) [431](#) [432](#) [433](#) [434](#) [435](#) [436](#) [437](#) [438](#) [439](#) [440](#) [Next >](#)

Skip pages [<100](#) [<10](#) [10>](#) [100>](#) [1,000>](#) [10,000>](#)

Figure 4.3: GUS pagination links for a large search

4.2.3 Embedding Mímir in a Grails application

Both the XML web service and the GUS interface ultimately use a Grails service provided by the Mímir plugin to search their indexes. If you install the Mímir plugin into your own Grails application this service will be your primary entry point to make use of Mímir functionality, so this section explains what you need to know to use it effectively.

The `searchService` is a normal Grails service which can be autowired into your own services, controllers, etc. The service itself is very simple, offering only the following methods:

postQuery(index, queryString) start running a query against the given index. The index can be specified either as a string containing the indexId (the last component of the index URL, typically a UUID) or as an `Index` domain object (the database object representing a local, remote or federated index). Returns a *query ID* string.

getQueryRunner(queryId) retrieves the `QueryRunner` for the given running query ID. `QueryRunner` is the interface through which you can interact with the running query.

closeQueryRunner(queryId) indicates that the given query runner is no longer required. It is important to call this method when you have finished with a query runner, as each runner owns resources such as background threads which need to be properly cleaned up.

Once a query has been started, its `QueryRunner` provides access to the statistics, the hits themselves, and the text in the matched documents. The most important methods are summarised below, but for full details you should look at the interface definition itself, in the `gate.mimir.search` package of `mimir-core`. Note that the search itself is performed in a background thread so many of the methods of `QueryRunner` will return different values over time as the search progresses.

isComplete() Checks whether the search has completely finished and there are definitely no more hits to be found.

isActive() Checks whether the runner is currently actively searching. By default a query runner keeps searching until it has either exhausted all the possible hits, has found a million hits since it was last restarted, or has been running for 30 seconds without hitting this limit. If `isActive()` and `isComplete()` both return false, it means that the search hit one of these limits and was suspended, it can be restarted by calling `getMoreHits()`.

getHitsCount() Gets the number of hits obtained so far. This number may increase at any time if the query is currently active.

getDocumentsCount() Gets the number of distinct documents that have so far been found to contain hits. This number may increase at any time if the query is currently active.

getHits(start, max) Gets the details of some hits found by the query. Conceptually, a `QueryRunner` can be thought of as holding a flat list of hits numbered from zero upwards, containing all the hits from the first matched document, followed by all the hits from the second matched document, etc. This method retrieves a sub-list from that list, starting at the (zero-based) index *start* and containing a maximum of *max* entries (it may contain fewer if not enough hits

have yet been found). The return value from this method is a list of **Binding** objects, each representing one hit.

getDocumentHitsCount(documentIndex) Gets the number of hits in the *n*th document that matched this query (zero-based index). To retrieve the hits for a particular document you would need to sum up all the `getDocumentHitsCount` values for the preceding documents and pass that sum as the *start* parameter to `getHits`.

getDocumentID(documentIndex) Gets the ID in the underlying index of the *n*th document that matched this query. This ID is needed to get the document text and metadata.

getDocumentTitle/URI(id) Gets metadata about the document with the given ID.

getDocumentText(id, start, length) Gets the text of the document with the given ID, starting at the *start*th token and extending for *length* tokens. The return value is a pair of parallel string arrays, one containing the text of the tokens and the other containing the text between each token and the following one.

renderDocument(id, Appendable) Render the document content, with hits highlighted, using the document renderer configured for the index. The content is written to the specified Appendable (a `StringBuilder`, `Writer`, etc.).

The `getHits` method returns a list of **Binding** objects, which provide several methods, the most important ones being **getDocumentId** (the document that contains the hit), **getTermPosition** (the offset of the first token covered by the hit) and **getLength** (the number of tokens it covers).

