



REFERENCE GUIDE

For Solr 4.1



Table of Contents

About This Guide	7
Getting Started	9
Installing Solr	10
Running Solr	12
A Quick Overview	18
A Step Closer	21
Upgrading to Solr 4.1	22
Using the Solr Administration User Interface	23
Overview of the Solr Admin UI	24
Getting Assistance	26
Logging	27
Cloud Screens	29
Core Admin	31
Java Properties	32
Thread Dump	33
Core-Specific Tools	35
Documents, Fields, and Schema Design	49
Overview of Documents, Fields, and Schema Design	50
Solr Field Types	52
Defining Fields	66
Copying Fields	68
Dynamic Fields	69
Other Schema Elements	70
Putting the Pieces Together	73
Understanding Analyzers, Tokenizers, and Filters	75
Overview of Analyzers, Tokenizers, and Filters	76
What Is An Analyzer?	77
What Is A Tokenizer?	80
What Is a Filter?	83
Tokenizers	85
Filter Descriptions	95
CharFilterFactories	126
Language Analysis	130
Phonetic Matching	158
Running Your Analyzer	159
Indexing and Basic Data Operations	164
What Is Indexing?	165
Uploading Data with Solr Cell using Apache Tika	167

Uploading Data with Index Handlers	177
Uploading Structured Data Store Data with the Data Import Handler	193
De-Duplication	215
Detecting Languages During Indexing	218
Content Streams	222
UIMA Integration	224
Searching	227
Overview of Searching in Solr	229
Relevance	233
Query Syntax and Parsing	235
Highlighting	275
MoreLikeThis	280
Faceting	283
Result Grouping	298
Spell Checking	306
Suggester	316
Function Queries	322
Spatial Search	340
The Terms Component	345
The Term Vector Component	350
The Stats Component	353
The Query Elevation Component	357
Response Writers	360
Near Real Time Searching	366
RealTime Get	369
The Well-Configured Solr Instance	371
Configuring solrconfig.xml	372
Core Admin and Configuring solr.xml	398
Solr Plugins	406
JVM Settings	407
Managing Solr	409
Running Solr on Tomcat	410
Running Solr on Jetty	413
Configuring Logging	414
Backing Up	418
Using JMX with Solr	420
SolrCloud	423
Getting Started with SolrCloud	425
How SolrCloud Works	434
SolrCloud Configuration and Parameters	444
SolrCloud Glossary	459
Legacy Scaling and Distribution	460

Introduction to Scaling and Distribution	461
Distributed Search with Index Sharding	462
Index Replication	469
Combining Distribution and Replication	484
Merging Indexes	486
Client APIs	488
Introduction to Client APIs	489
Choosing an Output Format	490
Using JavaScript	491
Using Python	492
Client API Lineup	494
Using SolrJ	495
Using Solr From Ruby	500
MBean Request Handler	502
Further Assistance	503
Solr Glossary	505
B	505
C	505
D	506
E	506
F	506
I	506
L	507
M	507
N	507
O	507
Q	507
R	508
S	508
T	509
W	510
Z	510

Apache Solr Reference Guide

This reference guide describes Apache Solr, an open source solution for search. You can download Apache Solr from the Solr website at <http://lucene.apache.org/solr/>.

This Guide contains the following sections:

Getting Started: This section guides you through the installation and setup of Solr.

Using the Solr Administration User Interface: This section introduces the Solr Web-based user interface. From your browser you can view configuration files, submit queries, view logfile settings and Java environment settings, and monitor and control distributed configurations.

Documents, Fields, and Schema Design: This section describes how Solr organizes its data for indexing. It explains how a Solr schema defines the fields and field types which Solr uses to organize data within the document files it indexes.

Understanding Analyzers, Tokenizers, and Filters: This section explains how Solr prepares text for indexing and searching. Analyzers parse text and produce a stream of tokens, lexical units used for indexing and searching. Tokenizers break field data down into tokens. Filters perform other transformational or selective work on token streams.

Indexing and Basic Data Operations: This section describes the indexing process and basic index operations, such as commit, optimize, and rollback.

Searching: This section presents an overview of the search process in Solr. It describes the main components used in searches, including request handlers, query parsers, and response writers. It lists the query parameters that can be passed to Solr, and it describes features such as boosting and faceting, which can be used to fine-tune search results.

The Well-Configured Solr Instance: This section discusses performance tuning for Solr. It begins with an overview of the `solrconfig.xml` file, then tells you how to configure cores with `solr.xml`, how to configure the Lucene index writer, and more.

Managing Solr: This section discusses important topics for running and monitoring Solr. It describes running Solr in the Apache Tomcat servlet runner and Web server. Other topics include how to back up a Solr instance, and how to run Solr with Java Management Extensions (JMX).

SolrCloud: This section describes the newest and most exciting of Solr's new features, SolrCloud, which provides comprehensive distributed capabilities.

Legacy Scaling and Distribution: This section tells you how to grow a Solr distribution by dividing a large index into sections called shards, which are then distributed across multiple servers, or by replicating a single index across multiple services.

[Client APIs](#): This section tells you how to access Solr through various client APIs, including JavaScript, JSON, and Ruby.

About This Guide

This guide describes all of the important features and functions of Apache Solr. It is free to download from <http://lucene.apache.org/solr/>.

Designed to provide high-level documentation, this guide is intended to be more encyclopedic and less of a cookbook. It is structured to address a broad spectrum of needs, ranging from new developers getting started to well-experienced developers extending their application or troubleshooting. It will be of use at any point in the application life cycle, for whenever you need authoritative information about Solr.

The material as presented assumes that you are familiar with some basic search concepts and that you can read XML. It does not assume that you are a Java programmer, although knowledge of Java is helpful when working directly with Lucene or when developing custom extensions to a Lucene/Solr installation.

Special notes are included throughout these pages.

Note Type	Look & Description
Information	<p> Notes with a blue background are used for information that is important for you to know.</p>
Notes	<p> Yellow notes are further clarifications of important points to keep in mind while using Solr.</p>
Tip	<p> Notes with a green background are Helpful Tips.</p>
Warning	<p> Notes with a red background are warning messages.</p>

-  The default port configured for Solr during the install process is 8983. The samples, URLs and screenshots in this guide may show different ports, because the port number that Solr uses is configurable. If you have not customized your installation of Solr, please make sure that you use port 8983 when following the examples, or configure your own installation to use the port numbers shown in the examples. For information about configuring port numbers used by Tomcat or Jetty, see [Managing Solr](#).

Getting Started

Solr makes it easy for programmers to develop sophisticated, high-performance search applications with advanced features such as faceting (arranging search results in columns with numerical counts of key terms). Solr builds on another open source search technology: Lucene, a Java library that provides indexing and search technology, as well as spellchecking, hit highlighting and advanced analysis/tokenization capabilities. Both Solr and Lucene are managed by the Apache Software Foundation (www.apache.org).

The Lucene search library currently ranks among the top 15 open source projects and is one of the top 5 Apache projects, with installations at over 4,000 companies. Lucene/Solr downloads have grown nearly ten times over the past three years, with a current run-rate of over 6,000 downloads a day. The Solr search server, which provides application builders a ready-to-use search platform on top of the Lucene search library, is the fastest growing Lucene sub-project. Apache Lucene/Solr offers an attractive alternative to the proprietary licensed search and discovery software vendors.

This section helps you get Solr up and running quickly, and introduces you to the basic Solr architecture and features. It covers the following topics:

[Installing Solr](#): A walkthrough of the Solr installation process.

[Running Solr](#): An introduction to running Solr. Includes information on starting up the servers, adding documents, and running queries.

[A Quick Overview](#): A high-level overview of how Solr works.

[A Step Closer](#): An introduction to Solr's home directory and configuration options.

Installing Solr

This section describes how to install Solr. You can install Solr anywhere that a suitable Java Runtime Environment (JRE) is available, as detailed below. Currently this includes Linux, OS X, and Microsoft Windows. The instructions in this section should work for any platform, with a few exceptions for Windows as noted.

Got Java?

You will need the Java Runtime Environment (JRE) version 1.6 or higher. At a command line, check your Java version like this:

```
$ *java -version*
java version "1.6.0_0"
IcedTea6 1.3.1 (6b12-0ubuntu6.1) Runtime Environment (build 1.6.0_0-b12)
OpenJDK Client VM (build 1.6.0_0-b12, mixed mode, sharing)
```

The output will vary, but you need to make sure you have version 1.6 or higher. If you don't have the required version, or if the java command is not found, download and install the latest version from Sun at <http://java.sun.com/javase/downloads/>.

Installing Solr

Solr is available from the Solr website at <http://lucene.apache.org/solr/>.

For Linux/Unix/OSX systems, download the .gzip file. For Microsoft Windows systems, download the .zip file.

Solr runs inside a Java servlet container such as Tomcat, Jetty, or Resin. The Solr distribution includes a working demonstration server in the Example directory that runs in Jetty. You can use the example server as a template for your own installation, whether or not you are using Jetty as your servlet container. For more information about the demonstration server, see the [Solr Tutorial](#).

To install Solr

1. Unpack the Solr distribution to your desired location.
2. Stop your Java servlet container.
3. Copy the solr.war file from the Solr distribution to the webapps directory of your servlet container. Do not change the name of this file: it must be named solr.war.
4. Copy the Solr Home directory apache-solr-3.x.0/example/solr/ from the distribution to your desired Solr Home location.

5. Start your servlet container, passing to it the location of your Solr Home in one of these ways:

- Set the Java system property `solr.solr.home` to your Solr Home. (for example, using the example jetty setup: `java -Dsolr.solr.home=/some/dir -jar start.jar`).
- Configure the servlet container so that a JNDI lookup of `java:comp/env/solr/home` by the Solr webapp will point to your Solr Home.
- Start the servlet container in the directory containing `./solr`: the default Solr Home is `solr` under the JVM's current working directory (`$CWD/solr`).

To confirm your installation, go to the [Solr Admin page](#) at `http://_hostname_:8983/solr/admin/`. Note that your servlet container may have started on a different port: check the documentation for your servlet container to troubleshoot that issue. Also note that if that port is already in use, Solr will not start. In that case, shut down the servlet container running on that port, or change your Solr port.

For more information about installing and running Solr on different Java servlet containers, see the [SolrInstall](#) page on the [Solr Wiki](#).

Related Topics

- [SolrInstall](#)



Running Solr

This section describes how to run Solr with an example schema, how to add documents, and how to run queries.

Start the Server

If you didn't start Solr after installing it, you can start it by running `start.jar` from the Solr example directory.

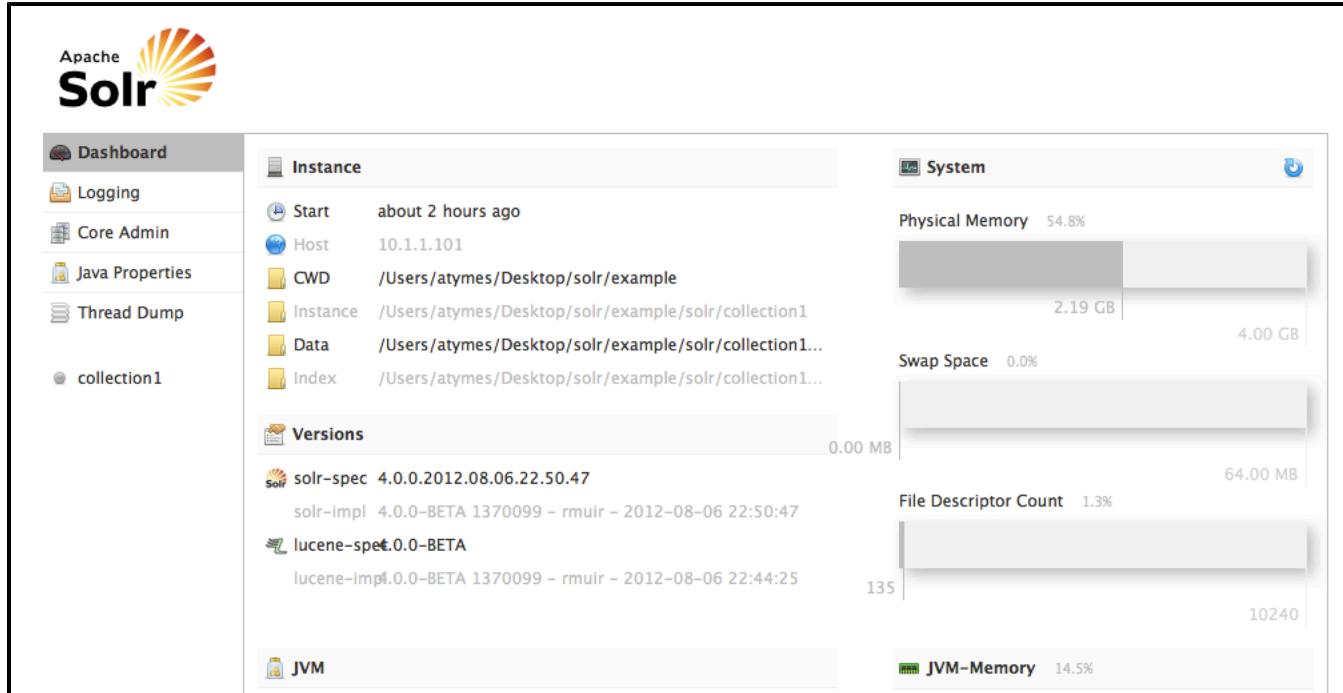
```
$ java -jar start.jar
```

If you are running Windows, you can start the Web server by running `start.bat` instead.

```
C:\Applications\Solr\example > start.bat
```

That's it! Solr is running. If you need convincing, use a Web browser to see the Admin Console.

<http://localhost:8983/solr/admin>



The Solr Admin interface.

If Solr is not running, your browser will complain that it cannot connect to the server. Check your port number and try again.

Add Documents

Solr is built to find documents that match queries. Solr's schema provides an idea of how content is structured (more on the schema [later](#)), but without documents there is nothing to find. Solr needs input before it can do anything.

You may want to add a few sample documents before trying to index your own content. The Solr installation comes with example documents located in the `example/exampledocs` directory of your installation.

In the `exampledocs` directory is the `SimplePostTool`, a Java-based command line tool, `post.jar`, which can be used to index the documents. Do not worry too much about the details for now. The [Indexing and Basic Data Operations](#) section has all the details on indexing.

To see some information about the usage of `post.jar`, use the `-help` option.

```
$ java -jar post.jar -help
```

The `SimplePostTool` is a simple command line tool for POSTing raw XML to a Solr port. XML data can be read from files specified as command line arguments, as raw command line `arg` strings, or via STDIN.

Examples:

```
java -Ddata=files -jar post.jar *.xml  
java -Ddata=args -jar post.jar '<delete><id>42</id></delete>'  
java -Ddata=stdin -jar post.jar < hd.xml
```

Other options controlled by System Properties include the Solr URL to POST to, and whether a commit should be executed. These are the defaults for all System Properties:

```
-Ddata=files  
-Durl=http://localhost:8983/solr/update  
-Dcommit=yes
```

Go ahead and add all the documents in the directory as follows:

```
$ *java -Durl=http://localhost:8983/solr/update -jar post.jar *.xml*
SimplePostTool: version 1.2
SimplePostTool: WARNING: Make sure your XML documents are encoded in UTF-8, other
encodings are not currently supported
SimplePostTool: POSTing files to http://10.211.55.8:8983/solr/update..
SimplePostTool: POSTing file hd.xml
SimplePostTool: POSTing file ipod_other.xml
SimplePostTool: POSTing file ipod_video.xml
SimplePostTool: POSTing file mem.xml
SimplePostTool: POSTing file monitor.xml
SimplePostTool: POSTing file monitor2.xml
SimplePostTool: POSTing file mp500.xml
SimplePostTool: POSTing file sd500.xml
SimplePostTool: POSTing file solr.xml
SimplePostTool: POSTing file spellchecker.xml
SimplePostTool: POSTing file utf8-example.xml
SimplePostTool: POSTing file vidcard.xml
SimplePostTool: COMMITting Solr index changes..
$
```

That's it! Solr has indexed the documents contained in the files.

Ask Questions

Now that you have indexed documents, you can perform queries. The simplest way is by building a URL that includes the query parameters. This is exactly the same as building any other HTTP URL.

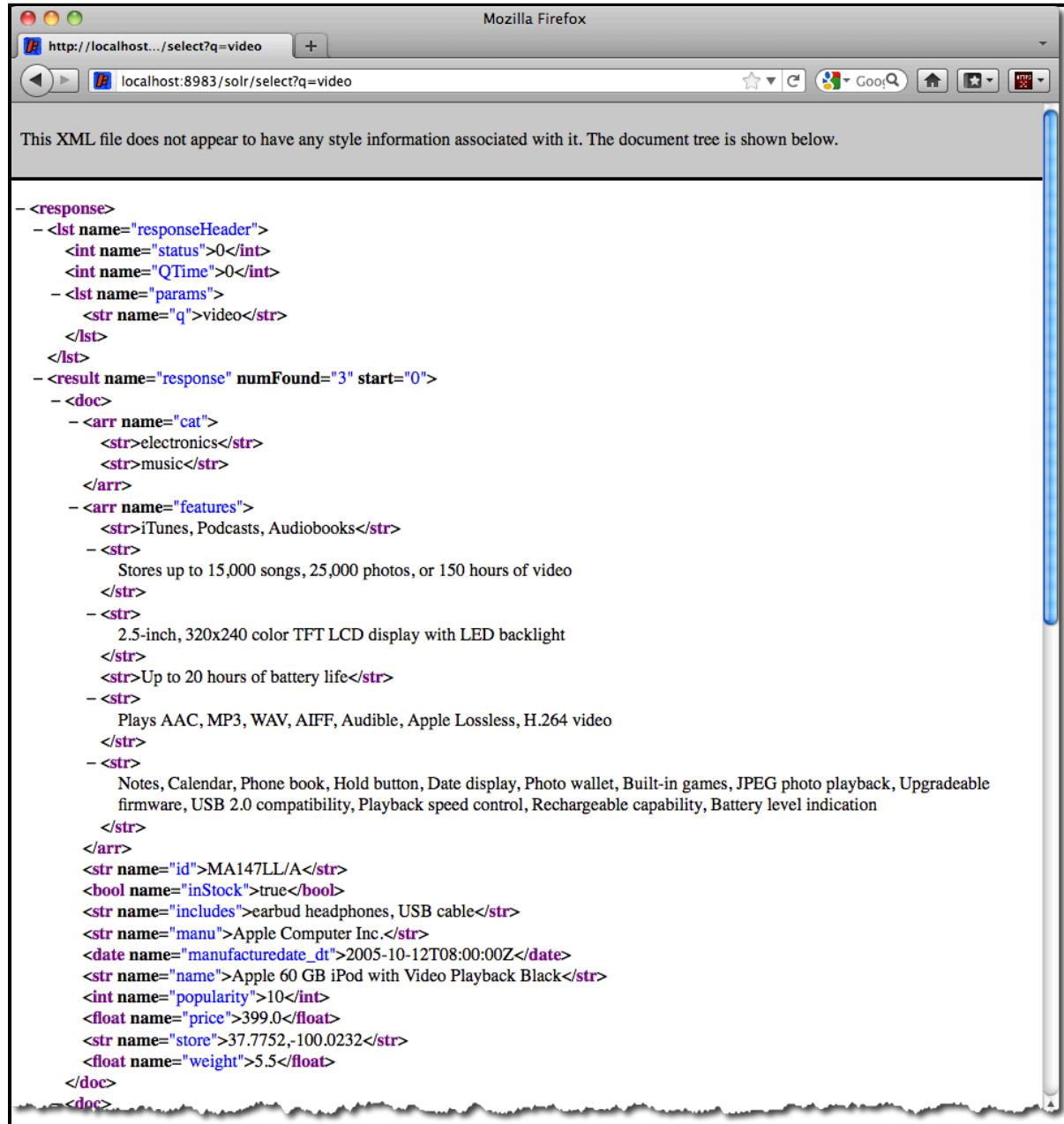
For example, the following query searches all document fields for "video":

```
http://localhost:8983/solr/select?q=video
```

Notice how the URL includes the host name (`localhost`), the port number where the server is listening (8983), the application name (`solr`), the request handler for queries (`select`), and finally, the query itself (`q=video`).

The results are contained in an XML document, which you can examine directly by clicking on the link above. The document contains two parts. The first part is the `responseHeader`, which contains information about the response itself. The main part of the reply is in the `result` tag, which contains one or more `doc` tags, each of which contains fields from documents that match the query. You can use standard XML transformation techniques to mold Solr's results into a form that is suitable for displaying to users. Alternatively, Solr can output the results in JSON, PHP, Ruby and even user-defined formats.

Just in case you are not running Solr as you read, the following screen shot shows the result of a query (the next example, actually) as viewed in Mozilla Firefox. The top-level response contains a 1st named `responseHeader` and a result named `response`. Inside `result`, you can see the three docs that represent the search results.



This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
- <response>
  - <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">0</int>
  - <lst name="params">
    <str name="q">video</str>
  </lst>
</lst>
- <result name="response" numFound="3" start="0">
  - <doc>
    - <arr name="cat">
      <str>electronics</str>
      <str>music</str>
    </arr>
    - <arr name="features">
      <str>iTunes, Podcasts, Audiobooks</str>
      - <str>
        Stores up to 15,000 songs, 25,000 photos, or 150 hours of video
      </str>
      - <str>
        2.5-inch, 320x240 color TFT LCD display with LED backlight
      </str>
      <str>Up to 20 hours of battery life</str>
      - <str>
        Plays AAC, MP3, WAV, AIFF, Audible, Apple Lossless, H.264 video
      </str>
      - <str>
        Notes, Calendar, Phone book, Hold button, Date display, Photo wallet, Built-in games, JPEG photo playback, Upgradeable firmware, USB 2.0 compatibility, Playback speed control, Rechargeable capability, Battery level indication
      </str>
    </arr>
    <str name="id">MA147LL/A</str>
    <bool name="inStock">true</bool>
    <str name="includes">earbud headphones, USB cable</str>
    <str name="manu">Apple Computer Inc.</str>
    <date name="manufacturedate_dt">2005-10-12T08:00:00Z</date>
    <str name="name">Apple 60 GB iPod with Video Playback Black</str>
    <int name="popularity">10</int>
    <float name="price">399.0</float>
    <str name="store">37.7752,-100.0232</str>
    <float name="weight">5.5</float>
  </doc>
</doc>
```

An XML response to a query.

Once you have mastered the basic idea of a query, it is easy to add enhancements to explore the query syntax. This one is the same as before but the results only contain the ID, name, and price for each returned document. If you don't specify which fields you want, all of them are returned.

```
http://localhost:8983/solr/select?q=video&fl=id,name,price
```

Here is another example which searches for "black" in the `name` field only. If you do not tell Solr which field to search, it will search default fields, as specified in the schema.

```
http://localhost:8983/solr/select?q=name:black
```

You can provide ranges for fields. The following query finds every document whose price is between \$0 and \$400.

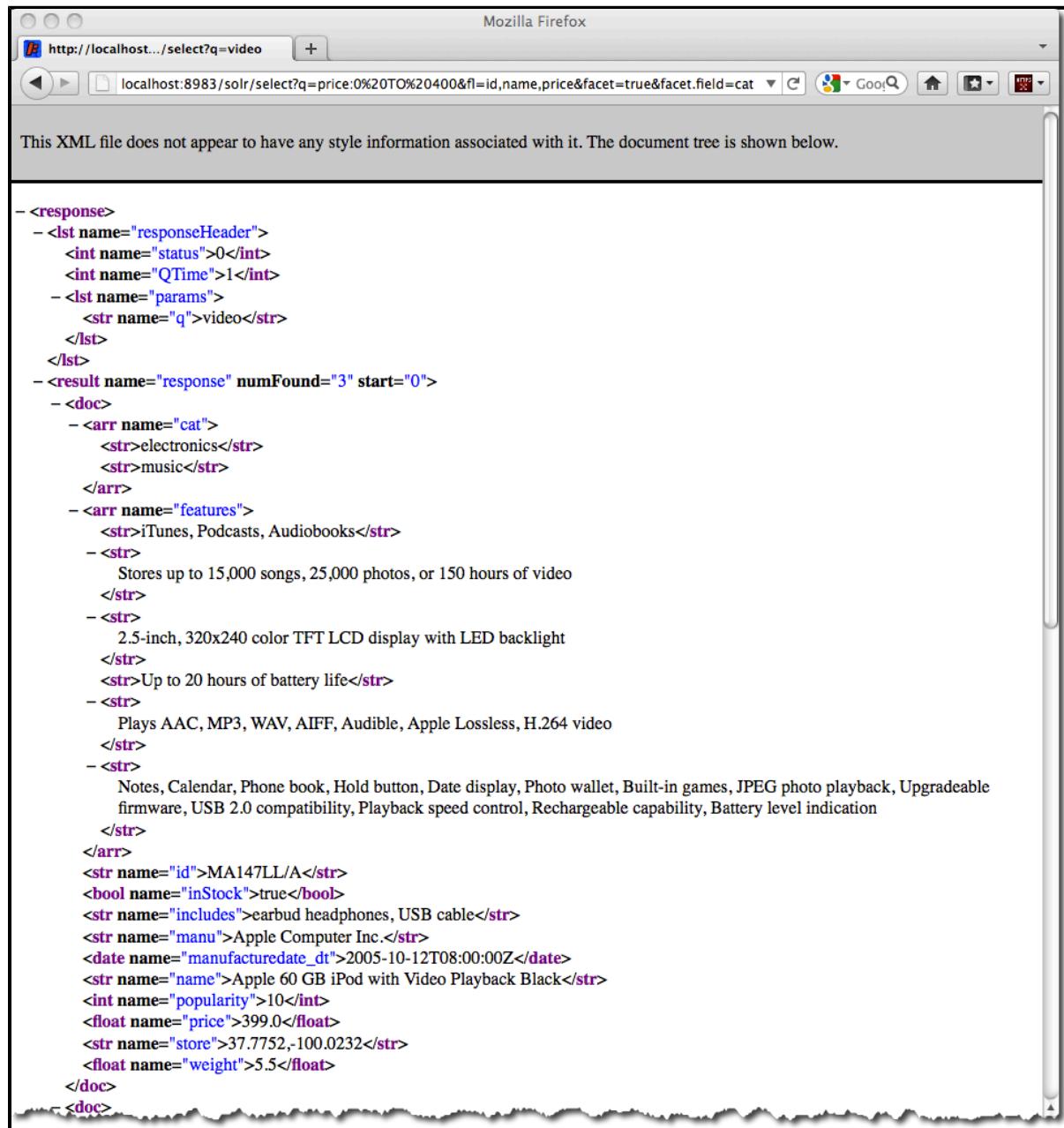
```
http://localhost:8983/solr/select?q=price:0%20TO%20400&fl=id,name,price
```

[Faceted browsing](#) is one of Solr's key features. It allows users to narrow search results in ways that are meaningful to your application. For example, a shopping site could provide facets to narrow search results by manufacturer or price.

Faceting information is returned as a third part of Solr's query response. To get a taste of this power, take a look at the following query. It adds `facet=true` and `facet.field=cat`.

```
http://localhost:8983/solr/select?q=price:0%20TO%20400&fl=id,name,price&facet=true&f
```

In addition to the familiar `responseHeader` and `response` from Solr, a `facet_counts` element is also present. Here is a view with the `responseHeader` and `response` collapsed so you can see the faceting information clearly.



This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
- <response>
- <lst name="responseHeader">
  <int name="status">0</int>
  <int name="QTime">1</int>
- <lst name="params">
  <str name="q">video</str>
</lst>
</lst>
- <result name="response" numFound="3" start="0">
- <doc>
  - <arr name="cat">
    <str>electronics</str>
    <str>music</str>
  </arr>
  - <arr name="features">
    <str>iTunes, Podcasts, Audiobooks</str>
    - <str>
      Stores up to 15,000 songs, 25,000 photos, or 150 hours of video
    </str>
    - <str>
      2.5-inch, 320x240 color TFT LCD display with LED backlight
    </str>
    <str>Up to 20 hours of battery life</str>
    - <str>
      Plays AAC, MP3, WAV, AIFF, Audible, Apple Lossless, H.264 video
    </str>
    - <str>
      Notes, Calendar, Phone book, Hold button, Date display, Photo wallet, Built-in games, JPEG photo playback, Upgradeable firmware, USB 2.0 compatibility, Playback speed control, Rechargeable capability, Battery level indication
    </str>
  </arr>
  <str name="id">MA147LL/A</str>
  <bool name="inStock">true</bool>
  <str name="includes">earbud headphones, USB cable</str>
  <str name="manu">Apple Computer Inc.</str>
  <date name="manufacturedate_dt">2005-10-12T08:00:00Z</date>
  <str name="name">Apple 60 GB iPod with Video Playback Black</str>
  <int name="popularity">10</int>
  <float name="price">399.0</float>
  <str name="store">37.7752,-100.0232</str>
  <float name="weight">5.5</float>
</doc>
<doc>
```

An XML Response with faceting.

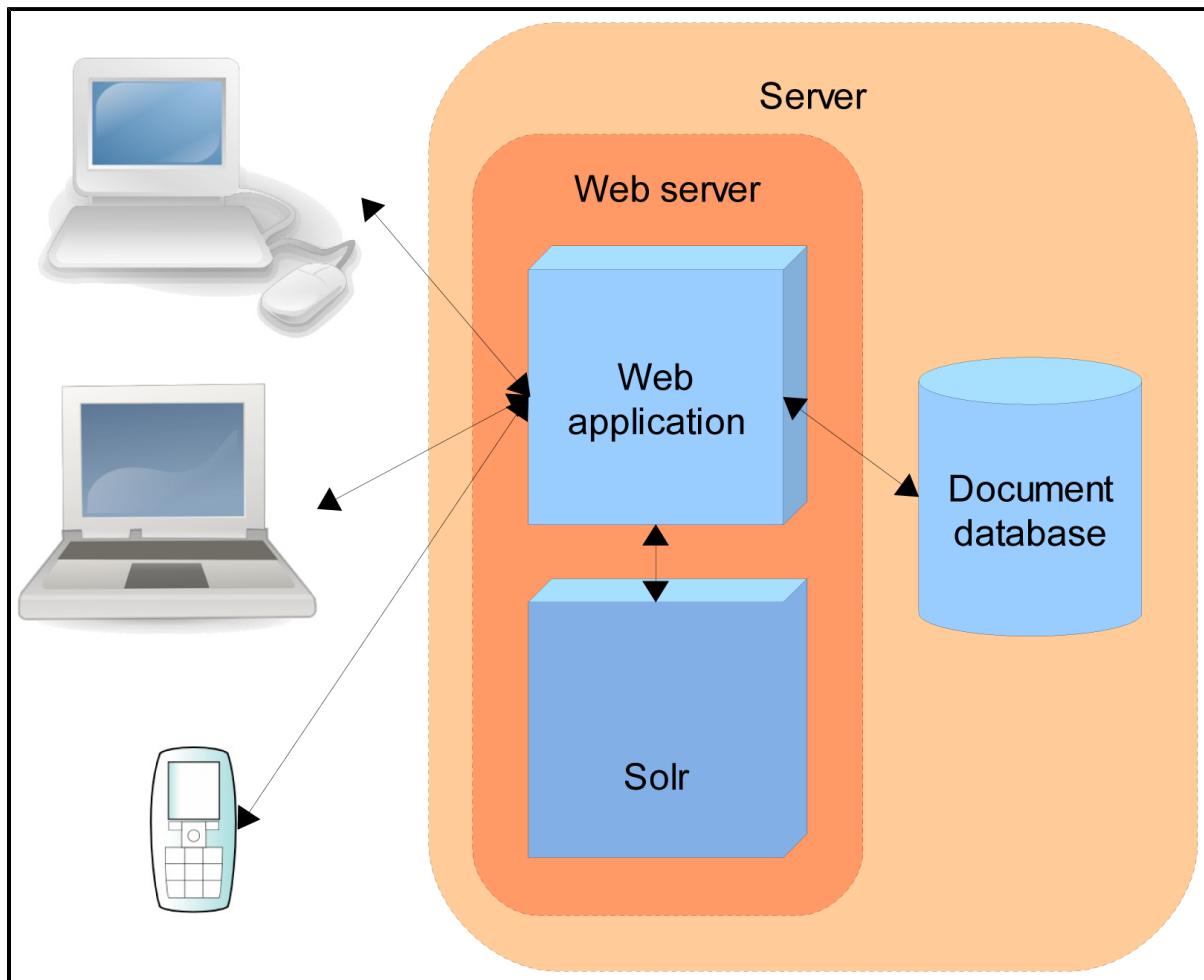
The facet information shows how many of the query results have each possible value of the `cat` field. You could easily use this information to provide users with a quick way to narrow their query results. You can filter results by adding one or more filter queries to the Solr request. Here is a request further constraining the request to documents with a category of "software".

`http://localhost:8983/solr/select?q=price:0%20TO%20400&fl=id,name,price&facet=true&f`

A Quick Overview

Having had some fun with Solr, you will now learn about all the cool things it can do.

Here is a typical configuration:



In the scenario above, Solr runs alongside another application in a Web server. For example, an online store application would provide a user interface, a shopping cart, and a way to make purchases. The store items would be kept in some kind of database.

Solr makes it easy to add the capability to search through the online store through the following steps:

1. Define a *schema*. The schema tells Solr about the contents of documents it will be indexing. In the online store example, the schema would define fields for the product name, description, price, manufacturer, and so on. Solr's schema is powerful and flexible and allows you to tailor Solr's behavior to your application. See [Documents, Fields, and Schema Design](#) for all the details.

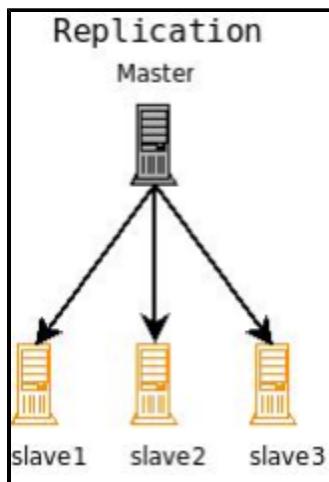
2. Deploy Solr to your application server.
3. Feed Solr the document for which your users will search.
4. Expose search functionality in your application.

Because Solr is based on open standards, it is highly extensible. Solr queries are RESTful, which means, in essence, that a query is a simple HTTP request URL and the response is a structured document: mainly XML, but it could also be JSON, CSV, or some other format. This means that a wide variety of clients will be able to use Solr, from other web applications to browser clients, rich client applications, and mobile devices. Any platform capable of HTTP can talk to Solr. See [Client APIs](#) for details on client APIs.

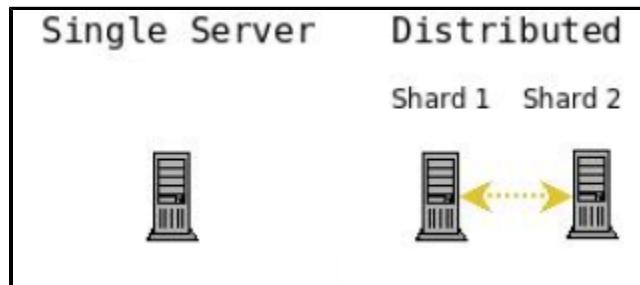
Solr is based on the Apache Lucene project, a high-performance, full-featured search engine. Solr offers support for the simplest keyword searching through to complex queries on multiple fields and faceted search results. [Searching](#) has more information about searching and queries.

If Solr's capabilities are not impressive enough, its ability to handle very high-volume applications should do the trick.

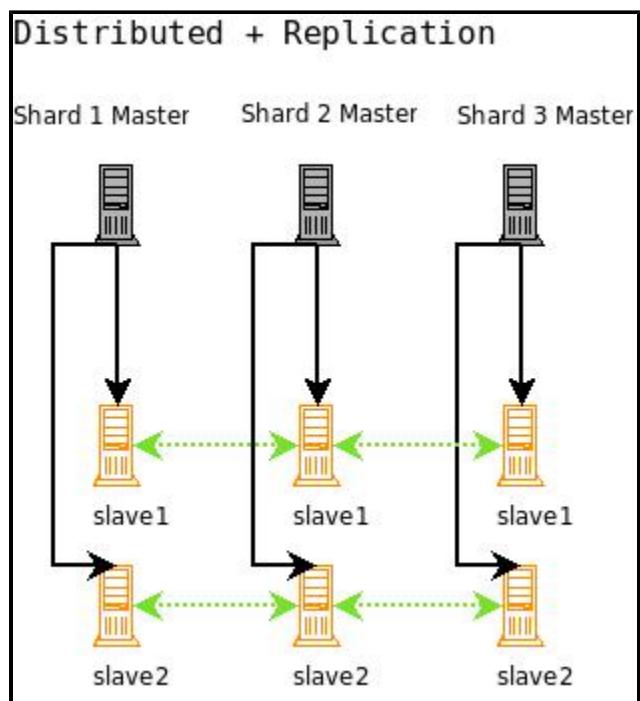
A relatively common scenario is that you have so many queries that the server is unable to respond fast enough to each one. In this case, you can make copies of an index. This is called replication. Then you can distribute incoming queries among the copies in any way you see fit. A round-robin mechanism is one simple way to do this.



Another useful technique is sharding. If you have so many documents that you simply cannot fit them all on a single box for RAM or index size reasons, you can split an index into multiple pieces, called *shards*. Each shard lives on its own physical server. An incoming query is sent to all the shard servers, which respond with matching results.



If you have huge numbers of documents and users, you might need to combine the techniques of sharding and replication. In this case, Solr's new SolrCloud functionality may be more effective for your needs. Solrcloud includes a number of features to simplify the process of distributing the index and the queries, and manage the resulting nodes.



For full details on sharding and replication, see [Legacy Scaling and Distribution](#). We've split the SolrCloud information into it's own section, called [SolrCloud](#).

Best of all, this talk about high-volume applications is not just hypothetical: some of the famous Internet sites that use Solr today are Macy's, EBay, and Zappo's.

For more information, take a look at <https://wiki.apache.org/solr/PublicServers>.

A Step Closer

You already have some idea of Solr's schema. This section describes Solr's home directory and other configuration options.

When Solr runs in an application server, it needs access to a home directory. The home directory contains important configuration information and is the place where Solr will store its index.

The crucial parts of the Solr home directory are shown here:

```
<solr-home-directory>/  
  solr.xml  
  conf/  
    solrconfig.xml  
    schema.xml  
  data/
```

You supply `solr.xml`, `solrconfig.xml`, and `schema.xml` to tell Solr how to behave. By default, Solr stores its index inside `data`.

`solr.xml` specifies configuration options for your Solr core, and also allows you to configure multiple cores. For more information on `solr.xml` see [The Well-Configured Solr Instance](#).

`solrconfig.xml` controls high-level behavior. You can, for example, specify an alternate location for the `data` directory. For more information on `solrconfig.xml`, see [The Well-Configured Solr Instance](#).

`schema.xml` describes the documents you will ask Solr to index. Inside `schema.xml`, you define a document as a collection of fields. You get to define both the field types and the fields themselves. Field type definitions are powerful and include information about how Solr processes incoming field values and query values. For more information on `schema.xml`, see [Documents, Fields, and Schema Design](#).

Upgrading to Solr 4.1

If you are already using Solr 4.0, moving to Solr 4.1 should not present any major problems. However, you should review the `CHANGES.txt` file found in your Solr package for changes and updates that may effect your existing implementation.

During updates for this guide, we found one thing that you may want to be aware of, which is a change to the way Solr identifies node names for SolrCloud. If you are using SolrCloud, you may have issues with unknown or lost nodes. If this occurs, you can manually set the `host` parameter either in `solr.xml` or as a system variable. More information can be found in the section on [SolrCloud](#).

With any Solr update, it is recommended that you re-index your content for use with Solr 4.1.