Chapter 5

Standard Mimir Plugins

Mimir uses the GATE Embedded *CREOLE plugin* mechanism to load semantic annotation helper classes. A number of plugins are supplied by default with the Mimir distribution, and those plugins are described in this chapter. Information on how to create new plugins to provide user-defined helper classes can be found in section 6.1.

5.1 The db-h2 plugin

The **db-h2** plugin is one of two plugins (the other being the **sesame** plugin) that provides a *generic* semantic annotation helper implementation that can be configured for any annotation type with any features. The helper provided by **db-h2** uses an embedded relational database engine (<http://www.h2database.com/>) to store the annotation data, and generally provides the best performance of the standard generic helpers.

`gate.mimir.db.DBSemanticAnnotationHelper` is the helper class provided by the **db-h2** plugin. It has a constructor that takes a `Map` of configuration parameters, and Groovy provides special “named argument” support for `Map`-valued method and constructor parameters, allowing the following idiom in the index template DSL:

```
1 // note the “import X as Y”, which is another Groovy feature to create an
2 // alias for an imported class name
3 import gate.mimir.db.DBSemanticAnnotationHelper as DefaultHelper
4 // ...
5 semanticAnnotations = {
6     index {
7         annotation helper:new DefaultHelper(annType:'Person',
8                                             nominalFeatures:["gender"])
9     }
10 }
```

The supported constructor arguments are:

annType: the annotation type which the helper is to process.

nominalFeatures: the names of the features to be indexed that have nominal values. An annotation feature is said to be nominal if the range of possible values is clearly defined and limited in size. There is no hard rule regarding the size of the set of permitted values, but, for optimal results, this should not exceed a few tens of values.

integerFeatures: the names of the features to be indexed that have integer values (i.e. values that can be converted to a Java `long` value).

floatFeatures: the names of the features to be indexed that have floating-point numeric values (i.e. values that can be converted to a Java `double` value).

textFeatures: the names of the features to be indexed that have arbitrary text values (as opposed to the nominal case of a fixed list of possible values).

uriFeatures: the names of the features to be indexed that have URIs as values.

The DB-based helper does not distinguish between text- and URI-valued features, indexing both types in the same way, but it accepts both kinds as arguments to be consistent with the Sesame helper (described below). When using the `import as` mechanism, switching an index template from the DB helper to the Sesame helper or vice versa is simply a matter of changing the name of the class being imported.

5.2 The sesame plugin

The `sesame` plugin provides the other standard generic helper, `gate.mimir.sesame.SesameSemanticAnnotationHelper`. This helper stores annotation data as triples in a Sesame repository (by default using OWLIM as the underlying storage engine).

The `SesameSemanticAnnotationHelper` constructor takes the same named arguments as the DB helper for the annotation type and feature names (see the previous section for details). Internally all numeric features (integer and float) are stored as floating point numbers, but the helper accepts both “integerFeatures” and “floatFeatures” parameters for consistency with the DB helper.

For advanced users it also provides a mechanism to specify the location of the Sesame repository configuration template file (in the standard Turtle RDF format), as a parameter named either “absolutePath” (for an absolute path to the file) or “relativePath” (for a path relative to the `sesame` plugin directory). The default configuration file is `resources/owlim.ttl` under the `sesame` plugin. When specifying a custom repository configuration it is important to retain the repositoryID of “owlim” as this is assumed by the helper code.

5.2.1 Searching in the knowledge base

For annotations that have a URI-valued feature named “inst”, the Sesame helper provides an additional “synthetic” feature named `semanticConstraint` that provides a way to query for annotations based on information in the knowledge base.

```
1 {Organization semanticConstraint = "?inst <http://proton.
   semanticweb.org/2005/04/protont#locatedIn> <http://example.com#
   Sheffield> ."}}
```

The value of the `semanticConstraint` feature is a SPARQL fragment that has access to the `?inst` variable, referring to the URI in the “inst” feature of the annotation. This mechanism can only query for one specific feature (inst) and only against triples stored in the same semantic repository as the Mimir annotation data. For a more flexible approach to semantic queries that can make use of semantic data stored in a remote knowledge base, see the `sparql` plugin (section 5.4).