Using the Solr Administration User Interface

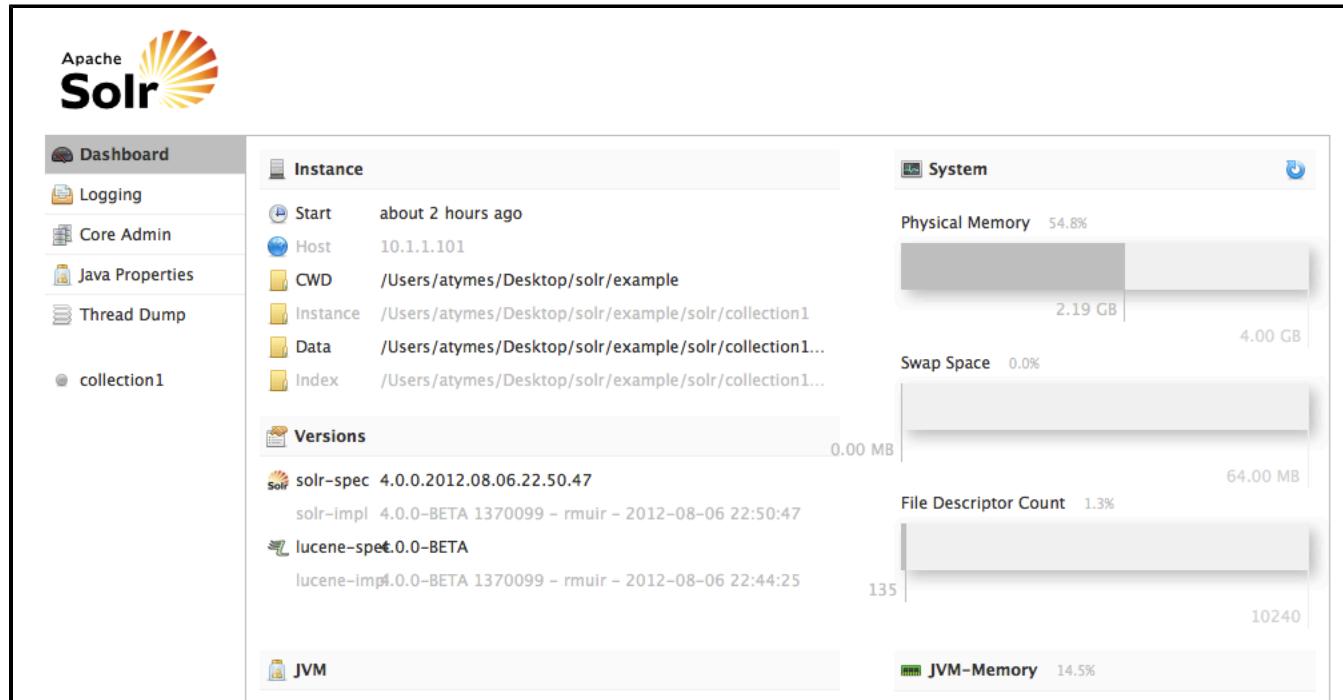
This section discusses the Solr Administration User Interface ("Admin UI").

The [Overview of the Solr Admin UI](#) explains how the features of the user interface that are new with Solr 4, what's on the initial Admin UI page, and how to configure the interface. In addition, there are pages describing each screen of the Admin UI:

- [Getting Assistance](#) shows you how to get more information about the UI.
- [Logging](#) explains the various logging levels available and how to invoke them.
- [Cloud Screens](#) display information about nodes when running in SolrCloud mode.
- [Core Admin](#) explains how to get management information about each core.
- [Java Properties](#) shows the Java information about each core.
- [Thread Dump](#) lets you see detailed information about each thread, along with state information.
- [Core-Specific Tools](#) is a section explaining additional screens available for each named core.
 - [Ping](#) lets you ping a named core and determine whether the core is active.
 - [Query Screen](#) lets you submit a structured query about various elements of a core.
 - [Schema Screen](#) describes the schema.xml file for the core.
 - [Config Screen](#) shows the current configuration of the solrconfig.xml file for the core.
 - [Analysis Screen](#) lets you analyze the data found in specific fields.
 - [Schema Browser](#) displays schema data in a browser window.
 - [Replication Screen](#) shows you the current replication status for the core, and lets you enable/disable replication.
 - [Plugins & Stats Screen](#) shows statistics for plugins and other installed components.
 - [Dataimport Screen](#) shows you information about the current status of the Data Import Handler.

Overview of the Solr Admin UI

Solr features a Web interface that makes it easy for Solr administrators and programmers to view [Solr configuration](#) details, run [queries and analyze](#) document fields in order to fine-tune a Solr configuration and access [online documentation](#) and other help.



With Solr 4, the Solr Admin has been completely redesigned. The redesign was completed with these benefits in mind:

- load pages quicker
- access and control functionality from the Dashboard
- re-use the same servlets that access Solr-related data from an external interface, and
- ignore any differences between working with one or multiple cores.

Accessing the URL <http://hostname:8983/solr/admin/> (if running Jetty on the default port of 8983), will show the main dashboard, which is divided into two parts.

A left-side of the screen is a menu under the Solr logo that provides the navigation through the screens of the UI. The first set of links are for system-level information and configuration and provide access to Logging, Core Admin and Java Properties, among other things. At the end of this information is a list of Solr cores configured for this instance. Clicking on a core name shows a secondary menu of information and configuration options for the core specifically. Items in this list include the Schema, Config, Plugins, and an ability to perform Queries on indexed data.

The center of the screen shows the detail of the option selected. This may include a sub-navigation for the option or text or graphical representation of the requested data. See the sections in this guide for each screen for more details.

 If you are running Solr on a Macintosh, you should access the Admin UI in a browser other than Safari, because Safari will not display raw XML content, such as the contents of the `schema.xml` file.

Configuring the Admin UI in `solrconfig.xml`

You can configure the Solr Admin UI by editing the file `solrconfig.xml`.

The `<admin>` block in the `solrconfig.xml` file determines the default query to be displayed in the Query section of the core-specific pages. The default is `*:*`, which is to find all documents. In this example, we have changed the default to the term `solr`.

```
<admin>
  <defaultQuery>solr</defaultQuery>
</admin>
```

Related Topics

- [Configuring `solrconfig.xml`](#)

Getting Assistance

At the bottom of each screen of the Admin UI is a set of links that can be used to get more assistance with configuring and using Solr.

[Documentation](#)[Issue Tracker](#)[IRC Channel](#)[Community forum](#)[Solr Query Syntax](#)

Assistance icons

These icons include the following links.

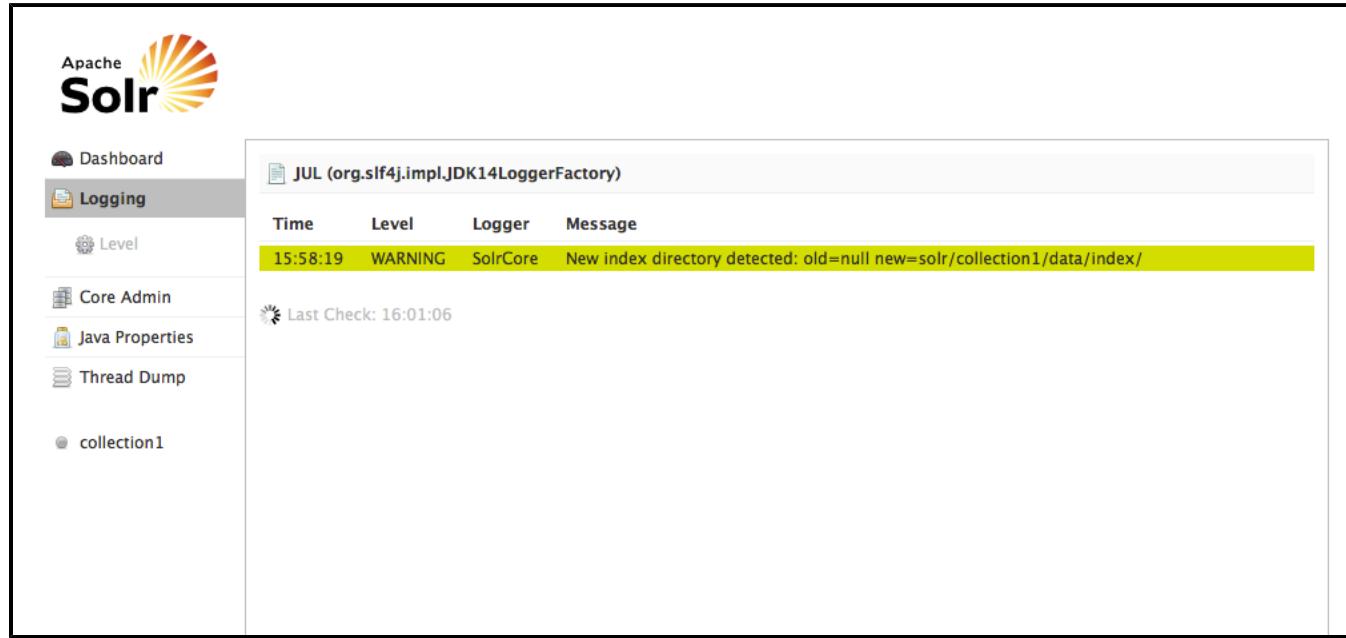
Link	Description
Documentation	Navigates to the Apache Solr documentation hosted on http://lucene.apache.org/solr/ .
Issue Tracker	Navigates to the JIRA issue tracking server for the Apache Solr project. This server resides at http://issues.apache.org/jira/browse/SOLR .
IRC Channel	Connects you to the web interface for Solr's IRC channel. This channel is found on Irc.freenode.net , Port 7000, #solr channel.
Community forum	Connects you to the Solr community forum , which at the current time is a set of mailing lists and their archives.
Solr Query Syntax	Navigates to the Apache Wiki page describing the Solr query syntax: http://wiki.apache.org/solr/SolrQuerySyntax .

These links cannot be modified without editing the `admin.html` in the `solr.war` that contains the Admin UI files.

Logging

The Logging page shows messages from Solr's log files.

When you click the link for "Logging", a page similar to the one below will be displayed:



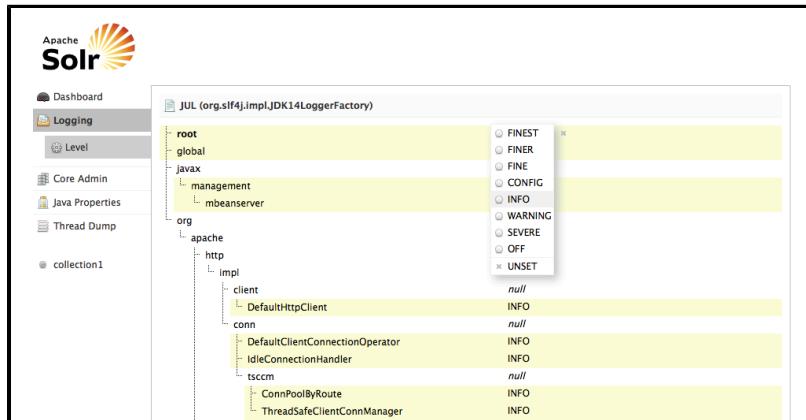
The screenshot shows the Apache Solr admin interface with the "Logging" tab selected in the sidebar. The main content area displays a table of log entries under the heading "JUL (org.slf4j.impl.JDK14LoggerFactory)". The table has columns for Time, Level, Logger, and Message. A single entry is shown: "15:58:19 WARNING SolrCore New index directory detected: old=null new=solr/collection1/data/index/". Below the table, a message indicates "Last Check: 16:01:06".

Time	Level	Logger	Message
15:58:19	WARNING	SolrCore	New index directory detected: old=null new=solr/collection1/data/index/

The Main Logging Screen

While this example shows logged messages for only one core, if you have multiple cores in a single instance, they will each be listed, with the level for each.

Selecting a Logging Level



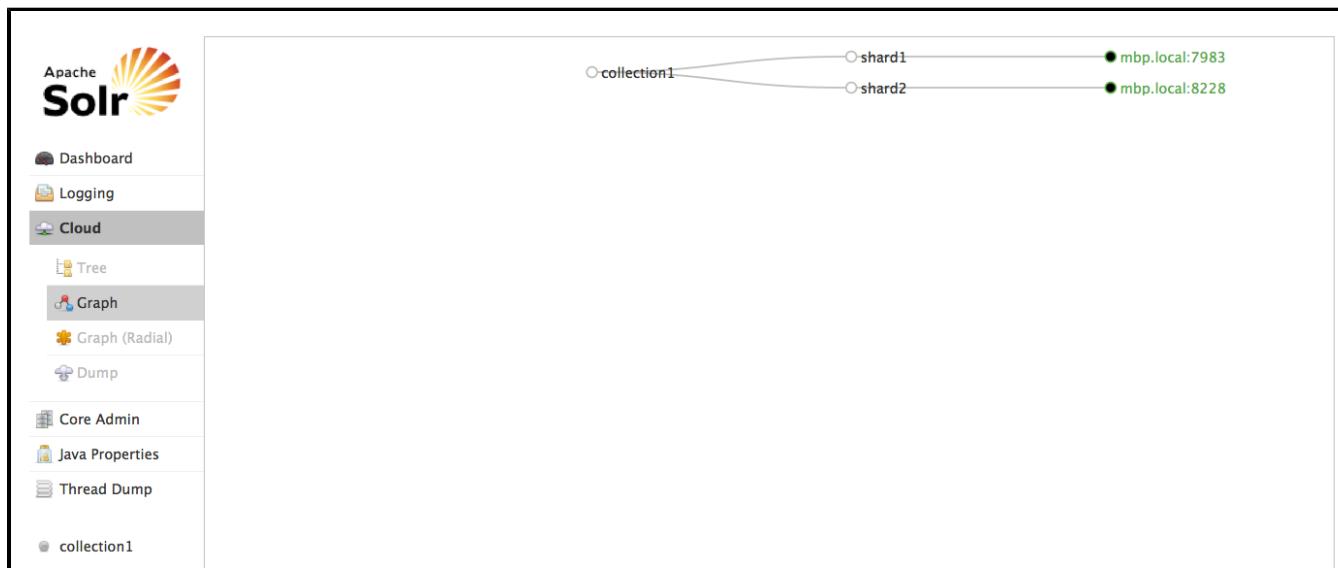
When you select the **Level** link on the left, you see the hierarchy of classpaths and classnames for your instance. A row highlighted in yellow indicates that the class has logging capabilities. Click on a highlighted row, and a menu will appear to allow you to change the log level for that class. Characters in boldface indicate that the class will not be affected by level changes to root.

For an explanation of the various logging levels, see [Configuring Logging](#).

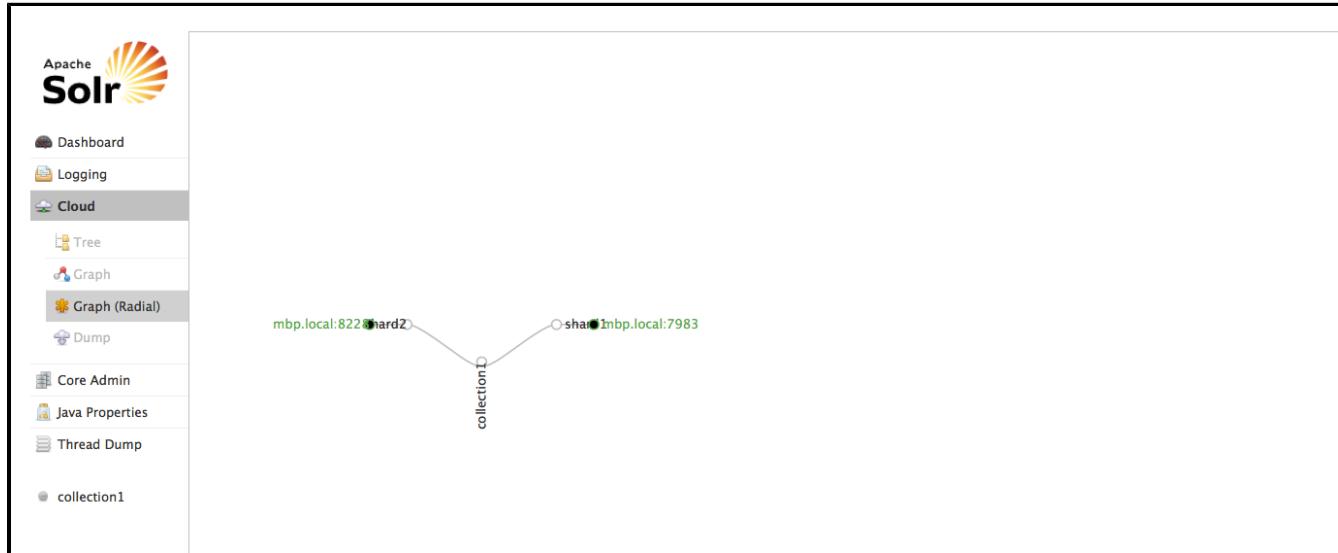
Cloud Screens

When running in SolrCloud mode, an option will appear in the Admin UI between Logging and Core Admin for Cloud. It's not possible at the current time to manage the nodes of the SolrCloud cluster, but you can view them and open the Solr Admin UI on each node to view the status and statistics for the node and each core on each node.

Click on the Cloud option in the left-hand navigation, and a small sub-menu appears with options called "Tree", "Graph", "Graph (Radial)" and "Dump". The default view (which is "Tree") shows a graph of each core and the addresses of each node. This example shows a very simple two-node cluster with a single core:



The "Graph (Radial)" option provides a different visual view of each node. Using the same simple two-node cluster, the radial graph view looks like:



The "Tree" option shows a directory structure of the files in ZooKeeper, including `clusterstate.json`, configuration files, and other status and information files. In this example, we show the leader definition files for the core named "collection1":



The final option is "Dump", which allows you to download an XML file with all the ZooKeeper configuration files.

Core Admin

The Core Admin screen lets you manage your cores.

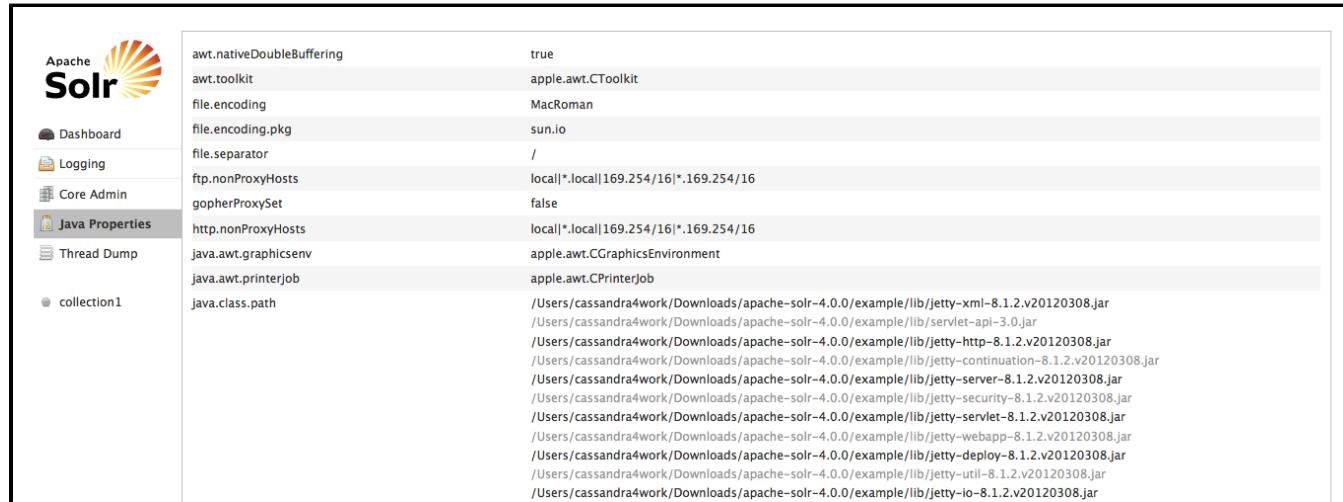
The buttons at the top of the screen let you add a new core, unload the core displayed, rename the currently displayed core, swap the existing core with one that you specify in a drop-down box, reload the current core, and optimize the current core.

The main display and available actions correspond to the commands used with the [CoreAdminHandler](#), but provide another way of working with your cores.

The screenshot shows the Solr Core Admin interface. On the left is a sidebar with links: Dashboard, Logging, Cloud, Core Admin (which is selected and highlighted in grey), Java Properties, Thread Dump, collection1, and test_core_shard1_r... . The main content area has a header with buttons: Add Core, Unload (red), Rename, Swap, Reload, and Optimize. Below this, there are two sections: 'Core' and 'Index'. The 'Core' section contains fields: startTime (less than a minute ago), instanceDir (solr/collection1/), and dataDir (solr/collection1/data/). The 'Index' section contains fields: lastModified (-), version (1), numDocs (0), maxDoc (0), deletedDocs (-), optimized (✓ checked), current (✓ checked), and directory (org.apache.lucene.store.NRTCachingDirectory:NRTCachingDirectory/org.apache.lucene.store.NIOFSDirectory@/Applications/solr-4.1.0/shard1/solr/collection1/data/index lockFactory=org.apache.lucene.store.NativeFSLockFactory@5f14a3c6; maxCacheMB=48.0 maxMergeSizeMB=4.0).

Java Properties

The Java Properties screen provides easy access to one of the most essential components of a top-performing Solr systems. With the Java Properties screen, you can see all the properties of the JVM running Solr, including the class paths, file encodings, JVM memory settings, operating system, and more.



The screenshot shows the Apache Solr Java Properties screen. On the left is a navigation sidebar with links: Dashboard, Logging, Core Admin, Java Properties (which is selected and highlighted in grey), Thread Dump, and collection1. The main content area displays a table of system properties:

Property	Value
awt.nativeDoubleBuffering	true
awt.toolkit	apple.awt.CToolkit
file.encoding	MacRoman
file.encoding.pkg	sun.io
file.separator	/
ftp.nonProxyHosts	local *.local 169.254/16 *.169.254/16
gopherProxySet	false
http.nonProxyHosts	local *.local 169.254/16 *.169.254/16
java.awt.graphicsenv	apple.awt.CGraphicsEnvironment
java.awt.printerjob	apple.awt.CPrinterJob
java.class.path	/Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-xml-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/servlet-api-3.0.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-http-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-continuation-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-server-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-security-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-servlet-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-webapp-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-deploy-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-util-8.1.2.v20120308.jar /Users/cassandra4work/Downloads/apache-solr-4.0.0/example/lib/jetty-io-8.1.2.v20120308.jar

Thread Dump

The Thread Dump screen lets you inspect the currently active threads on your server. Each thread is listed and access to the stacktraces is available where applicable. Icons to the left indicate the state of the thread: for example, threads with a green check-mark in a green circle are in a "RUNNABLE" state. On the right of the thread name, a down-arrow means you can expand to see the stacktrace for that thread.

name		cpuTime / userTime
<input checked="" type="checkbox"/> DestroyJavaVM (26)		3536.3050ms 3393.5180ms
<input checked="" type="checkbox"/> pool-1-thread-1 (25)		65.8300ms 63.7670ms
<input checked="" type="checkbox"/> HashSessionScavenger-0 (23)		1.3790ms 0.7880ms
<input checked="" type="checkbox"/> Poller SunPKCS11-Darwin (22)		18.1240ms 14.1830ms
<input checked="" type="checkbox"/> qtp1566301264-21 Acceptor0 SocketConnector@0.0.0.0:8983 (21)		30.8460ms 29.2750ms
<input checked="" type="checkbox"/> qtp1566301264-20 (20)		4.4950ms 3.0650ms
<input checked="" type="checkbox"/> qtp1566301264-19 (19)		93.9410ms 85.0270ms

When you move your cursor over a thread name, a box floats over the name with the state for that thread. Thread states can be:

State	Meaning
NEW	A thread that has not yet started.
RUNNABLE	A thread executing in the Java virtual machine.
BLOCKED	A thread that is blocked waiting for a monitor lock.
WAITING	A thread that is waiting indefinitely for another thread to perform a particular action.
TIMED_WAITING	A thread that is waiting for another thread to perform an action for up to a specified waiting time.

TERMINATED	A thread that has exited.
------------	---------------------------

When you click on one of the threads that can be expanded, you'll see the stacktrace, as in the example below:

The screenshot shows the Apache Solr interface with the "Thread Dump" section selected. It displays a list of threads with their names and execution times. Two threads are expanded to show their full stacktraces:

name	cpuTime / userTime
DestroyJavaVM (26)	3536.3050ms 3393.5180ms
pool-1-thread-1 (25)	65.8300ms 63.7670ms
HashSessionScavenger-0 (23)	1.3790ms 0.7880ms

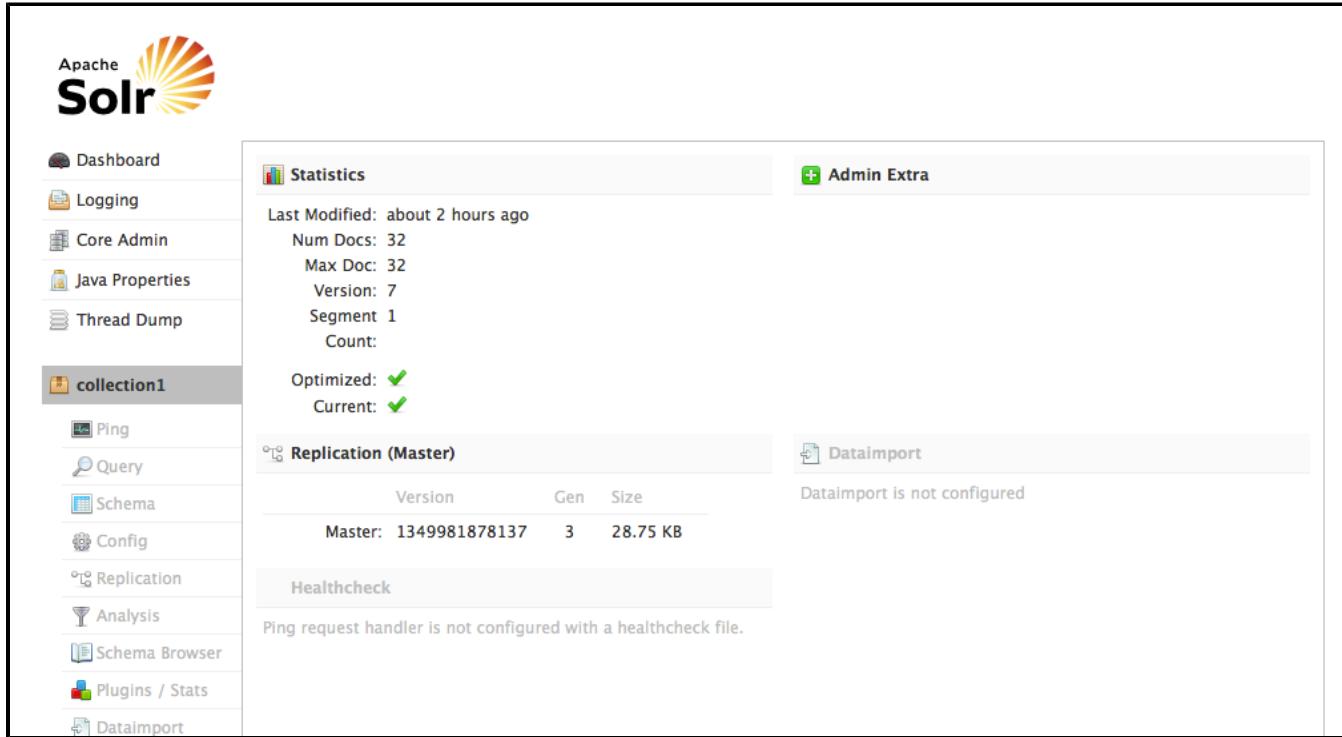
Expanded Stacktraces:

- pool-1-thread-1 (25)**
 - java.util.concurrent.locks.AbstractQueuedSynchronizer\$ConditionObject@6cf84b0a
 - sun.misc.Unsafe.park(Native Method)
 - java.util.concurrent.locks.LockSupport.park(LockSupport.java:156)
 - java.util.concurrent.locks.AbstractQueuedSynchronizer\$ConditionObject.await(AbstractQueuedSynchronizer.java:1987)
 - java.util.concurrent.LinkedBlockingQueue.take(LinkedBlockingQueue.java:399)
 - java.util.concurrent.ThreadPoolExecutor.getTask(ThreadPoolExecutor.java:947)
 - java.util.concurrent.ThreadPoolExecutor\$Worker.run(ThreadPoolExecutor.java:907)
 - java.lang.Thread.run(Thread.java:680)
- HashSessionScavenger-0 (23)**
 - java.util.TaskQueue@2106b56a
 - java.lang.Object.wait(Native Method)
 - java.util.TimerThread.mainLoop(Timer.java:509)

Inspecting a thread

You can also check the **Show all Stacktraces** button to automatically enable expansion for all threads.

Core-Specific Tools



When you click the name of a core, such as **collection1** in the example, a secondary menu opens under the core name with the administration options available for that particular core.

After selecting the core, the central part of the screen shows Statistics and other information about the core you chose. You can define a file called `admin-extra.html` that includes links or other information you would like to display in the "Admin Extra" part of this main screen.

On the left side, under the core name, are links to other screens that display information or provide options for the specific core chosen. The core-specific options are listed below, with a link to the section of this Guide to find out more:

- [Ping](#)
- [Query](#)
- [Schema](#)
- [Config](#)
- [Replication](#)
- [Analysis](#)
- [Schema Browser](#)
- [Plugins/Stats](#)
- [Dataimport](#)

Ping

Choosing Ping under a core name issues a ping request to check whether a server is up.

Ping is configured using a `requestHandler` in the `solrconfig.xml` file:

```
<!-- ping/healthcheck -->
<requestHandler name="/admin/ping" class="solr.PingRequestHandler">
  <lst name="invariants">
    <str name="q">solrpingingquery</str>
  </lst>
  <lst name="defaults">
    <str name="echoParams">all</str>
  </lst>
  <!-- An optional feature of the PingRequestHandler is to configure the
      handler with a "healthcheckFile" which can be used to enable/disable
      the PingRequestHandler.
      relative paths are resolved against the data dir
  -->
  <!-- <str name="healthcheckFile">server-enabled.txt</str> -->
</requestHandler>
```

The Ping option doesn't open a page, but the status of the request can be seen on the core overview page shown when clicking on a collection name. The length of time the request has taken is displayed next to the Ping option, in milliseconds.

Query Screen

You can use **Query**, shown under the name of each core, to submit a search query to a Solr server and analyze the results. In the example in the screenshot, a query has been submitted, and the screen shows the query results sent to the browser as XML.

The screenshot shows the Solr Admin UI for the "collection1" core. The "Query" section is selected in the sidebar. The main form has the following values:

- Request-Handler (qt): /select
- q: *:*
- fq: (empty)
- sort: (empty)
- start, rows: 0, 10
- fl: (empty)
- df: (empty)
- wt: xml
- checkboxes: Indent, debugQuery, dismax, edismax, hl, facet, spatial, spellcheck

The XML response is as follows:

```

<response>
  <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">1</int>
  </lst>
  <result name="response" numFound="32" start="0">
    <doc>
      <str name="id">GB18030TEST</str>
      <str name="name">Test with some GB18030 encoded characters</str>
      <arr name="features">
        <str>No accents here</str>
        <str>这是一个功能</str>
        <str>This is a feature (translated)</str>
        <str>这份文件是很有光泽</str>
        <str>This document is very shiny (translated)</str>
      </arr>
      <float name="price">0.0</float>
      <str name="price_c">0.USD</str>
      <bool name="inStock">true</bool>
      <long name="_version_">1417559430731399168</long>
    </doc>
    <doc>
      <str name="id">SP2514N</str>
      <str name="name">
        Samsung SpinPoint P120 SP2514N - hard drive - 250 GB - ATA-133
      </str>
    </doc>
  </result>
</response>

```

The query was sent to a core named "collection1". We used Solr's default query for this screen (as defined in `solrconfig.xml`), which is `*:*`. This query will find all records in the index for this core. We kept the other defaults, but the table below explains these options, which are also covered in detail in later parts of this Guide.

The response is shown to the right of the form. Requests to Solr are simply HTTP requests, and the query submitted is shown in light type above the results; if you click on this it will open a new browser window with just this request and response (without the rest of the Solr Admin UI). The rest of the response is shown in XML, which is part of the request (see the `wt=xml` part at the end).

The response has at least two sections, but may have several more depending on the options chosen. The two sections it always has are the `responseHeader` and the `response`. The `responseHeader` includes the status of the search (`status`), the processing time (`QTime`), and the parameters (`params`) that were used to process the query.

The `response` includes the documents that matched the query, in `doc` sub-sections. The fields return depend on the parameters of the query (and the defaults of the request handler used). The number of results is also included in this section.

This screen allows you to experiment with different query options, and inspect how your documents were indexed. The query parameters available on the form are some basic options that most users want to have available, but there are dozens more available which could be simply added to the basic request by hand (if opened in a browser). The table below explains the parameters available:

Field	Description
Request-handler	Specifies the query handler for the request. If a query handler is not specified, Solr processes the response with the standard query handler.
q	The query event. See Searching for an explanation of this parameter.
fq	The filter query. See Common Query Parameters for more information on this parameter.
sort	Sorts the response to a query in either ascending or descending order based on the response's score or another specified characteristic.
start, rows	<code>start</code> is the offset into the query result starting at which documents should be returned. The default value is 0, meaning that the query should return results starting with the first document that matches. This field accepts the same syntax as the start query parameter, which is described in Searching . <code>rows</code> is the number of rows to return.
fl	Defines the fields to return for each document. You can explicitly list the stored fields you want to have returned by separating them with either a comma or a space. In Solr 4, the results of functions can also be included in the <code>fl</code> list.
wt	Specifies the Response Writer to be used to format the query response. Normally XML.
indent	Click this button to request that the XML Response Writer use indentation to make the XML response more readable.
debugQuery	Click this button to augment the query response with debugging information, including "explain info" for each document returned. This debugging information is intended to be intelligible to the administrator or programmer.
dismax	Click this button to enable the Dismax query parser. See The DisMax Query Parser for further information.
edismax	Click this button to enable the Extended query parser. See The Extended DisMax Query Parser for further information.
hl	Click this button to enable highlighting in the query response. See Highlighting for more information.

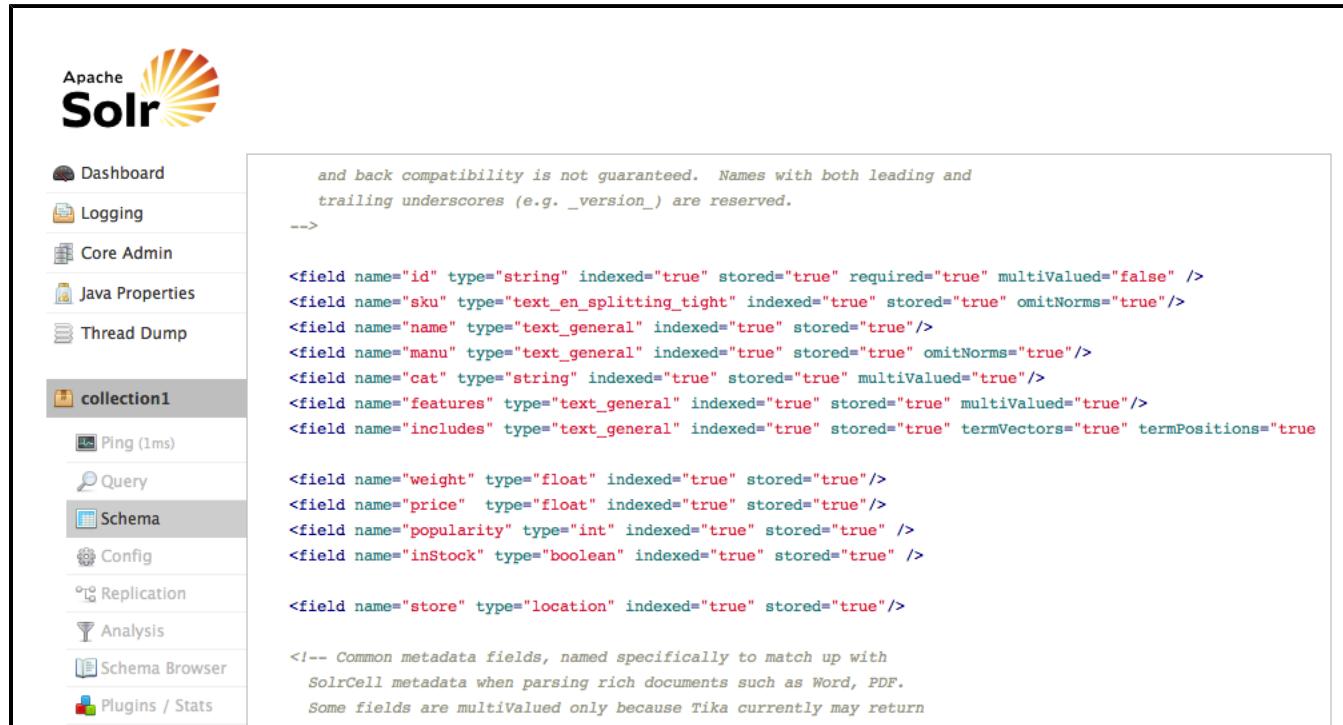
facet	Enables faceting, the arrangement of search results into categories based on indexed terms. See Faceting for more information.
spatial	Click to enable using location data for use in spatial or geospatial searches. See Spatial Search for more information.
spellcheck	Click this button to enable the Spellchecker, which provides inline query suggestions based on other, similar, terms. See Spell Checking for more information.

Related Topics

- [Searching](#)

Schema Screen

The Schema option displays the `schema.xml` file, a configuration file that describes the data to be indexed and searched.



The screenshot shows the Apache Solr interface. On the left, there's a sidebar with various links: Dashboard, Logging, Core Admin, Java Properties, Thread Dump, collection1 (which is selected), Ping (1ms), Query, Schema (which is also selected), Config, Replication, Analysis, Schema Browser, and Plugins / Stats. The main content area is titled "Schema" and contains the XML code for `schema.xml`. The code defines several fields like `id`, `sku`, `name`, etc., with various types and storage options. There are also comments explaining field naming conventions and metadata fields.

```
<!-- Compatibility note -->  
and back compatibility is not guaranteed. Names with both leading and  
trailing underscores (e.g. _version_) are reserved.  
-->  
  
<field name="id" type="string" indexed="true" stored="true" required="true" multiValued="false" />  
<field name="sku" type="text_en_splitting_tight" indexed="true" stored="true" omitNorms="true"/>  
<field name="name" type="text_general" indexed="true" stored="true"/>  
<field name="manu" type="text_general" indexed="true" stored="true" omitNorms="true"/>  
<field name="cat" type="string" indexed="true" stored="true" multiValued="true"/>  
<field name="features" type="text_general" indexed="true" stored="true" multiValued="true"/>  
<field name="includes" type="text_general" indexed="true" stored="true" termVectors="true" termPositions="true"  
  
<field name="weight" type="float" indexed="true" stored="true"/>  
<field name="price" type="float" indexed="true" stored="true"/>  
<field name="popularity" type="int" indexed="true" stored="true" />  
<field name="inStock" type="boolean" indexed="true" stored="true" />  
  
<field name="store" type="location" indexed="true" stored="true"/>  
  
<!-- Common metadata fields, named specifically to match up with  
SolrCell metadata when parsing rich documents such as Word, PDF.  
Some fields are multiValued only because Tika currently may return
```

The `schema.xml` file cannot be edited from this screen, but it provides easy access to view the file if needed. Editing is done by modifying the file with a text editor. As described in detail in a later section, the `schema.xml` allows you to define the types of data in your content (field types), the fields your documents will be broken into, and any dynamic fields that should be generated based on patterns of field names in the incoming documents. These options are described in the section called [Documents, Fields, and Schema Design](#).