5.3 The measurements plugin

The GATE `Tagger Measurements` plugin, introduced in GATE 6.1, is able to recognise many different kinds of measurement expressions in text. It normalizes the value and unit of each measurement into the SI system of measurements and stores these values as features of the Measurement annotation. For example, the text “45 cm” would be annotated with a normalized unit of metres and a normalized value of 0.45, the text “18 in” would also be normalized to metres, in this case with a normalized value of 0.4572.

The Mimir `measurements` plugin provides a SAH that implements the same normalization on queries. It processes queries for a “synthetic” feature called “spec” which represents a measurement specification in a controlled language and converts constraints on this feature into the corresponding constraints on the real normalized value and unit features that have been indexed. For example, a search for `{Measurement spec="1 to 3 feet"}` would be treated as a query for measurements whose normalized unit is metres and whose normalized value is between 0.3048 and 0.9144, which would match both the “45 cm” and “18 in” examples above.

5.3.1 Configuring the Measurements SAH

To use the measurements helper you need to first ensure that the `measurements` plugin is loaded into your Mimir instance, then create an index template that specifies an instance of the helper:

```

1  import gate.mimir.measurements.MeasurementAnnotationHelper
2
3  // ...
4  semanticAnnotations = {
5    index {
6      // Measurement helper with default settings
7      annotation helper:new MeasurementAnnotationHelper(
8        delegateHelperType:DefaultHelper)
9    }
10 }
```

Note that the measurement helper does not need any “annType” or “xxxFeatures” parameters, as it is hard-coded to work only for annotations of type “Measurement” with the features that are produced by the measurement tagger PR. However the constructor does take a `Map` with other named arguments:

```

6  // Example of how to configure a custom “units” file
7  annotation helper:new MeasurementAnnotationHelper(
```

```
8         delegateHelperType:DefaultHelper ,
9         unitsFile:'resources/americanUnits.dat',
10        locale:'en_US')
```

The following parameters are supported:

delegateHelperType (required) a `Class` object representing the type of generic helper that the Measurements helper should delegate to. This class must provide a 6-argument constructor taking the annotation type (a `String`) and five `String` arrays for the nominal, integer, float, text and URI feature names respectively. Both the DB and Sesame helpers provide this constructor.

unitsFile the location of the `units.dat` file used to configure the measurements parser. If not specified, a default file provided with the `measurements` plugin is used. This value can be an absolute URL (`file:/path/to/units.dat`) or a relative path which will be resolved against the `measurements` plugin directory.

commonWords the location of the common words file used by the measurements parser. As with the `unitsFile` parameter, if omitted a default file bundled with the plugin is used.

locale the locale under which the measurements will be parsed. Defaults to “en_GB” if unspecified.

encoding the character encoding used to read the configuration files. Defaults to “UTF-8” if unspecified.

The measurements SAH is pre-configured with the feature names that the measurements tagger produces, and can only be used for annotations of type `Measurement`.

Measurements helper implementation

The `MeasurementAnnotationHelper` extends the `DelegatingSemanticAnnotationHelper` base class described above. It does not add any behaviour at indexing time, simply passing all the annotations through directly to its delegate. However it overrides the `getMentions` search method to support the “spec” feature.

When a query including a spec feature constraint is received, the helper parses this spec using the measurements parser to obtain a normalized unit and value or values for the measurement sought. It then constructs a number of new constraint sets that match annotations compatible with the spec and then for each of these alternatives, runs these constraints in combination with the other non-spec constraints of the original query against the delegate helper. The final set of URIs returned is the union of the results obtained from the delegate for all the alternative reformulations of the spec constraint.

As well as being useful in its own right for Measurement annotations, the measurements helper serves as an example of how to implement your own special-purpose helper based on the delegating base class. Feel free to use it as a template for your own helper implementations.

5.4 The sparql plugin

The `sparql` plugin provides a semantic annotation helper that wraps another helper, adding flexible *semantic query* support. It is intended to be used with annotations that have one or more URI-valued features whose values refer to entities in an external knowledge base (accessible at a standard *SPARQL endpoint*). The SPARQL helper has no effect at indexing time, simply delegating all calls through to the underlying helper, but at search time it allows queries for the synthetic feature “`sparql`”. This feature value is taken to be a SPARQL “SELECT” query, which is posted to the configured SPARQL endpoint. The variables in the SELECT query must correspond to the names of features that have been indexed by the underlying helper, and each row in the result set becomes a standard Mimir query to the underlying helper. Any annotations that match any of these new queries will be treated as a match for the `sparql` constraint. This process is described in detail below.

For example, given a helper configured for the public DBpedia endpoint `http://dbpedia.org/sparql`, the following Mimir query:

```
{Person sparql = "SELECT DISTINCT ?inst WHERE {
  ?inst <http://dbpedia.org/ontology/birthPlace>
    <http://dbpedia.org/resource/Sheffield> }"}}
```

would search for all Person annotations that have an “`inst`” feature containing the URI of an entity in DBpedia that represents a person born in Sheffield.

5.4.1 Creating a SPARQL helper

The SPARQL semantic annotation helper class is called `gate.mimir.sparql.SPARQLSemanticAnnotationHelper`. It has a `Map` constructor taking the following parameters:

- delegate (required)** the underlying semantic annotation helper that this SPARQL helper should wrap.
- sparqlEndpoint (required)** the address of the SPARQL endpoint that this helper should use when making SPARQL queries.
- queryPrefix** an optional prefix that will be prepended to the string specified in the `sparql` synthetic feature to form the actual SPARQL query that will be sent to the endpoint. Typically this would be used to define appropriate namespace prefixes.
- querySuffix** an optional suffix that will be appended to the end of the SPARQL queries. Thus the actual SPARQL query submitted to the endpoint is `queryPrefix + sparql feature + querySuffix`.
- sparqlEndpointUser and sparqlEndpointPassword** username and password used to authenticate to the SPARQL endpoint (only HTTP basic authentication is supported). May be omitted if your endpoint does not require authentication.

The helper also accepts the usual “`annType`” and “`xxxFeatures`” parameters but these are not normally required – if the delegate helper is a subclass of

`AbstractSemanticAnnotationHelper` (which is the case for all the standard helpers) then the SPARQL helper will take its feature names from the delegate, so the only time the features need to be specified explicitly for the SPARQL helper is if the delegate is a custom helper type that does not extend `AbstractSemanticAnnotationHelper`.

For example, the following configuration would set up a helper for Person annotations operating against DBpedia, to support the example query above:

```

1  import gate.mimir.db.DBSemanticAnnotationHelper as DBH
2  import gate.mimir.sparql.SPARQLSemanticAnnotationHelper as
                                     SPARQLHelper
3  // ...
4  semanticAnnotations = {
5      index {
6          annotation helper:new SPARQLHelper(
7              sparqlEndpoint:'http://dbpedia.org/sparql',
8              delegate:new DBH(annType:"Person", uriFeatures:['inst']))
9      }
10 }
```

Alternatively, the helper could be configured with a `queryPrefix` to set up some useful namespace prefixes:

```

1  semanticAnnotations = {
2      index {
3          annotation helper:new SPARQLHelper(
4              sparqlEndpoint:'http://dbpedia.org/sparql',
5              queryPrefix:
6              'PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#> \
7              PREFIX xsd:<http://www.w3.org/2001/XMLSchema#> \
8              PREFIX owl:<http://www.w3.org/2002/07/owl#> \
9              PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#> \
10             PREFIX dbo:<http://dbpedia.org/ontology/> \
11             PREFIX dbr:<http://dbpedia.org/resource/> ',
12              delegate:new DBH(annType:"Person", uriFeatures:['inst']))
13      }
14 }
```