This screen is related to the [Schema Browser Screen](#) screen, in that they both display information from the schema, but the Schema Browser provides a way to drill into the analysis chain and displays linkages between field types, fields, and dynamic field rules.

Config Screen

The Config screen shows you the current `solrconfig.xml` for the core you selected. This screenshot shows the beginning of the Query section of `solrconfig.xml`.

```
<!-- ----->
Query section - these settings control query time things like caches
----- -->

<query>
<!-- Max Boolean Clauses

Maximum number of clauses in each BooleanQuery, an exception
is thrown if exceeded.

** WARNING **

This option actually modifies a global Lucene property that
will affect all SolrCores. If multiple solrconfig.xml files
disagree on this property, the value at any given moment will
be based on the last SolrCore to be initialized.

-->
<maxBooleanClauses>1024</maxBooleanClauses>

<!-- Solr Internal Query Caches

There are two implementations of cache available for Solr,
LRUCache, based on a synchronized LinkedHashMap, and
FastLRUCache, based on a ConcurrentHashMap.

FastLRUCache has faster gets and slower puts in single
threaded operation and thus is generally faster than LRUCache
when the hit ratio of the cache is high (> 75%), and may be
faster under other scenarios on multi-cpu systems.

-->
```

The `solrconfig.xml` file cannot be edited with this screen, so a text editor of some kind must be used. While the `schema.xml` defines the structure of your content, `solrconfig.xml` defines the behavior of Solr as it indexes content and responds to queries. Many of the options defined with `solrconfig.xml` are described throughout the rest of this Guide. In particular, you will want to review these sections:

- [Indexing and Basic Data Operations](#)
- [Searching](#)
- [The Well-Configured Solr Instance](#)

Replication Screen

The Replication screen shows you the current replication state for the named core you have specified. In Solr, replication is for the index only. SolrCloud has supplanted much of this functionality, but if you are still using index replication, you can see the replication state, as shown below:

The screenshot shows the Apache Solr interface for the 'collection1' core. The left sidebar lists various management options: Dashboard, Logging, Core Admin, Java Properties, Thread Dump, Ping, Query, Schema, Config, Replication (which is selected), Analysis, Schema Browser, and Plugins / Stats. The main content area displays the replication status for the master node. A table shows the index details: Master ID is 1349981878137, Version is 3, and Size is 28.75 KB. Below the table, configuration settings are listed: Settings (Master) show replication.enable: checked (green checkmark) and replicateAfter: commit.

In this example, replication is enabled and will be done after each commit. Because this server is the Master, it is showing only the config settings for the master. On the master, you can disable replication by clicking the **Disable Replication** button.

In Solr, the replication is initiated by the slave servers so there is more value by looking at the Replication screen on the slave nodes. This screenshot shows the Replication screen for a slave:

The screenshot shows the Apache Solr admin interface for a collection named "collection1". The left sidebar has a tree view with "collection1" selected. The main content area has three tabs: "Iterations", "Settings", and "Replication". The "Iterations" tab displays a table of replication iterations:

Index	Version	Gen	Size
Master:	0	1	65 bytes
Slave:	0	1	65 bytes

The "Settings" tab shows configuration details:

- master url: http://localhost:7888/solr/collection1
- polling enable:
- (Master):
 - replication enable:
 - replicateAfter: commit
- confFiles:
admin-extra.html, admin-extra.menu-bottom.html, admin-extra.menu-top.html, elevate.xml, LucidStemRules_en.txt, protwords.txt, schema.xml, solrconfig.xml, stopwords.txt, synonyms.txt

You can click the **Refresh Status** button to show the most current replication status, or choose to get a new snapshot from the master server.

More details on how to configure replication is available in the section called [Index Replication](#).

Analysis Screen

The Analysis screen lets you inspect how data will be handled according to the field, field type and dynamic rule configurations found in schema.xml. You can analyze how content would be handled during indexing or during query processing and view the results separately or at the same time. Ideally, you would want content to be handled consistently, and this screen allows you to validate the settings in the field type or field analysis chains.

Enter content in one or both boxes at the top of the screen, and then choose the field or field type definitions to use for analysis.

The standard output (shown if "Verbose Output" is not checked) will display the step of the analysis and the output based on the current settings. If you click the **Verbose Output** check box, you see more information, including transformations to the input (such as, convert to lower case, strip extra characters, etc.) and the bytes, type and detailed position information.

The information displayed will vary depending on the settings of the field or field type. Each step of the process is displayed in a separate section, with an abbreviation for the tokenizer or filter that is applied in that step. Hover or click on the abbreviation, and you'll see the name and path of the tokenizer or filter.

In the examples on the right (click either screenshot for a larger image), several transformations are applied to the text string "Running is a sport." We've used the field "text", which has rules that remove the "is" and "a" and the word "running" has been changed to its basic form, "run". This is because we have defined the field type, `text_en` in this scenario, to remove stop words (small words that usually do not provide a great deal of context) and "stem" terms when possible to find more possible matches (this is particularly helpful with plural forms of words). If you click the question mark next to the **Analyze Fieldname/Field Type** pull-down menu, the Schema Browser window will open, showing you the settings for the field specified.

The section [Understanding Analyzers, Tokenizers, and Filters](#) describes in detail what each option is and how it may transform your data.

Schema Browser Screen

The Schema Browser screen lets you see schema data in a browser window. If you have accessed this window from the Analysis screen, it will be opened to a specific field, dynamic field rule or field type. If there is nothing chosen, use the pull-down menu to choose the field or field type.

The screenshot shows the Apache Solr Schema Browser interface. On the left is a navigation sidebar with links like Dashboard, Logging, Core Admin, Java Properties, Thread Dump, collection1 (Ping, Query, Schema, Config, Replication, Analysis, Schema Browser, Plugins / Stats), Unique Key Field (id), and Default Search Field. The main area has a search bar at the top with 'text' selected. Below it, a table shows the 'Field' column with 'text' and the 'Copied from' column with 'author', 'cat', 'content', 'content_type', 'description', 'features', 'includes', 'keywords', 'manu', 'name', 'resourcename', 'title', 'url'. To the right of the table, detailed information for the 'text' field is shown: Field-Type: org.apache.solr.schema.TextField, Properties: Indexed, Tokenized, Multivalued, Schema: Indexed, Tokenized, Multivalued, Index: (unstored field), PI Gap: 100, Docs: 21. Below this, two sections are visible: 'Index Analyzer: org.apache.solr.analysis.TokenizerChain' and 'Query Analyzer: org.apache.solr.analysis.TokenizerChain'. A 'Load Term Info' button is present, followed by a list of top terms: 10 / 398 Top-Terms: electronics (14), inc (8), and (6), usb (5), lcd (4), notes (2.0), memory (2), one (1), x (1). An 'Autoload' checkbox is checked. A histogram titled 'Histogram:' shows the frequency distribution of terms across bins: 1 (296), 2 (65), 4 (27), 8 (9), 16 (1).

The screen provides a great deal of useful information about each particular field. In the example above, we have chosen the `text` field. On the right side of the center window, we see the field name, and a list of fields that populate this field because they are defined to be copied to the `text` field. Click on one of those field names, and you can see the definitions for that field. We can also see the field type, which would allow us to inspect the type definitions as well.

In the left part of the center window, we see the field type again, and the defined properties for the field. We can also see how many documents have populated this field. Then we see the analyzer used for indexing and query processing. Click the icon to the left of either of those, and you'll see the definitions for the tokenizers and/or filters that are used. The output of these processes is the information you see when testing how content is handled for a particular field with the [Analysis Screen](#).

Under the analyzer information is a button to **Load Term Info**. Clicking that button will show the top N terms that are in the index for that field. Click on a term, and you will be taken to the [Query Screen](#) to see the results of a query of that term in that field. If you want to always see the term information for a field, choose **Autoload** and it will always appear when there are terms for a field. A histogram shows the number of terms with a given frequency in the field.

Plugins & Stats Screen

The Plugins screen shows information and statistics about Solr's status and performance. You can find information about the performance of Solr's caches, the state of Solr's searchers, and the configuration of searchHandlers and requestHandlers.

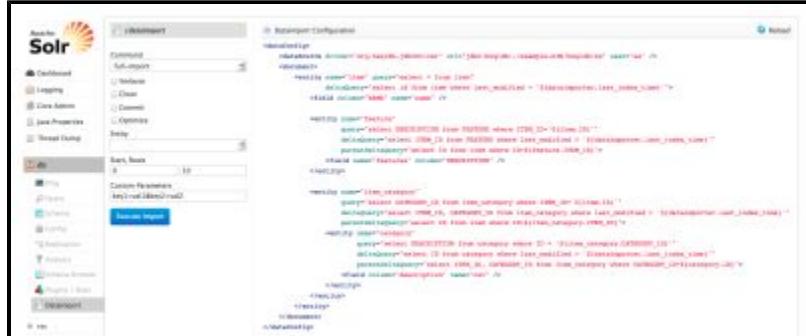
Choose an area of interest on the right, and then drill down into more specifics by clicking on one of the names that appear in the central part of the window. In this example, we've chosen to look at the Searcher stats, from the Core area:

The screenshot shows the Apache Solr interface. On the left is a sidebar with various administrative links: Dashboard, Logging, Core Admin, Java Properties, Thread Dump, collection1 (with sub-links: Ping, Query, Schema, Config, Replication, Analysis, Schema Browser), and Plugins / Stats. The Plugins / Stats link is highlighted. The main content area has a tree view on the left with nodes: CACHE, CORE (which is selected and highlighted in grey), HIGHLIGHTING, OTHER, QUERYHANDLER, UPDATEHANDLER, Watch Changes, and Refresh Values. The CORE node is expanded, showing a detailed view of a Searcher named 'Searcher@e2942da main'. This view includes fields like class, version, description, src, stats, core, and searcher. The 'stats' section provides specific metrics such as searcherName, caching, numDocs, maxDoc, reader, readerDir, indexVersion, openedAt, registeredAt, and warmupTime. The 'core' and 'searcher' sections are also visible but collapsed.

Searcher Statistics

The display is a snapshot taken when the page is loaded. You can get updated status by choosing to either **Watch Changes** or **Refresh Values**. Watching the changes will highlight those areas that have changed, while refreshing the values will reload the page with updated information.

Dataimport Screen



The Dataimport screen shows the configuration of the DataImportHandler (DIH) and allows you to start indexing data, as defined by the options selected on the screen and defined in the configuration file. Click the screenshot on the right to see a larger image of this screen.

The configuration file defines the location of the data and how to perform the SQL queries for the data you want. The options on the screen control how the data is imported to Solr. For more information about data importing with DIH, see the section on [Uploading Structured Data Store Data with the Data Import Handler](#).

Documents, Fields, and Schema Design

This section discusses how Solr organizes its data into documents and fields, as well as how to work with the Solr schema file, `schema.xml`. It includes the following topics:

[Overview of Documents, Fields, and Schema Design](#): An introduction to the concepts covered in this section.

[Solr Field Types](#): Detailed information about field types in Solr, including the field types in the default Solr schema.

[Defining Fields](#): Describes how to define fields in Solr.

[Copying Fields](#): Describes how to populate fields with data copied from another field.

[Dynamic Fields](#): Information about using dynamic fields in order to catch and index fields that do not exactly conform to other field definitions in your schema.

[Other Schema Elements](#): Describes other important elements in the Solr schema: Unique Key, Default Search Field, and the Query Parser Operator.

[Putting the Pieces Together](#): A higher-level view of the Solr schema and how its elements work together.

Overview of Documents, Fields, and Schema Design

The fundamental premise of Solr is simple. You give it a lot of information, then later you can ask it questions and find the piece of information you want. The part where you feed in all the information is called *indexing* or *updating*. When you ask a question, it's called a *query*.

One way to understand how Solr works is to think of a loose-leaf book of recipes. Every time you add a recipe to the book, you update the index at the back. You list each ingredient and the page number of the recipe you just added. Suppose you add one hundred recipes. Using the index, you can very quickly find all the recipes that use garbanzo beans, or artichokes, or coffee, as an ingredient. Using the index is much faster than looking through each recipe one by one. Imagine a book of one thousand recipes, or one million.

Solr allows you to build an index with many different fields, or types of entries. The example above shows how to build an index with just one field, ingredients. You could have other fields in the index for the recipe's cooking style, like Asian, Cajun, or vegan, and you could have an index field for preparation times. Solr can answer questions like "What Cajun-style recipes that have blood oranges as an ingredient can be prepared in fewer than 30 minutes?"

The schema is the place where you tell Solr how it should build indexes from input documents.

How Solr Sees the World

Solr's basic unit of information is a *document*, which is a set of data that describes something. A recipe document would contain the ingredients, the instructions, the preparation time, the cooking time, the tools needed, and so on. A document about a person, for example, might contain the person's name, biography, favorite color, and shoe size. A document about a book could contain the title, author, year of publication, number of pages, and so on.

In the Solr universe, documents are composed of *fields*, which are more specific pieces of information. Shoe size could be a field. First name and last name could be fields.

Fields can contain different kinds of data. A name field, for example, is text (character data). A shoe size field might be a floating point number so that it could contain values like 6 and 9.5. Obviously, the definition of fields is flexible (you could define a shoe size field as a text field rather than a floating point number, for example), but if you define your fields correctly, Solr will be able to interpret them correctly and your users will get better results when they perform a query.

You can tell Solr about the kind of data a field contains by specifying its *field type*. The field type tells Solr how to interpret the field and how it can be queried.

When you add a document, Solr takes the information in the document's fields and adds that information to an index. When you perform a query, Solr can quickly consult the index and return the matching documents.

Field Analysis

Field analysis tells Solr what to do with incoming data when building an index. A more accurate name for this process would be *processing* or even *digestion*, but the official name is *analysis*.

Consider, for example, a biography field in a person document. Every word of the biography must be indexed so that you can quickly find people whose lives have had anything to do with ketchup, or dragonflies, or cryptography.

However, a biography will likely contain lots of words you don't care about and don't want clogging up your index—words like "the," "a," "to," and so forth. Furthermore, suppose the biography contains the word "Ketchup," capitalized at the beginning of a sentence. If a user makes a query for "ketchup," you want Solr to tell you about the person even though the biography contains the capitalized word.

The solution to both these problems is field analysis. For the biography field, you can tell Solr how to break apart the biography into words. You can tell Solr that you want to make all the words lower case, and you can tell Solr to remove accents marks.

Field analysis is an important part of a field type. [Understanding Analyzers, Tokenizers, and Filters](#) is a detailed description of field analysis.

Solr Field Types

The field type defines how Solr should interpret data in a field and how the field can be queried. There are many field types included with Solr by default, and they can also be defined locally.

Topics covered in this section:

- [Field Type Definitions and Properties](#)
- [Field Types Included with Solr](#)
- [Working with Currencies and Exchange Rates](#)
- [Working with Dates](#)
- [Working with External Files and Processes](#)
- [Field Properties by Use Case](#)

Related Topics

- [SchemaXML-DataTypes](#)

- [FieldType Javadoc](#)


Field Type Definitions and Properties

A field type includes four types of information:

- The name of the field type
- An implementation class name
- If the field type is `TextField`, a description of the field analysis for the field type
- Field attributes

Field Type Definitions in schema.xml

Field types are defined in `schema.xml`, with the `types` element. Each field type is defined between `fieldType` elements. Here is an example of a field type definition for a type called `text_general`:

```
<fieldType name="text_general" class="solr.TextField" positionIncrementGap="100">
    <analyzer type="index">
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.StopFilterFactory" ignoreCase="true" words="stopwords.txt"
enablePositionIncrements="true" />
        <!-- in this example, we will only use synonyms at query time
        <filter class="solr.SynonymFilterFactory" synonyms="index_synonyms.txt"
ignoreCase="true" expand="false"/>
        -->
        <filter class="solr.LowerCaseFilterFactory"/>
    </analyzer>
    <analyzer type="query">
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.StopFilterFactory" ignoreCase="true" words="stopwords.txt"
enablePositionIncrements="true" />
        <filter class="solr.SynonymFilterFactory" synonyms="synonyms.txt"
ignoreCase="true" expand="true"/>
        <filter class="solr.LowerCaseFilterFactory"/>
    </analyzer>
</fieldType>
```

The first line in the example above contains the field type name, `text_general`, and the name of the implementing class, `solr.TextField`. The rest of the definition is about field analysis, described in [Understanding Analyzers, Tokenizers, and Filters](#).

The implementing class is responsible for making sure the field is handled correctly. In the class names in `schema.xml`, the string `solr` is shorthand for `org.apache.solr.schema` or `org.apache.solr.analysis`. Therefore, `solr.TextField` is really `org.apache.solr.schema.TextField`.

Field Type Properties

The field type `class` determines most of the behavior of a field type, but optional properties can also be defined. For example, the following definition of a date field type defines two properties, `sortMissingLast` and `omitNorms`.

```
<fieldType name="date" class="solr.DateField"
    sortMissingLast="true" omitNorms="true"/>
```

Here are some commonly used properties. Most properties are either true or false. In addition, the filters or tokenizers defined for field analysis may have properties that can be defined.

Field Property	Description	Values
<code>indexed</code>	If true, the value of the field can be used in queries to retrieve matching documents	true or false
<code>stored</code>	If true, the actual value of the field can be retrieved by queries	true or false
<code>sortMissingFirst</code> <code>sortMissingLast</code>	Control the placement of documents when a sort field is not present. As of Solr 3.5, these work for all numeric fields, including Trie and date fields.	true or false
<code>multiValued</code>	If true, indicates that a single document might contain multiple values for this field type	true or false
<code>positionIncrementGap</code>	For multivalued fields, specifies a distance between multiple values, which prevents spurious phrase matches	integer
<code>omitNorms</code>	If true, omits the norms associated with this field (this disables length normalization and index-time boosting for the field, and saves some memory). Defaults to true for all primitive (non-analyzed) field types, such as int, float, data, bool, and string. Only full-text fields or fields that need an index-time boost need norms.	true or false
<code>omitTermFreqAndPositions</code>	If true, omits term frequency, positions, and payloads from postings for this field. This can be a performance boost for fields that don't require that information. It also reduces the storage space required for the index. Queries that rely on position that are issued on a field with this option will silently fail to find documents. This property defaults to true for all fields that are not text fields.	true or false

autoGeneratePhraseQueries	For text fields. If true, Solr automatically generates phrase queries for adjacent terms. If false, terms must be enclosed in double-quotes to be treated as phrases.	true or false
postingsFormat	Defines a custom codec to use. The codec must be schema-aware, such as the <code>SchemaCodecFactory</code> and must exist in any Solr <code>lib</code> directories.	n/a

Field Types Included with Solr

The following table lists the field types that are available in Solr. The `org.apache.solr.schema` package includes all the classes listed in this table.

Class	Description
BCDIntField	Binary-coded decimal (BCD) integer. BCD is a relatively inefficient encoding that offers the benefits of quick decimal calculations and quick conversion to a string.
BCDLongField	Binary-coded decimal long integer.
BCDStrField	Binary-coded decimal string.
BinaryField	Binary data.
BoolField	Contains either true or false. Values of "1", "t", or "T" in the first character are interpreted as true. Any other values in the first character are interpreted as false.
ByteField	Contains an array of bytes.
CurrencyField	Supports currencies and exchange rates. See the section Working with Currencies and Exchange Rates .
DateField	Represents a point in time with millisecond precision. See the section Working with Dates .
DoubleField	Double (64-bit IEEE floating point).
ExternalFileField	Pulls values from a file on disk. See the section Working with External Files and Processes .
FloatField	Floating point (32-bit IEEE floating point).
IntField	Integer (32-bit signed integer).
LatLonType	For spatial search: a latitude/longitude coordinate pair. The latitude is specified first in the pair.
LongField	Long integer (64-bit signed integer).
PointType	For spatial search: An arbitrary n-dimensional point, useful for searching sources such as blueprints or CAD drawings.