Note the backslashes which are a Groovy feature to permit a String literal to be broken across several lines, and also the trailing space before the closing quotation mark – the helper simply concatenates the prefix, query and suffix without any additional space, so required spaces must be part of the prefix or suffix string itself. This would allow the example query above to be rewritten more succinctly as:

```
{Person sparql = "SELECT DISTINCT ?inst WHERE {
  ?inst dbo:birthPlace dbr:Sheffield }"}
```

For annotation types that have only one URI-valued feature it may be desirable to include the “`SELECT DISTINCT ?inst WHERE {` ” in the prefix and add a suffix of “`}`”, which would reduce the query down to

```
{Person sparql = "?inst dbo:birthPlace dbr:Sheffield"}
```

If your index template includes several ontology-based annotation types sharing the same SPARQL endpoint and prefixes then listing these in full for each annotation type will result in a large and unwieldy template. However, since the index template is itself a Groovy script it is possible to declare methods to factor out the common code. Method declarations must be placed *outside* the `semanticAnnotations` block:

```

1  import gate.mimir.db.DBSemanticAnnotationHelper as DBH
2  import gate.mimir.sparql.SPARQLSemanticAnnotationHelper as
                                   SPARQLHelper
3
4  def standardHelper(type) {
5      return new SPARQLHelper(
6          sparqlEndpoint: 'http://dbpedia.org/sparql',
7          queryPrefix:
8              'PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#> \
9              PREFIX xsd:<http://www.w3.org/2001/XMLSchema#> \
10             PREFIX owl:<http://www.w3.org/2002/07/owl#> \
11             PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#> \
12             PREFIX dbo:<http://dbpedia.org/ontology/> \
13             PREFIX dbr:<http://dbpedia.org/resource/> ',
14          delegate: new DBH(annType: type, uriFeatures: ['inst']))
15  }
16
17  // ...
18  semanticAnnotations = {
19      index {
20          annotation helper: standardHelper('Person')
21          annotation helper: standardHelper('Location')
22          annotation helper: standardHelper('Organization')
23      }
24  }

```

5.4.2 Format of SPARQL queries

This section describes in more detail how the SPARQL queries relate to the annotations indexed by the underlying semantic annotation helper. As a simple example we consider the query for people born in Sheffield:

```
{Person sparql = "SELECT DISTINCT ?inst WHERE {
    ?inst dbo:birthPlace dbr:Sheffield }"}
```

The helper will submit the SPARQL query to its configured endpoint, and receive a response of the form:

inst
http://dbpedia.org/resource/Gordon_Banks
http://dbpedia.org/resource/Michael_Palin
http://dbpedia.org/resource/David_Blunkett
...

This will then generate a series of queries to the underlying helper of the form:

```
{Person inst = "http://dbpedia.org/resource/Gordon_Banks"}
```

```
{Person inst = "http://dbpedia.org/resource/Michael_Palin"}  
...
```

and any annotation that matches any of these queries will be returned as a match for the `sparql` constraint.

The SPARQL query can bind more than one variable, and the values of the variable bindings can be RDF literals as well as URIs, they convert to queries on the underlying helper in the same way.

Chapter 6

Extending and Customising Mimir

The standard semantic annotation helpers provided by Mimir are adequate for many use cases, but if your application needs more functionality that they cannot provide it is easy to add your own custom helper implementations using a plugin mechanism. This process is described in section 6.1. Similarly, the basic Mimir demo application shows the simplest way to use the Mimir Grails plugin, but it provides no authentication or security, for example. To add these kinds of features you can install the Mimir plugin into your own custom Grails application, as described in section 6.2.

6.1 Creating new Semantic Annotation Helpers

Semantic annotation helpers (SAHs) are the mechanism that Mimir uses to store information about annotations and allow this information to be queried at search time. A SAH is associated with a particular annotation type in the Mimir index configuration, and performs two functions:

During indexing for each annotation of the relevant type, store information about that annotation in some persistent form and return to Mimir one or more URIs that represent that annotation. These URIs are included in the main MG4J index and associated with the location in the document where the annotation was found.

During searching given a set of feature value constraints, use the persistent store created during indexing to determine the URIs associated with annotations that satisfy the constraints.

Conceptually, SAH implementations can be divided into two types. *Generic* helpers are those that can index any annotation types and features, and *special-purpose* helpers are those that are designed to work with specific types of annotation. There are two generic SAH implementations provided with Mimir by default. You would create a new generic SAH implementation if you wanted to store annotation data in

a different underlying storage format. Mimir provides one example special-purpose SAH for Measurement annotations, which can serve as a template for how to implement your own helpers for other annotation types.

6.1.1 The `SemanticAnnotationHelper` interface

The `gate.mimir.SemanticAnnotationHelper` interface is the contract that all helpers must implement. It specifies three groups of methods that must be implemented:

Lifecycle methods

The interface includes two pairs of `init/close` lifecycle methods, one pair taking an `Indexer` parameter (used when the helper is indexing annotations) and the other pair taking a `QueryEngine` parameter (used when the helper is searching). Both the `Indexer` and `QueryEngine` provide access to an `IndexConfig` object which defines the configuration of the index, including the location of the index files on disk, and provides a mutable “context” map that can be used to share objects among the various SAH objects (for example the Sesame helper uses the context to share a single connection to the semantic repository among all the helpers associated with the index). The appropriate `init` method is called by Mimir when the index is opened (in whichever mode), before any other requests are passed to the helper, and the corresponding `close` method is called when the index is shut down.

Indexing methods

When indexing annotations Mimir calls the following methods:

`documentStart(document)` Called when the indexer starts processing a particular document to allow the helper to perform any per-document setup tasks. This method is guaranteed to be called once per document, before any calls to `getMentionUris`.

`getMentionUris(annotation, length, indexer)` Called once for each semantic annotation of this helper's type in the document. The helper is expected to use the annotation's length (in tokens) and feature values to determine the relevant URI or URIs that represent this annotation, and return them.

`documentEnd()` Called after all the annotation for a particular document have been processed.

These methods are always called from a single thread, as long as the same helper object is not used for more than one annotation type.

Note that the annotation length passed to `getMentionUris` is measured in *tokens*, not characters. Because Mimir operates on streams of tokens, semantic annotations that partially overlap a token will be considered by Mimir to cover the whole token. I.e. given the hypothetical example:

```
... started on 10/05/1987 by John Smith ...
-----
```

where tokens are represented by ---, an annotation that covers just the “87” would be indexed as if it covered the whole “1987” token.

Search methods

The final method in the interface is `getMentions(annotationType, constraints, queryEngine)`. This method is called by Mímir when searching for annotations, and the helper must use its stored data to determine all the possible mention URIs that satisfy the provided constraints, and return them along with their lengths (in tokens) as provided to `getMentionUris` when the annotations were indexed.

There is a second overloading of this method specified in the interface which is a convenience for callers when all the constraints are simple feature value equality constraints, but generally implementors of new SAH types can ignore this as Mímir provides an abstract base class that converts the Map form of constraints into the more general List<Constraint> form and calls the other method.

The `getMentions()` methods may be called from multiple threads at the same time, so implementations should be thread-safe.

6.1.2 Abstract base classes

Mímir provides an abstract class `AbstractSemanticAnnotationHelper` which, as described above, implements the Map version of `getMentions` in terms of the List<Constraint> version, and also provides empty implementations of `documentStart` and `documentEnd`. As well as this, it provides accessor methods to access the list of feature names of each of the five types (nominal, integer, float, text and URI) that a particular helper object supports. `AbstractSemanticAnnotationHelper` enjoys special support in the Mímir Grails plugin, allowing clients to determine what feature names an index supports for each annotation type whose helper extends `AbstractSemanticAnnotationHelper`. All the standard helper implementations provided with Mímir extend this base class.