PreAnalyzedField	Provides a way to send to Solr serialized token streams, optionally with independent stored values of a field, and have this information stored and indexed without any additional text processing. Useful if you want to submit field content that was already processed by some existing external text processing pipeline (e.g. tokenized, annotated, stemmed, inserted synonyms, etc.), while using all the rich attributes that Lucene's TokenStream provides via token attributes.
RandomSortField	Does not contain a value. Queries that sort on this field type will return results in random order. Use a dynamic field to use this feature.
ShortField	Short integer.
SortableDoubleField	The Sortable fields provide correct numeric sorting. If you use the plain types (DoubleField, IntField, and so on) sorting will be lexicographical instead of numeric.
SortableFloatField	Numerically sorted floating point.
SortableIntField	Numerically sorted integer.
SortableLongField	Numerically sorted long integer.
StrField	String (UTF-8 encoded string or Unicode).
TextField	Text, usually multiple words or tokens.
TrieDateField	Date field accessible for Lucene TrieRange processing.
TrieDoubleField	Double field accessible Lucene TrieRange processing.
TrieField	If this type is used, a "type" attribute must also be specified, with a value of either: integer, long, float, double, date. Using this field is the same as using any of the Trie fields.
TrieFloatField	Floating point field accessible Lucene TrieRange processing.
TrieIntField	Int field accessible Lucene TrieRange processing.
TrieLongField	Long field accessible Lucene TrieRange processing.
UUIDField	Universally Unique Identifier (UUID). Pass in a value of "NEW" and Solr will create a new UUID.

The `MultiTermAwareComponent` has been added to relevant `solr.TextField` entries in `schema.xml` (e.g., wildcards, regex, prefix, range, etc.) to allow automatic lowercasing for multi-term queries.

Further, you can now optionally specify `analyzerType="multiterm"` in `schema.xml`; if you don't, `analyzer` will process the fields according to their specific attributes.

Working with Currencies and Exchange Rates

The `currency` FieldType provides support for monetary values to Solr/Lucene with query-time currency conversion and exchange rates. The following features are supported:

- Point queries
- Range queries
- Sorting
- Currency parsing by either currency code or symbol
- Symmetric & asymmetric exchange rates (asymmetric exchange rates are useful if there are fees associated with exchanging the currency)

Configuring Currencies

The `currency` field type is defined in `schema.xml`. This is the default configuration of this type:

```
<fieldType name="currency" class="solr.CurrencyField" precisionStep="8"  
defaultCurrency="USD" currencyConfig="currency.xml" />
```

In this example, we have defined the name and class of the field type, and defined the `defaultCurrency` as "USD", for U.S. Dollars. We have also defined a `currencyConfig` to use a file called "currency.xml". This is a file of exchange rates between our default currency to other currencies. There is an alternate implementation that would allow regular downloading of currency data. See [Exchange Rates](#) below for more.

At indexing time, money fields can be indexed in a native currency. For example, if a product on an e-commerce site is listed in Euros, indexing the price field as "1000,EUR" will index it appropriately. The price should be separated from the currency by a comma, and the price must be encoded with a floating point value (a decimal point).

During query processing, range and point queries are both supported.

Exchange Rates

You configure exchange rates by specifying a provider. Natively, two provider types are supported: `FileExchangeRateProvider` or `openExchangeRatesOrgProvider`.

FileExchangeRateProvider

This provider requires you to provide a file of exchange rates. It is the default, meaning that to use this provider you only need to specify the file path and name as a value for `currencyConfig` in the definition for this type.

There is a sample `currency.xml` file included with Solr, found in the same directory as the `schema.xml` file. Here is a small snippet from this file:

```
<currencyConfig version="1.0">
  <rates>
    <!-- Updated from http://www.exchangerate.com/ at 2011-09-27 -->
    <rate from="USD" to="ARS" rate="4.333871" comment="ARGENTINA Peso" />
    <rate from="USD" to="AUD" rate="1.025768" comment="AUSTRALIA Dollar" />
    <rate from="USD" to="EUR" rate="0.743676" comment="European Euro" />
    <rate from="USD" to="CAD" rate="1.030815" comment="CANADA Dollar" />

    <!-- Cross-rates for some common currencies -->
    <rate from="EUR" to="GBP" rate="0.869914" />
    <rate from="EUR" to="NOK" rate="7.800095" />
    <rate from="GBP" to="NOK" rate="8.966508" />

    <!-- Asymmetrical rates -->
    <rate from="EUR" to="USD" rate="0.5" />
  </rates>
</currencyConfig>
```

OpenExchangeRatesOrgProvider

With Solr 4, you can configure Solr to download exchange rates from [OpenExchangeRates.Org](#), with updates rates between USD and 158 currencies hourly. These rates are symmetrical only.

In this case, you need to specify the `providerClass` in the definitions for the field type. Here is an example:

```
<fieldType name="currency" class="solr.CurrencyField" precisionStep="8"
  providerClass="solr.OpenExchangeRatesOrgProvider"
  refreshInterval="60"
  ratesFileLocation="http://internal.server/rates.json"/>
```

The `refreshInterval` is minutes, so the above example will download the newest rates every 60 minutes.

Working with Dates

`DateField` represents a point in time with millisecond precision. The format is:

`YYYY-MM-DDThh:mm:ssZ`

`YYYY` is the year.

`MM` is the month.

`DD` is the day of the month.

`hh` is the hour of the day as on a 24-hour clock.

`mm` is minutes.

`ss` is seconds.

Note that no time zone can be specified; the time given should be expressed in Coordinated Universal Time (UTC). Here is an example value:

`1972-05-20T17:33:18Z`

You can include fractional seconds if you wish, although trailing zeros are not allowed and any precision beyond milliseconds will be ignored. Here is another example value with milliseconds included:

`1972-05-20T17:33:18.772Z`

In addition, `DateField` also supports *date math*. This makes it easy to create times relative to the current time. This represents a point in time two months from now:

`+2MONTHS`

This is one day ago:

`-1DAY`

Use a slash to indicate rounding. This represents the beginning of the current hour:

`/HOUR`

You can combine terms. The following is six months and three days in the future, at the beginning of the day:

`+6MONTHS+3DAYS/DAY`

Working with External Files and Processes

The ExternalFileField Type

The `ExternalFileField` type makes it possible to specify the values for a field in a file outside the Solr index. For such a field, the file contains mappings from a key field to the field value. Another way to think of this is that, instead of specifying the field in documents as they are indexed, Solr finds values for this field in the external file.

 External fields are not searchable. They can be used only for function queries or display. For more information on function queries, see the section on [Function Queries](#).

The `ExternalFileField` type is handy for cases where you want to update a particular field in many documents more often than you want to update the rest of the documents. For example, suppose you have implemented a document rank based on the number of views. You might want to update the rank of all the documents daily or hourly, while the rest of the contents of the documents might be updated much less frequently. Without `ExternalFileField`, you would need to update each document just to change the rank. Using `ExternalFileField` is much more efficient because all document values for a particular field are stored in an external file that can be updated as frequently as you wish.

In `schema.xml`, the definition of this field type might look like this:

```
<fieldType name="entryRankFile" keyField="pkId" defVal="0" stored="false"
indexed="false"
class="solr.ExternalFileField" valType="pfloa"/>
```

The `keyField` attribute defines the key that will be defined in the external file. It is usually the unique key for the index, but it doesn't need to be as long as the `keyField` can be used to identify documents in the index. A `defVal` defines a default value that will be used if there is no entry in the external file for a particular document.

The `valType` attribute specifies the actual type of values that will be found in the file. The type specified must be either a float field type, so valid values for this attribute are `pfloa`, `float` or `tfloa`. This attribute can be omitted.

Format of the External File

The file itself is located in Solr's index directory, which by default is `$SOLR_HOME/data`. The name of the file should be `external_fieldname` or `external_fieldname.*`. For the example above, then, the file could be named `external_entryRankFile` or `external_entryRankFile.txt`.

- ✓ If any files using the name pattern `.*` (such as `.txt`) appear, the last (after being sorted by name) will be used and previous versions will be deleted. This behavior supports implementations on systems where one may not be able to overwrite a file (for example, on Windows, if the file is in use).

The file contains entries that map a key field, on the left of the equals sign, to a value, on the right. Here are a few example entries:

```
doc33=1.414  
doc34=3.14159  
doc40=42
```

The keys listed in this file do not need to be unique. The file does not need to be sorted, but Solr will be able to perform the lookup faster if it is.

Reloading an External File

As of Solr 4.1, it's possible to define an event listener to reload an external file when either a searcher is reloaded or when a new searcher is started. See the section [Query Related Listeners](#) for more information, but a sample definition in `solrconfig.xml` might look like this:

```
<listener event="newSearcher"  
class="org.apache.solr.schema.ExternalFileFieldReloaded" />  
<listener event="firstSearcher"  
class="org.apache.solr.schema.ExternalFileFieldReloaded" />
```

Pre-Analyzing a Field Type

The `PreAnalyzedField` type provides a way to send to Solr serialized token streams, optionally with independent stored values of a field, and have this information stored and indexed without any additional text processing applied in Solr. This is useful if user wants to submit field content that was already processed by some existing external text processing pipeline (e.g., it has been tokenized, annotated, stemmed, synonyms inserted, etc), while using all the rich attributes that Lucene's [TokenStream](#) provides (per-token attributes).

The serialization format is pluggable using implementations of [PreAnalyzedParser](#) interface. There are two out-of-the-box implementations:

- [JsonPreAnalyzedParser](#): as the name suggests, it parses content that uses JSON to represent field's content. This is the default parser to use if the field type is not configured otherwise.
- [SimplePreAnalyzedParser](#): uses a simple strict plain text format, which in some situations may be easier to create than JSON.

There is only one configuration parameter, `parserImpl`. The value of this parameter should be a fully qualified class name of a class that implements [PreAnalyzedParser](#) interface. The default value of this parameter is `org.apache.solr.schema.JsonPreAnalyzedParser`.

Field Properties by Use Case

Here is a summary of common use cases, and the attributes the fields or field types should have to support the case. An entry of true or false in the table indicates that the option must be set to the given value for the use case to function correctly. If no entry is provided, the setting of that attribute has no impact on the case.

Use Case	indexed	stored	multiValued	omitNorms	termVectors	termPositions
search within field	true					
retrieve contents		true				
use as unique key	true		false			
sort on field	true		false	true ¹		
use field boosts ⁵				false		
document boosts affect searches within field				false		
highlighting	true ⁴	true			2	true ³
faceting ⁵	true					
add multiple values, maintaining order			true			
field length affects doc score				false		
MoreLikeThis ⁵					true ⁶	

Notes:

¹ Recommended but not necessary.

² Will be used if present, but not necessary.

³ (if termVectors=true)

⁴ A tokenizer must be defined for the field, but it doesn't need to be indexed.

⁵ Described in [Understanding Analyzers, Tokenizers, and Filters](#).

⁶ Term vectors are not mandatory here. If not true, then a stored field is analyzed. So term vectors are recommended, but only required if stored=false.

Defining Fields

Once you have the field types set up, defining the fields themselves is simple. All you do is supply a name and a field type. If you wish, you can also provide options that will override the options for the field type.

Fields are defined in the `fields` element of `schema.xml`. The following example defines a field named `price` with a type of `sfloat`.

```
<field name="price" type="sfloat" indexed="true" stored="true" />
```

Fields can have the same options as field types. The field type options serve as defaults which can be overridden by options defined per field. Included below is the table of field type properties from the section [Field Type Definitions and Properties](#):

Field Property	Description	Values
indexed	If true, the value of the field can be used in queries to retrieve matching documents	true or false
stored	If true, the actual value of the field can be retrieved by queries	true or false
sortMissingFirst sortMissingLast	Control the placement of documents when a sort field is not present. As of Solr 3.5, these work for all numeric fields, including Trie and date fields.	true or false
multiValued	If true, indicates that a single document might contain multiple values for this field type	true or false
positionIncrementGap	For multivalued fields, specifies a distance between multiple values, which prevents spurious phrase matches	integer
omitNorms	If true, omits the norms associated with this field (this disables length normalization and index-time boosting for the field, and saves some memory). Defaults to true for all primitive (non-analyzed) field types, such as int, float, data, bool, and string. Only full-text fields or fields that need an index-time boost need norms.	true or false

omitTermFreqAndPositions	If true, omits term frequency, positions, and payloads from postings for this field. This can be a performance boost for fields that don't require that information. It also reduces the storage space required for the index. Queries that rely on position that are issued on a field with this option will silently fail to find documents. This property defaults to true for all fields that are not text fields.	true or false
autoGeneratePhraseQueries	For text fields. If true, Solr automatically generates phrase queries for adjacent terms. If false, terms must be enclosed in double-quotes to be treated as phrases.	true or false
postingsFormat	Defines a custom codec to use. The codec must be schema-aware, such as the <code>SchemaCodecFactory</code> and must exist in any Solr <code>lib</code> directories.	n/a

Related Topics

- [SchemaXML-Fields](#)

- [Field Options by Use Case](#)


Copying Fields

You might want to interpret some document fields in more than one way. Solr has a mechanism for making copies of fields so that you can apply several distinct field types to a single piece of incoming information.

The name of the field you want to copy is the *source*, and the name of the copy is the *destination*. In schema.xml, it's very simple to make copies of fields:

```
<copyField source="cat" dest="text" maxChars="30000" />
```

If the text field has data of its own in input documents, the contents of cat will be added to the index for text. The maxChars parameter, an int parameter, establishes an upper limit for the number of characters to be copied. This limit is useful for situations in which you want to control the size of index files.

Both the source and the destination of copyField can contain asterisks, which will match anything. For example, the following line will copy the contents of all incoming fields that match the wildcard pattern *t to the text field.:

```
<copyField source="*_t" dest="text" maxChars="25000" />
```

 The copyField command can use a wildcard (*) character in the dest parameter only if the source parameter contains one as well. copyField uses the matching glob from the source field for the dest field name into which the source content is copied.

Related Topics

- [SchemaXML-Copy Fields](#)



Dynamic Fields

Dynamic fields allow Solr to index fields that you did not explicitly define in your schema. This is useful if you discover you have forgotten to define one or more fields. Dynamic fields can make your application less brittle by providing some flexibility in the documents you can add to Solr.

A dynamic field is just like a regular field except it has a name with a wildcard in it. When you are indexing documents, a field that does not match any explicitly defined fields can be matched with a dynamic field.

For example, suppose your schema includes a dynamic field with a name of `*_i`. If you attempt to index a document with a `cost_i` field, but no explicit `cost_i` field is defined in the schema, then the `cost_i` field will have the field type and analysis defined for `*_i`.

Dynamic fields are also defined in the `fields` element of `schema.xml`. Like fields, they have a name, a field type, and options.

```
<dynamicField name="*_i" type="int" indexed="true" stored="true" />
```

LucidWorks recommends that you include basic dynamic field mappings (like that shown above) in your `schema.xml`. The mappings can be very useful.

Related Topics

- [SchemaXML-Dyanmic Fields](#)



Other Schema Elements

This section describes several other important elements of `schema.xml`.

Unique Key

The `uniqueKey` element specifies which field is a unique identifier for documents. Although `uniqueKey` is not required, it is nearly always warranted by your application design. For example, `uniqueKey` should be used if you will ever update a document in the index.

You can define the unique key field by naming it:

```
<uniqueKey>id</uniqueKey>
```

Starting with Solr 4, schema defaults and `copyFields` cannot be used to populate the `uniqueKey` field. You also can't use `UUIDUpdateProcessorFactory` to have `uniqueKey` values generated automatically.

Further, the operation will fail if the `uniqueKey` field is used, but is multivalued (or inherits the multivalueness from the `fieldtype`). However, `uniqueKey` will continue to work, as long as the field is properly used.

Default Search Field

If you are using the Lucene query parser, queries that don't specify a field name will use the `defaultSearchField`. The DisMax and Extended DisMax query parsers do not use this value.

- ⌚ Use of the `defaultSearchKey` is deprecated in Solr versions 3.6 and higher. Instead, you should use the `df` parameter. At some point, the `defaultSearchKey` parameter may be removed.

For more information about query parsers, see the section on [Query Syntax and Parsing](#).

Query Parser Default Operator

In queries with multiple terms, Solr can either return results where all conditions are met or where one or more conditions are met. The *operator* controls this behavior. An operator of AND means that all conditions must be fulfilled, while an operator of OR means that one or more conditions must be true.

In `schema.xml`, the `solrQueryParser` element controls what operator is used if an operator is not specified in the query. The default operator setting only applies to the Lucene query parser, not the DisMax or Extended DisMax query parsers, which internally hard-code their operators to OR.

- ➥ The query parser default operator parameter has been deprecated in Solr versions 3.6 and higher. You are instead encouraged to specify the query parser `q.op` parameter in your request handler.

Similarity

Similarity is a Lucene class used to score a document in searching. This class can be changed in order to provide a more custom sorting. With Solr 4, you can configure a different `similarity` for each field, meaning that scoring a document will differ depending on what's in each field. However, you can still configure a global `similarity` is configured in the `schema.xml` file, where an implicit instance of `DefaultSimilarityFactory` is used.

A global `<similarity>` declaration can be used to specify a custom similarity implementation that you want Solr to use when dealing with your index. A similarity can be specified either by referring directly to the name of a class with a no-argument constructor:

```
<similarity/>...
```

or by referencing a `SimilarityFactory` implementation, which may take optional initialization parameters:

```
<similarity>
  <str name="basicModel">P</str>
  <str name="afterEffect">L</str>
  <str name="normalization">H2</str>
  <float name="c">7</float>
</similarity>
```

Beginning with Solr 4, Similarity factories such as `SchemaSimilarityFactory` can also support specifying specific `similarity` implementations on individual field types:

```
<fieldType name="text_ib">
  <analyzer/>
  <similarity>
    <str name="distribution">SPL</str>
    <str name="lambda">DF</str>
    <str name="normalization">H2</str>
  </similarity>
</fieldType>
```

Related Topics

- [SchemaXML-Miscellaneous Settings](#)

- [UniqueKey](#)


Putting the Pieces Together

At the highest level, `schema.xml` is structured as follows. This example is not real XML, but it gives you an idea of the structure of the file.

```
<schema>
  <types>
  <fields>
    <uniqueKey>
    <defaultSearchField>
    <solrQueryParser defaultOperator>
      <copyField>
    </schema>
```

Obviously, most of the excitement is in types and fields, where the field types and the actual field definitions live. These are supplemented by `copyFields`. Sandwiched between fields and the `copyField` section are the unique key, default search field, and the default query operator.

Choosing Appropriate Numeric Types

For general numeric needs, use the sortable field types, `SortableIntField`, `SortableLongField`, `SortableFloatField`, and `SortableDoubleField`. These field types will sort numerically instead of lexicographically, which is the main reason they are preferable over their simpler cousins, `IntegerField`, `LongField`, `FloatField`, and `DoubleField`.

If you expect users to make frequent range queries on numeric types, consider using `TrieField`. It offers faster speed for range queries at the expense of increasing index size.

Working With Text

Handling text properly will make your users happy by providing them with the best possible results for text searches.

One technique is using a text field as a catch-all for keyword searching. Most users are not sophisticated about their searches and the most common search is likely to be a simple keyword search. You can use `copyField` to take a variety of fields and funnel them all into a single text field for keyword searches. In the example schema representing a store, `copyField` is used to dump the contents of `cat`, `name`, `manu`, `features`, and `includes` into a single field, `text`. In addition, it could be a good idea to copy `ID` into `text` in case users wanted to search for a particular product by passing its product number to a keyword search.

Another technique is using `copyField` to use the same field in different ways. Suppose you have a field that is a list of authors, like this:

Schildt, Herbert; Wolpert, Lewis; Davies, P.

For searching by author, you could tokenize the field, convert to lower case, and strip out punctuation:

```
schildt / herbert / wolpert / lewis / davies / p
```

For sorting, just use an untokenized field, converted to lower case, with punctuation stripped:

```
schildt herbert wolpert lewis davies p
```

Finally, for facetting, use the primary author only via a `StringField`:

Schildt, Herbert

Related Topics

- [SchemaXML](#)


Understanding Analyzers, Tokenizers, and Filters

This sections describes how Solr breaks down and works with textual data. It covers the following topics:

[Overview of Analyzers, Tokenizers, and Filters](#): A conceptual introduction to Solr's analyzers, tokenizers, and filters.