Special-purpose helpers for particular annotation types typically operate by mapping the features of their target annotations and/or the feature constraints in a query into a different set of features or constraints which can then be handled by a generic helper. The Measurement helper described in section 5.3 operates in this way. To support this pattern, Mímir provides an abstract `DelegatingSemanticAnnotationHelper` which implements all the SAH interface methods to simply delegate to another helper instance. Subclasses can then override the methods as appropriate to map their features or constraints into terms that the underlying helper can understand and then call the `super` method to pass these parameters on to the delegate.

`DelegatingSemanticAnnotationHelper` extends `AbstractSemanticAnnotationHelper` so it advertises the features it supports in the usual way. However it is important to note that the various `get*FeatureNames` methods of the delegating helper do *not* call their counterparts in the delegate, which allows a delegating helper to advertise different features from those supported by its delegate.

The SAH lifecycle

Semantic annotation helper objects go through a specific lifecycle in Mímir. When creating a new index for indexing, the helpers defined for each semantic annotation type are instantiated by calling their constructors from the Groovy DSL (see the measurements example below). Once instantiated, the `init(Indexer)` methods of each helper in turn are called (one after the other, in a single thread, so if you are sharing objects among your helpers through the context you can be sure that you have exclusive access to the context map during the call to your init method).

The actual indexing process takes place in several threads in a pipelined manner. When a document arrives for indexing it is first processed by the token indexer (to index the token features), then the semantic indexers specified by the `index { ... }` blocks in the DSL in turn. Each indexer operates in its own thread, with documents passing from one to the next via queues. So each document is only processed by one thread at a time but under load you may have the token indexer dealing with document 3 at the same time as semantic indexers are dealing with documents 2 and 1.

When indexing is complete the helpers' `close(Indexer)` methods are called (again, in sequence in one thread). The index is now closed and the SAH objects can be garbage collected.

The index configuration, including all the SAH objects, is serialized to XML using XStream (<http://xstream.codehaus.org>). Therefore it is important to mark as `transient` any fields of your helper class that should not be serialized (e.g. temporary in-memory caches, etc.).

When an index is opened for searching the XML configuration is deserialized to recreate the helper objects, and their `init(QueryEngine)` methods are called. Note that as with Java object serialization XStream does *not* call object constructors when deserializing, so any initialization must happen in the init method or in a `readResolve` method, and not in the constructor.

Annotation queries result in calls to the relevant helper's `getMentions` method, which has been discussed in detail above.

Finally, when the index is shut down the `close(QueryEngine)` methods of the helpers are called in sequence.

6.1.3 Packaging new helper types for use with Mímir

To create a new helper type and make it available to Mímir you need to first write your Java class that implements `SemanticAnnotationHelper`, typically via one of the abstract base classes. To make your helper available to Mímir requires a couple of extra steps.

Helper implementations are loaded into Mímir using the standard GATE CREOLE plugin mechanism. The helper class or classes must be packaged up into a JAR file which is placed in a directory containing a simple `creole.xml` file.

```
2 <JAR>my-plugin.jar</JAR>
3 </CREOLE-DIRECTORY>
```

This directory is now a CREOLE plugin that can be loaded into Mímir, for example by using the configuration options described in Section 2.2.3.

6.2 Embedding Mímir in a custom Grails application

The demo Grails application provided with Mímir provides access to the functionality of the Mímir Grails plugin but is deliberately kept simple. It does not address concerns such as security and authentication, but these can be handled by installing the Mímir plugin into a custom Grails application and adding these functions there.

6.2.1 Installing the Mímir plugin

The Mímir plugin is installed like any other Grails plugin, by running `grails package-plugin` in the `grails-plugin-mimir` directory to create a plugin zip file and then `grails install-plugin grails-mimir-web-{version}.zip` in your application's base directory to install the plugin into your application. On installation the plugin adds a number of files to the host application:

grails-app/conf/MimirConfig.groovy a default configuration file for Mímir settings, which can be customised as described in section 2.2.3.

grails-app/conf/MimirUrlMappings.groovy a default set of URL mappings for Mímir controllers. These are discussed in more detail below.

web-app/WEB-INF/gate-home a default GATE home directory, containing user and site configuration files and the default query tokeniser application, as required by the template configuration file.

6.2.2 URL mappings

The standard URL mappings for Mímir adhere to the following conventions:

- All administrative actions have URLs that start with `/admin`. This is intended to simplify the job of providing security via an external mechanism, such as the default container-provided security defined by the servlet specification.
- The URLs for all actions referring to a given index have a common prefix (by default `/indexId/`). This is used to support remote access through the `mimir-client` library, which relies on a base URL to which it adds known suffixes to access given functionality (e.g. the search actions are all found under `/indexId/search/{actionName}`). The `mimir-client` library is used, for example, in the implementation of remote indexes.

These default mappings may be modified, but the protocols used by client code to push documents into an index, to search using the XML API, and to interact with

remote indexes all make assumptions about certain relative paths. Therefore, if you need to allow remote access via the `mimir-client` library, then you must make sure that:

- there is a URL prefix that is common to all actions provided by the `IndexManagementController` and the `SearchController`. Let us call this common prefix the `{indexUrl}`;
- all actions provided by the `IndexManagementController` should be published under `{indexUrl}/manage/{actionName}`;
- all actions provided by the `SearchController` should be published under `{indexUrl}/search/{actionName}`;
- all actions of the search controller must have `parseRequest:true`;
- the names of the actions in the controllers above must not be changed.

6.2.3 Security considerations

There are a number of established and mature Grails plugins that implement user registration, authentication and security in various ways for Grails applications. By design, the Mímir plugin and demo app do not provide any security of their own, but for production deployments you may wish to consider any or all of the following.

- Place your application behind an Apache HTTPD or other similar front-end server, and proxy only those URLs that need to be public (`/css`, `/gwt`, `/images`, `/plugins`, `/yourindex}/search`, etc.)
- Install a security plugin such as `spring-security-core` or `nimble`, and restrict the Mímir pages to certain users. For example the `/admin` pages could be restricted to administrative users, the `gus` pages to registered searchers, etc.
- If your Mímir instance is a slave in a federated index, use firewall rules or a Grails filter to restrict access to the `*/search` URLs to only the master Mímir server.

Note that if you use user authentication on the XML search API or the remote protocol then you will need to modify the default client code to authenticate itself appropriately. For example if you require HTTP basic authentication for the search or manage URLs then you will need to configure the master Mímir server that uses this as a remote to make it pass the relevant `Authorization` HTTP header. The remote protocol uses the `WebUtils` class in `mimir-client` to make its HTTP calls, and there is a subclass of this class that is able to provide basic authentication headers, you can wire this in by overriding the definition of the `webUtilsManager` Spring bean supplied by the Mímir plugin – see the plugin source code for full details.

Appendix A

Change Log

This appendix details the main changes in each Mimir release.

A.1 Version 3.3.0 (October 2011)

- Added support for marking documents as “deleted” (see section 2.4).
- Major changes to the format of the Index Template Groovy DSL (see section 3.1). The old format provided by Mimir 3.2.0 is still supported for existing semantic annotation helper types, but new helper types in future may not be supported in the old style DSL.
- Added the *SPARQL* semantic annotation helper (see section 5.4).
- Updated versions of a number of libraries (H2 database to 1.3.160, OWLIM to 3.5, MG4J to 4.0, fastutil to 6.4, dsiutils to 2.0).
- The `ordi` semantic annotation helper plugin is now deprecated. Use the `sesame` plugin instead, which supports the same on-disk format for its annotation storage but uses a different library to access it.
- Fixed various bugs and memory leaks (see subversion logs for full details).

A.2 Version 3.2.0 (May 2011)

First public release of Mimir.