[What Is An Analyzer?](#): Detailed conceptual information about Solr analyzers.

[What Is A Tokenizer?](#): Detailed conceptual information about Solr tokenizers.

[What Is a Filter?](#): Detailed conceptual information about Solr filters.

[Tokenizers](#): Information about configuring tokenizers, and about the tokenizer factory classes included in this distribution of Solr.

[Filter Descriptions](#): Information about configuring filters, and about the filter factory classes included in this distribution of Solr.

[CharFilterFactories](#): Information about filters for pre-processing input characters.

[Language Analysis](#): Information about tokenizers and filters for character set conversion or for use with specific languages.

[Running Your Analyzer](#): Detailed information about testing and running your Solr analyzer.

Overview of Analyzers, Tokenizers, and Filters

[Field analyzers](#) are used both during ingestion, when a document is indexed, and at query time. An analyzer examines the text of fields and generates a token stream. Analyzers may be a single class or they may be composed of a series of tokenizer and filter classes.

[Tokenizers](#) break field data into lexical units, or *tokens*. [Filters](#) examine a stream of tokens and keep them, transform or discard them, or create new ones. Tokenizers and filters may be combined to form pipelines, or *chains*, where the output of one is input to the next. Such a sequence of tokenizers and filters is called an *analyzer* and the resulting output of an analyzer is used to match query results or build indices.

Although the analysis process is used for both indexing and querying, the same analysis process need not be used for both operations. For indexing, you often want to simplify, or normalize, words. For example, setting all letters to lowercase, eliminating punctuation and accents, mapping words to their stems, and so on. Doing so can increase recall because, for example, "ram", "Ram" and "RAM" would all match a query for "ram". To increase query-time precision, a filter could be employed to narrow the matches by, for example, ignoring all-cap acronyms if you're interested in male sheep, but not Random Access Memory.

The tokens output by the analysis process define the values, or *terms*, of that field and are used either to build an index of those terms when a new document is added, or to identify which documents contain the terms your are querying for.

This section will show you how to configure field analyzers and also serves as a reference for the details of configuring each of the available tokenizer and filter classes. It also serves as a guide so that you can configure your own analysis classes if you have special needs that cannot be met with the included filters or tokenizers.

For more information on Solr's analyzers, tokenizers, and filters, see
<http://wiki.apache.org/solr/AnalyzersTokenizersTokenFilters>.

What Is An Analyzer?

An analyzer examines the text of fields and generates a token stream. Analyzers are specified as a child of the `<fieldType>` element in the `schema.xml` configuration file that can be found in the `solr/conf` directory, or wherever `solrconfig.xml` is located.

In normal usage, only fields of type `solr.TextField` will specify an analyzer. The simplest way to configure an analyzer is with a single `<analyzer>` element whose class attribute is a fully qualified Java class name. The named class must derive from `org.apache.lucene.analysis.Analyzer`. For example:

```
<fieldType name="nametext" class="solr.TextField">
  <analyzer class="org.apache.lucene.analysis.WhitespaceAnalyzer"/>
</fieldType>
```

In this case a single class, `WhitespaceAnalyzer`, is responsible for analyzing the content of the named text field and emitting the corresponding tokens. For simple cases, such as plain English prose, a single analyzer class like this may be sufficient. But it's often necessary to do more complex analysis of the field content.

Even the most complex analysis requirements can usually be decomposed into a series of discrete, relatively simple processing steps. As you will soon discover, the Solr distribution comes with a large selection of tokenizers and filters that covers most scenarios you are likely to encounter. Setting up an analyzer chain is very straightforward; you specify a simple `<analyzer>` element (no class attribute) with child elements that name factory classes for the tokenizer and filters to use, in the order you want them to run.

For example:

```
<fieldType name="nametext" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.StandardFilterFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.StopFilterFactory"/>
    <filter class="solr.EnglishPorterFilterFactory"/>
  </analyzer>
</fieldType>
```

Note that classes in the `org.apache.solr.analysis` package may be referred to here with the shorthand `solr.` prefix.

In this case, no Analyzer class was specified on the `<analyzer>` element. Rather, a sequence of more specialized classes are wired together and collectively act as the Analyzer for the field. The text of the field is passed to the first item in the list (`solr.StandardTokenizerFactory`), and the tokens that emerge from the last one (`solr.EnglishPorterFilterFactory`) are the terms that are used for indexing or querying any fields that use the "nametext" fieldType.

Analysis Phases

Analysis takes place in two contexts. At index time, when a field is being created, the token stream that results from analysis is added to an index and defines the set of terms (including positions, sizes, and so on) for the field. At query time, the values being searched for are analyzed and the terms that result are matched against those that are stored in the field's index.

In many cases, the same analysis should be applied to both phases. This is desirable when you want to query for exact string matches, possibly with case-insensitivity, for example. In other cases, you may want to apply slightly different analysis steps during indexing than those used at query time.

If you provide a simple `<analyzer>` definition for a field type, as in the examples above, then it will be used for both indexing and queries. If you want distinct analyzers for each phase, you may include two `<analyzer>` definitions distinguished with a type attribute. For example:

```
<fieldType name="nametext" class="solr.TextField">
    <analyzer *type="index"{*}>
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.LowerCaseFilterFactory"/>
        <filter class="solr.KeepWordFilterFactory" words="keepwords.txt"/>
        <filter class="solr.SynonymFilterFactory" synonyms="syns.txt"/>
    </analyzer>
    <analyzer *type="query"{*}>
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.LowerCaseFilterFactory"/>
    </analyzer>
</fieldType>
```

In this theoretical example, at index time the text is tokenized, the tokens are set to lowercase, any that are not listed in `keepwords.txt` are discarded and those that remain are mapped to alternate values as defined by the synonym rules in the file `syns.txt`. This essentially builds an index from a restricted set of possible values and then normalizes them to values that may not even occur in the original text.

At query time, the only normalization that happens is to convert the query terms to lowercase. The filtering and mapping steps that occur at index time are not applied to the query terms. Queries must then, in this example, be very precise, using only the normalized terms that were stored at index time.

What Is A Tokenizer?

The job of a [tokenizer](#) is to break up a stream of text into tokens, where each token is (usually) a sub-sequence of the characters in the text. An analyzer is aware of the field it is configured for, but a tokenizer is not. Tokenizers read from a character stream (a Reader) and produce a sequence of Token objects (a TokenStream).

Characters in the input stream may be discarded, such as whitespace or other delimiters. They may also be added to or replaced, such as mapping aliases or abbreviations to normalized forms. A token contains various metadata in addition to its text value, such as the location at which the token occurs in the field. Because a tokenizer may produce tokens that diverge from the input text, you should not assume that the text of the token is the same text that occurs in the field, or that its length is the same as the original text. It's also possible for more than one token to have the same position or refer to the same offset in the original text. Keep this in mind if you use token metadata for things like highlighting search results in the field text.

```
<fieldType name="text" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
  </analyzer>
</fieldType>
```

The class named in the tokenizer element is not the actual tokenizer, but rather a class that implements the `org.apache.solr.analysis.TokenizerFactory` interface. This factory class will be called upon to create new tokenizer instances as needed. Objects created by the factory must derive from `org.apache.lucene.analysis.TokenStream`, which indicates that they produce sequences of tokens. If the tokenizer produces tokens that are usable as is, it may be the only component of the analyzer. Otherwise, the tokenizer's output tokens will serve as input to the first filter stage in the pipeline.

A `TypeTokenFilterFactory` is available that creates a `TypeTokenFilter` that filters tokens based on their `TypeAttribute`, which is set in `factory.getStopTypes`.

When To use a CharFilter vs. a TokenFilter

There are several pairs of CharFilters and TokenFilters that have related (ie: `MappingCharFilter` and `ASCIIFoldingFilter`) or nearly identical (ie: `PatternReplaceCharFilterFactory` and `PatternReplaceFilterFactory`) functionality and it may not always be obvious which is the best choice.

The decision about which to use depends largely on which Tokenizer you are using, and whether you need to preprocess the stream of characters.

For example, suppose you have a tokenizer such as `StandardTokenizer` and although you are pretty happy with how it works overall, you want to customize how some specific characters behave. You could modify the rules and re-build your own tokenizer with `javacc`, but it might be easier to simply map some of the characters before tokenization with a `CharFilter`.

TokenizerFactories

Solr provides the following `TokenizerFactories` (`Tokenizers` and `TokenFilters`):

solr.KeywordTokenizerFactory

Creates `org.apache.lucene.analysis.core.KeywordTokenizer`.

Treats the entire field as a single token, regardless of its content. For example:

```
http://example.com/I-am+example?Text=-Hello" ==>  
"http://example.com/I-am+example?Text=-Hello
```

solr.LetterTokenizerFactory

Creates `org.apache.lucene.analysis.LetterTokenizer`.

Creates tokens consisting of strings of contiguous letters. Any non-letter characters will be discarded. For example:

```
"I can't"==>"I", "can", "t"
```

solrWhitespaceTokenizerFactory

Creates `org.apache.lucene.analysis.WhitespaceTokenizer`.

Creates tokens of characters separated by splitting on white space.

solrLowerCaseTokenizerFactory

Creates `org.apache.lucene.analysis.LowerCaseTokenizer`.

Creates tokens by lowercasing all letters and dropping non-letters. For example:

```
"I can't"==>"i", "can", "t"
```

solrStandardTokenizerFactory

Creates `org.apache.lucene.analysis.standard.StandardTokenizer`.

A good general purpose tokenizer that strips many extraneous characters and sets token types to meaningful values. Token types are only useful for subsequent token filters that are type-aware of the same token types. There aren't any filters that use [StandardTokenizer](#)'s types.

What Is a Filter?

Like [tokenizers](#), [filters](#) consume input and produce a stream of tokens. Filters also derive from `org.apache.lucene.analysis.TokenStream`. Unlike tokenizers, a filter's input is another `TokenStream`. The job of a filter is usually easier than that of a tokenizer since in most cases a filter looks at each token in the stream sequentially and decides whether to pass it along, replace it or discard it.

A filter may also do more complex analysis by looking ahead to consider multiple tokens at once, although this is less common. One hypothetical use for such a filter might be to normalize state names that would be tokenized as two words. For example, the single token "california" would be replaced with "CA", while the token pair "rhode" followed by "island" would become the single token "RI".

Because filters consume one `TokenStream` and produce a new `TokenStream`, they can be chained one after another indefinitely. Each filter in the chain in turn processes the tokens produced by its predecessor. The order in which you specify the filters is therefore significant. Typically, the most general filtering is done first, and later filtering stages are more specialized.

```
<fieldType name="text" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.StandardFilterFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.EnglishPorterFilterFactory"/>
  </analyzer>
</fieldType>
```

This example starts with Solr's standard tokenizer, which breaks the field's text into tokens. Those tokens then pass through Solr's standard filter, which removes dots from acronyms, and performs a few other common operations. All the tokens are then set to lowercase, which will facilitate case-insensitive matching at query time.

The last filter in the above example is a stemmer filter that uses the Porter stemming algorithm. A stemmer is basically a set of mapping rules that maps the various forms of a word back to the base, or *stem*, word from which they derive. For example, in English the words "hugs", "hugging" and "hugged" are all forms of the stem word "hug". The stemmer will replace all of these terms with "hug", which is what will be indexed. This means that a query for "hug" will match the term "hugged", but not "huge".

Conversely, applying a stemmer to your query terms will allow queries containing non stem terms, like "hugging", to match documents with different variations of the same stem word, such as "hugged". This works because both the indexer and the query will map to the same stem ("hug").

Word stemming is, obviously, very language specific. Solr includes several language-specific stemmers created by the [Snowball](#) generator that are based on the Porter stemming algorithm. The generic Snowball Porter Stemmer Filter can be used to configure any of these language stemmers. Solr also includes a convenience wrapper for the English Snowball stemmer. There are also several purpose-built stemmers for non-English languages. These stemmers are described in [Language Analysis](#).

Tokenizers

You configure the tokenizer for a text field type in `schema.xml` with a `<tokenizer>` element, as a child of `<analyzer>`:

```
<fieldType name="text" class="solr.TextField">
    <analyzer type="index">
        <tokenizer
            class="solr.StandardTokenizerFactory"/>
        <filter class="solr.StandardFilterFactory"/>
    </analyzer>
</fieldType>
```

The `class` attribute names a factory class that will instantiate a `Tokenizer` object when needed. `Tokenizer` factory classes implement the `org.apache.solr.analysis.TokenizerFactory`. A `TokenizerFactory`'s `create()` method accepts a `Reader` and returns a `TokenStream`. When Solr creates the tokenizer it passes a `Reader` object that provides the content of the text field.

Arguments may be passed to tokenizer factories by setting attributes on the `<tokenizer>` element.

```
<fieldType name="semicolonDelimited"
    class="solr.TextField">
    <analyzer type="query">
        <tokenizer class="solr.PatternTokenizerFactory"
            pattern="; "/>
    </analyzer>
</fieldType>
```

The following sections describe the tokenizer factory classes included in this release of Solr.

For more information about Solr's tokenizers, see
<http://wiki.apache.org/solr/AnalyzersTokenizersTokenFilters>.

Tokenizers discussed in this section:

- [Standard Tokenizer](#)
- [Classic Tokenizer](#)
- [Keyword Tokenizer](#)
- [Letter Tokenizer](#)
- [Lower Case Tokenizer](#)
- [N-Gram Tokenizer](#)
- [Edge N-Gram Tokenizer](#)
- [ICU Tokenizer](#)
- [Path Hierarchy Tokenizer](#)
- [Regular Expression Pattern Tokenizer](#)
- [Type Tokenizer](#)
- [UAX29 URL Email Tokenizer](#)
- [White Space Tokenizer](#)
- [Related Topics](#)

Standard Tokenizer

This tokenizer splits the text field into tokens, treating whitespace and punctuation as delimiters. Delimiter characters are discarded, with the following exceptions:

- Periods (dots) that are not followed by whitespace are kept as part of the token.
- Words are split at hyphens, unless there is a number in the word, in which case the token is not split and the numbers and hyphen(s) are preserved.
- Recognizes Internet domain names and email addresses and preserves them as a single token.

The Standard Tokenizer supports [Unicode standard annex UAX#29](#) word boundaries with the following token types: <ALPHANUM>, <NUM>, <SOUTHEAST_ASIAN>, <IDEOGRAPHIC>, and <HIRAGANA>.

Factory class: solr.StandardTokenizerFactory

Arguments:

maxTokenLength: (integer, default 255) Solr ignores tokens that exceed the number of characters specified by maxTokenLength.

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
</analyzer>
```

In: "Please, email john.doe@foo.com by 03-09, re: m37-xq."

Out: "Please", "email", "john.doe@foo.com", "by", "03-09", "re", "m37-xq"

Classic Tokenizer

The Classic Tokenizer preserves the same behavior as the Standard Tokenizer of Solr versions 3.1 and previous. It does not use the [Unicode standard annex UAX#29](#) word boundary rules that the Standard Tokenizer uses. This tokenizer splits the text field into tokens, treating whitespace and punctuation as delimiters. Delimiter characters are discarded, with the following exceptions:

- Periods (dots) that are not followed by whitespace are kept as part of the token.
- Words are split at hyphens, unless there is a number in the word, in which case the token is not split and the numbers and hyphen(s) are preserved.
- Recognizes Internet domain names and email addresses and preserves them as a single token.

Factory class: solr.ClassicTokenizerFactory

Arguments:

`maxTokenLength`: (integer, default 255) Solr ignores tokens that exceed the number of characters specified by `maxTokenLength`.

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
</analyzer>
```

In: "Please, email john.doe@foo.com by 03-09, re: m37-xq."

Out: "Please", "email", "john.doe@foo.com", "by", "03-09", "re", "m37-xq"

Keyword Tokenizer

This tokenizer treats the entire text field as a single token.

Factory class: solr.KeywordTokenizerFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.KeywordTokenizerFactory" />
</analyzer>
```

In: "Please, email john.doe@foo.com by 03-09, re: m37-xq."

Out: "Please, email john.doe@foo.com by 03-09, re: m37-xq."

Letter Tokenizer

This tokenizer creates tokens from strings of contiguous letters, discarding all non-letter characters.

Factory class: solr.LetterTokenizerFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.LetterTokenizerFactory" />
</analyzer>
```

In: "I can't."

Out: "I", "can", "t"

Lower Case Tokenizer

Tokenizes the input stream by delimiting at non-letters and then converting all letters to lowercase. Whitespace and non-letters are discarded.

Factory class: solr.LowerCaseTokenizerFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.LowerCaseTokenizerFactory" />
</analyzer>
```

In: "I just **LOVE** my iPhone!"

Out: "i", "just", "love", "my", "iphone"

N-Gram Tokenizer

Reads the field text and generates n-gram tokens of sizes in the given range.

Factory class: solr.NGramTokenizerFactory

Arguments:

minGramSize: (integer, default 1) The minimum n-gram size, must be > 0.

maxGramSize: (integer, default 2) The maximum n-gram size, must be >= minGramSize.

Example:

Default behavior. Note that this tokenizer operates over the whole field. It does not break the field at whitespace. As a result, the space character is included in the encoding.

```
<analyzer>
  <tokenizer class="solr.NGramTokenizerFactory" />
</analyzer>
```

In: "hey man"

Out: "h", "e", "y", " ", "m", "a", "n", "he", "ey", "y ", " m", "ma", "an"

Example:

With an n-gram size range of 4 to 5:

```
<analyzer>
  <tokenizer class="solr.NGramTokenizerFactory" minGramSize="4" maxGramSize="5" />
</analyzer>
```

In: "bicycle"

Out: "bicy", "icyc", "cycl", "ycle", "bicyc", "icycl", "cycle"

Edge N-Gram Tokenizer

Reads the field text and generates edge n-gram tokens of sizes in the given range.

Factory class: solr.EdgeNGramTokenizerFactory

Arguments:

minGramSize: (integer, default 1) The minimum n-gram size, must be > 0.

maxGramSize: (integer, default 1) The maximum n-gram size, must be >= minGramSize.

side: ("front" or "back", default "front") Whether to compute the n-grams from the beginning (front) of the text or from the end (back).

Example:

Default behavior (min and max default to 1):

```
<analyzer>
  <tokenizer class="solr.EdgeNGramTokenizerFactory" />
</analyzer>
```

In: "babaloo"

Out: "b"

Example:

Edge n-gram range of 2 to 5

```
<analyzer>
  <tokenizer class="solr.EdgeNGramTokenizerFactory" minGramSize="2" maxGramSize="5" />
</analyzer>
```

In: "babaloo"

Out: "ba", "bab", "baba", "babal"

Example:

Edge n-gram range of 2 to 5, from the back side:

```
<analyzer>
  <tokenizer class="solr.EdgeNGramTokenizerFactory" minGramSize="2" maxGramSize="5"
side="back" />
</analyzer>
```

In: "babaloo"

Out: "oo", "loo", "aloo", "baloo"

ICU Tokenizer

This tokenizer processes multilingual text and tokenizes it appropriately based on its script attribute.

Factory class: solr.ICUTokenizerFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.ICUTokenizerFactory" />
</analyzer>
```

In: "Testing ""

Out: "Testing", "", "", "" "

Path Hierarchy Tokenizer

This tokenizer creates synonyms from file path hierarchies.

Factory class: solr.PathHierarchyTokenizerFactory

Arguments:

delimiter: (character, no default) You can specify the file path delimiter and replace it with a delimiter you provide. This can be useful for working with backslash delimiters.

replace: (character, no default) Specifies the delimiter character Solr uses in the tokenized output.

Example:

```
<fieldType name="text_path" class="solr.TextField" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.PathHierarchyTokenizerFactory" delimiter="/" replace="/" />
  </analyzer>
</fieldType>
```

In: "c:\usr\local\apache"

Out: "c:", "c:/usr", "c:/usr/local", "c:/usr/local/apache"

Regular Expression Pattern Tokenizer

This tokenizer uses a Java regular expression to break the input text stream into tokens. The expression provided by the pattern argument can be interpreted either as a delimiter that separates tokens, or to match patterns that should be extracted from the text as tokens.

See the Javadocs for [java.util.regex.Pattern](#) for more information on Java regular expression syntax.

Factory class: solr.PatternTokenizerFactory

Arguments:

pattern: (Required) The regular expression, as defined by in [java.util.regex.Pattern](#).

group: (Optional, default -1) Specifies which regex group to extract as the token(s).The value -1 means the regex should be treated as a delimiter that separates tokens.Non-negative group numbers (≥ 0) indicate that character sequences matching that regex group should be converted to tokens. Group zero refers to the entire regex, groups greater than zero refer to parenthesized sub-expressions of the regex, counted from left to right.

Example:

A comma separated list. Tokens are separated by a sequence of zero or more spaces, a comma, and zero or more spaces.

```
<analyzer>
  <tokenizer class="solr.PatternTokenizerFactory" pattern="\s*,\s*"/>
</analyzer>
```

In: "fee,fie, foe , fum, foo"

Out: "fee", "fie", "foe", "fum", "foo"

Example:

Extract simple, capitalized words. A sequence of at least one capital letter followed by zero or more letters of either case is extracted as a token.

```
<analyzer>
  <tokenizer class="solr.PatternTokenizerFactory" pattern="\[A-Z\]\[A-Za-z\]"
group="0"/>
</analyzer>
```

In: "Hello. My name is Inigo Montoya. You killed my father. Prepare to die."

Out: "Hello", "My", "Inigo", "Montoya". "You", "Prepare"

Example:

Extract part numbers which are preceded by "SKU", "Part" or "Part Number", case sensitive, with an optional semi-colon separator. Part numbers must be all numeric digits, with an optional hyphen. Regex capture groups are numbered by counting left parenthesis from left to right. Group 3 is the subexpression "[0-9-]+", which matches one or more digits or hyphens.

```
<analyzer>
  <tokenizer class="solr.PatternTokenizerFactory"
pattern="(SKU|Part(\sNumber)?):?\s(\[0-9-\]+)" group="3"/>
</analyzer>
```

In: "SKU: 1234, Part Number 5678, Part: 126-987"

Out: "1234", "5678", "126-987"

Type Tokenizer

This tokenizer filters tokens by its type, with either an exclude or include list.

Factory class: solr.TypeTokenFilterFactory

Arguments:

types: Defines the location of a file of types to filter.

enablePositionIncrements: If **true**, the token will be incremented by position.

useWhiteList: If **true**, the file defined in types should be used as include list.

Example:

```
<analyzer>
  <filter class="solr.TypeTokenFilterFactory" types="stotypes.txt"
    enablePositionIncrements="true" useWhiteList="false"/>
</analyzer>
```

UAX29 URL Email Tokenizer

This tokenizer splits the text field into tokens, treating whitespace and punctuation as delimiters. Delimiter characters are discarded, with the following exceptions:

- Periods (dots) that are not followed by whitespace are kept as part of the token.
- Words are split at hyphens, unless there is a number in the word, in which case the token is not split and the numbers and hyphen(s) are preserved.
- Recognizes top-level (.com) Internet domain names; email addresses; file:///, http(s)://, and ftp:// addresses; IPv4 and IPv6 addresses; and preserves them as a single token.

The UAX29 URL Email Tokenizer supports [Unicode standard annex UAX#29](#) word boundaries with the following token types: <ALPHANUM>, <NUM>, URL, EMAIL, <SOUTHEAST_ASIAN>, <IDEOGRAPHIC>, and <HIRAGANA>.

Factory class: solr.UAX29URLEmailTokenizerFactory

Arguments:

maxTokenLength: (integer, default 255) Solr ignores tokens that exceed the number of characters specified by maxTokenLength.

Example:

```
<analyzer>
  <tokenizer class="solr.UAX29URLEmailTokenizerFactory" />
</analyzer>
```

In: "Visit <http://accarol.com/contact.htm?from=external&a=10> or e-mail bob.cratchet@accarol.com"

Out: "Visit", "http://accarol.com/contact.htm?from=external&a=10", "or", "email", "bob.cratchet@accarol.com"

White Space Tokenizer

Simple tokenizer that splits the text stream on whitespace and returns sequences of non-whitespace characters as tokens. Note that any punctuation *will* be included in the tokenization.

Factory class: solrWhitespaceTokenizerFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory" />
</analyzer>
```

In: "To be, or what?"

Out: "To", "be,", "or", "what?"

Related Topics

- [TokenizerFactories](#)



Filter Descriptions

You configure each filter with a `<filter>` element in `schema.xml` as a child of `<analyzer>`, following the `<tokenizer>` element. Filter definitions should follow a tokenizer or another filter definition because they take a `TokenStream` as input. For example.

```
<fieldType name="text" class="solr.TextField">
    <analyzer type="index">
        <tokenizer
            class="solr.StandardTokenizerFactory"/>
        <filter
            class="solr.LowerCaseFilterFactory"/>...
    </analyzer>
</fieldType>
```

Filters discussed in this section:

The `class` attribute names a factory class that will instantiate a filter object as needed. Filter factory classes must implement the `org.apache.solr.analysis.TokenFilterFactory` interface. Like tokenizers, filters are also instances of `TokenStream` and thus are producers of tokens. Unlike tokenizers, filters also consume tokens from a `TokenStream`. This allows you to mix and match filters, in any order you prefer, downstream of a tokenizer.

Arguments may be passed to tokenizer factories to modify their behavior by setting attributes on the `<filter>` element. For example:

```
<fieldType name="semicolonDelimited"
    class="solr.TextField">
    <analyzer type="query">
        <tokenizer class="solr.PatternTokenizerFactory"
            pattern=";" />
        <filter class="solr.LengthFilterFactory" *min="2"
            max="7"/>
    </analyzer>
</fieldType>
```

The following sections describe the filter factories that are included in this release of Solr.

For more information about Solr's filters, see
<http://wiki.apache.org/solr/AnalyzersTokenizersTokenFilters>.

- [ASCII Folding Filter](#)
- [Beider-Morse Filter](#)
- [Classic Filter](#)
- [Common Grams Filter](#)
- [Collation Key Filter](#)
- [Edge N-Gram Filter](#)
- [English Minimal Stem Filter](#)
- [Hunspell Stem Filter](#)
- [Hyphenated Words Filter](#)
- [ICU Folding Filter](#)
- [ICU Normalizer 2 Filter](#)
- [ICU Transform Filter](#)
- [Keep Words Filter](#)
- [KStem Filter](#)
- [Length Filter](#)
- [Lower Case Filter](#)
- [N-Gram Filter](#)
- [Numeric Payload Token Filter](#)
- [Pattern Replace Filter](#)
- [Phonetic Filter](#)
- [Porter Stem Filter](#)
- [Position Filter Factory](#)
- [Remove Duplicates Token Filter](#)
- [Reversed Wildcard Filter](#)
- [Shingle Filter](#)
- [Snowball Porter Stemmer Filter](#)
- [Standard Filter](#)
- [Stop Filter](#)
- [Synonym Filter](#)
- [Token Offset Payload Filter](#)
- [Trim Filter](#)
- [Type As Payload Filter](#)
- [Word Delimiter Filter](#)
- [Related Topics](#)

ASCII Folding Filter

This filter converts alphabetic, numeric, and symbolic Unicode characters which are not in the Basic Latin Unicode block (the first 127 ASCII characters) to their ASCII equivalents, if one exists. This filter converts characters from the following Unicode blocks:

- [C1 Controls and Latin-1 Supplement](#) (PDF)
- [Latin Extended-A](#) (PDF)
- [Latin Extended-B](#) (PDF)
- [Latin Extended Additional](#) (PDF)
- [Latin Extended-C](#) (PDF)
- [Latin Extended-D](#) (PDF)
- [IPA Extensions](#) (PDF)
- [Phonetic Extensions](#) (PDF)
- [Phonetic Extensions Supplement](#) (PDF)
- [General Punctuation](#) (PDF)
- [Superscripts and Subscripts](#) (PDF)
- [Enclosed Alphanumerics](#) (PDF)
- [Dingbats](#) (PDF)
- [Supplemental Punctuation](#) (PDF)
- [Alphabetic Presentation Forms](#) (PDF)
- [Halfwidth and Fullwidth Forms](#) (PDF)

Factory class: solr.ASCIIFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.ASCIIFilterFactory"/>
</analyzer>
```

In: "á" (Unicode character 00E1)

Out: "á" (ASCII character 160)

Beider-Morse Filter

Implements the Beider-Morse Phonetic Matching (BMPM) algorithm, which allows identification of similar names, even if they are spelled differently or in different languages. More information about how this works is available in the section on [Phonetic Matching](#).

Factory class: solr.BeiderMorseFilterFactory

Arguments:

`nameType`: Types of names. Valid values are GENERIC, ASHKENAZI, or SEPHARDIC. If not processing Ashkenazi or Sephardic names, use GENERIC.

`ruleType`: Types of rules to apply. Valid values are APPROX or EXACT.

`concat`: Defines if multiple possible matches should be combined with a pipe ("|").

`languageSet`: The language set to use. The value "auto" will allow the Filter to identify the language, or a comma-separated list can be supplied.

Example:

```
<analyzer>
  <filter class="solr.BeiderMorseFilterFactory" nameType="GENERIC" ruleType="APPROX"
    concat="true" languageSet="auto"
  </filter>
</analyzer>
```

Classic Filter

This filter takes the output of the [Classic Tokenizer](#) and strips periods from acronyms and "'s" from possessives.

Factory class: solr.ClassicFilterFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.ClassicTokenizerFactory"/>
  <filter class="solr.ClassicFilterFactory"/>
</analyzer>
```

In: "I.B.M. cat's can't"

Tokenizer to Filter: "I.B.M", "cat's", "can't"

Out: "IBM", "cat", "can't"

Common Grams Filter

This filter creates word shingles by combining common tokens such as stop words with regular tokens. This is useful for creating phrase queries containing common words, such as "the cat." Solr normally ignores stop words in queried phrases, so searching for "the cat" would return all matches for the word "cat."

Factory class: solr.CommonGramsFilterFactory

Arguments:

words: (a common word file in .txt format) Provide the name of a common word file, such as stopwords.txt.

format: (optional) If the stopwords list has been formatted for Snowball, you can specify format="snowball" so Solr can read the stopwords file.

ignoreCase: (boolean) If true, the filter ignores the case of words when comparing them to the common word file.

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.CommonGramsFilterFactory" words="stopwords.txt"
ignoreCase="true" />
</analyzer>
```

In: "the Cat"

Tokenizer to Filter: "the", "Cat"

Out: "the_cat"

Collation Key Filter

Collation allows sorting of text in a language-sensitive way. It is usually used for sorting, but can also be used with advanced searches. We've covered this in much more detail in the section on [Unicode Collation](#).

Edge N-Gram Filter

This filter generates edge n-gram tokens of sizes within the given range.

Factory class: solr.EdgeNGramFilterFactory**Arguments:**

minGramSize: (integer, default 1) The minimum gram size.

maxGramSize: (integer, default 1) The maximum gram size.

Example:

Default behavior.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.EdgeNGramFilterFactory" />
</analyzer>
```

In: "four score and twenty"

Tokenizer to Filter: "four", "score", "and", "twenty"

Out: "f", "s", "a", "t"

Example:

A range of 1 to 4.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.EdgeNGramFilterFactory" minGramSize="1" maxGramSize="4" />
</analyzer>
```

In: "four score"

Tokenizer to Filter: "four", "score"

Out: "f", "fo", "fou", "four", "s", "sc", "sco", "scor"

Example:

A range of 4 to 6.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.EdgeNGramFilterFactory" minGramSize="4" maxGramSize="6" />
</analyzer>
```

In: "four score and twenty"

Tokenizer to Filter: "four", "score", "and", "twenty"

Out: "four", "sco", "scor"

English Minimal Stem Filter

This filter stems plural English words to their singular form.

Factory class: solr.EnglishMinimalStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.EnglishMinimalStemFilterFactory" />
</analyzer>
```

In: "dogs cats"

Tokenizer to Filter: "dogs", "cats"

Out: "dog", "cat"

Hunspell Stem Filter

The [Hunspell Stem Filter](#) provides support for several languages. You must provide the dictionary (.dic) and rules (.aff) files for each language you wish to use with the Hunspell Stem Filter. You can download those language files [here](#). Be aware that your results will vary widely based on the quality of the provided dictionary and rules files. For example, some languages have only a minimal word list with no morphological information. On the other hand, for languages that have no stemmer but do have an extensive dictionary file, the Hunspell stemmer may be a good choice.

Factory class: solr.HunspellStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrHunspellStemFilterFactory"
    dictionary="en_GB.dic"
    affix="en_GB.aff"
    ignoreCase="true" />
</analyzer>
```

In: "jump jumping jumped"

Tokenizer to Filter: "jump", "jumping", "jumped"

Out: "jump", "jump", "jump"

Hyphenated Words Filter

This filter reconstructs hyphenated words that have been tokenized as two tokens because of a line break or other intervening whitespace in the field test. If a token ends with a hyphen, it is joined with the following token and the hyphen is discarded. Note that for this filter to work properly, the upstream tokenizer must not remove trailing hyphen characters. This filter is generally only useful at index time.

Factory class: solr.HyphenatedWordsFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrHyphenatedWordsFilterFactory"/>
</analyzer>
```

In: "A hyphen- ated word"

Tokenizer to Filter: "A", "hyphen-", "ated", "word"

Out: "A", "hyphenated", "word"

ICU Folding Filter

This filter is a custom Unicode normalization form that applies the foldings specified in [Unicode Technical Report 30](#) in addition to the `NFKC_Casefold` normalization form as described in [ICU Normalizer 2 Filter](#). This filter is a better substitute for the combined behavior of the [ASCII Folding Filter](#), [Lower Case Filter](#), and [ICU Normalizer 2 Filter](#).

To use this filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

Factory class: `solr.ICUFoldingFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.ICUFoldingFilterFactory"/>
</analyzer>
```

For detailed information on this normalization form, see
<http://www.unicode.org/reports/tr30/tr30-4.html>.

ICU Normalizer 2 Filter

This filter factory normalizes text according to one of five Unicode Normalization Forms as described in [Unicode Standard Annex #15](#):

- NFC: (`name="nfc" mode="compose"`) Normalization Form C, canonical decomposition
- NFD: (`name="nfc" mode="decompose"`) Normalization Form D, canonical decomposition, followed by canonical composition
- NFKC: (`name="nfkc" mode="compose"`) Normalization Form KC, compatibility decomposition
- NFKD: (`name="nfkc" mode="decompose"`) Normalization Form KD, compatibility decomposition, followed by canonical composition
- NFKC_Casefold: (`name="nfkc_cf" mode="compose"`) Normalization Form KC, with additional Unicode case folding. Using the ICU Normalizer 2 Filter is a better-performing substitution for the [Lower Case Filter](#) and NFKC normalization.

Factory class: `solr.ICUNormalizer2FilterFactory`

Arguments:

`name`: (string) The name of the normalization form; `nfc`, `nfd`, `nfkc`, `nfkd`, `nfkc_cf`

`mode`: (string) The mode of Unicode character composition and decomposition; `compose` or `decompose`

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.ICUNormalizer2FilterFactory" />
</analyzer>
```

For detailed information about these Unicode Normalization Forms, see
<http://unicode.org/reports/tr15/>.

To use this filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

ICU Transform Filter

This filter applies [ICU Tranforms](#) to text. This filter supports only ICU System Transforms. Custom rule sets are not supported.

Factory class: `solr.ICUTransformFilterFactory`

Arguments:

`id`: (string) The identifier for the ICU System Transform you wish to apply with this filter. For a full list of ICU System Transforms, see

http://demo.icu-project.org/icu-bin/translit?TEMPLATE_FILE=data/translit_rule_main.html.

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.ICUTransformFilterFactory" id="Traditional-Simplified" />
</analyzer>
```

For detailed information about ICU Transforms, see
<http://userguide.icu-project.org/transforms/general>.

To use this filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

Keep Words Filter

This filter discards all tokens except those that are listed in the given word list. This is the inverse of the Stop Words Filter. This filter can be useful for building specialized indices for a constrained set of terms.

Factory class: `solr.KeepWordFilterFactory`

Arguments:

words: (required) Path of a text file containing the list of keep words, one per line. Blank lines and lines that begin with "#" are ignored. This may be an absolute path, or a simple filename in the Solr config directory.

ignoreCase: (true/false) If **true** then comparisons are done case-insensitively. If this argument is true, then the words file is assumed to contain only lowercase words. The default is **false**.

Example:

Where `keepwords.txt` contains:

happy

funny

silly

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.KeepWordFilterFactory" words="keepwords.txt" />
</analyzer>
```

In: "Happy, sad or funny"

Tokenizer to Filter: "Happy", "sad", "or", "funny"

Out: "funny"

Example:

Same `keepwords.txt`, case insensitive:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.KeepWordFilterFactory" words="keepwords.txt" ignoreCase="true" />
</analyzer>
```

In: "Happy, sad or funny"

Tokenizer to Filter: "Happy", "sad", "or", "funny"

Out: "Happy", "funny"

Example:

Using LowerCaseFilterFactory before filtering for keep words, no ignoreCase flag.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.LowerCaseFilterFactory" />
  <filter class="solr.KeepWordFilterFactory" words="keepwords.txt" />
</analyzer>
```

In: "Happy, sad or funny"

Tokenizer to Filter: "Happy", "sad", "or", "funny"

Filter to Filter: "happy", "sad", "or", "funny"

Out: "happy", "funny"

KStem Filter

KStem is an alternative to the Porter Stem Filter for developers looking for a less aggressive stemmer. KStem was written by Bob Krovetz, ported to Lucene by Sergio Guzman-Lara (UMASS Amherst). This stemmer is only appropriate for English language text.

Factory class: solr.KStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.KStemFilterFactory" />
</analyzer>
```

In: "jump jumping jumped"

Tokenizer to Filter: "jump", "jumping", "jumped"

Out: "jump", "jump", "jump"

Length Filter

This filter passes tokens whose length falls within the min/max limit specified. All other tokens are discarded.

Factory class: solr.LengthFilterFactory

Arguments:

min: (integer, required) Minimum token length. Tokens shorter than this are discarded.

max: (integer, required, must be \geq min) Maximum token length. Tokens longer than this are discarded.

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.LengthFilterFactory" min="3" max="7" />
</analyzer>
```

In: "turn right at Albuquerque"

Tokenizer to Filter: "turn", "right", "at", "Albuquerque"

Out: "turn", "right"

Lower Case Filter

Converts any uppercase letters in a token to the equivalent lowercase token. All other characters are left unchanged.

Factory class: solrLowerCaseFilterFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
</analyzer>
```

In: "Down With CamelCase"

Tokenizer to Filter: "Down", "With", "CamelCase"

Out: "down", "with", "camelcase"

N-Gram Filter

Generates n-gram tokens of sizes in the given range.

Factory class: solr.NGramFilterFactory

Arguments:

`minGramSize`: (integer, default 1) The minimum gram size.

`maxGramSize`: (integer, default 2) The maximum gram size.

Example:

Default behavior.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.NGramFilterFactory"/>
</analyzer>
```

In: "four score"

Tokenizer to Filter: "four", "score"

Out: "f", "o", "u", "r", "fo", "ou", "ur", "s", "c", "o", "r", "e", "sc", "co", "or", "re"

Example:

A range of 1 to 4.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.NGramFilterFactory" *minGramSize="1" maxGramSize="4"/>
</analyzer>
```

In: "four score"

Tokenizer to Filter: "four", "score"

Out: "f", "fo", "fou", "four", "s", "sc", "sco", "scor"

Example:

A range of 3 to 5.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.NGramFilterFactory" *minGramSize="3" maxGramSize="5"/>
</analyzer>
```

In: "four score"

Tokenizer to Filter: "four", "score"

Out: "fou", "our", "four", "sco", "cor", "ore", "scor", "core", "score"

Numeric Payload Token Filter

This filter adds a numeric floating point payload value to tokens that match a given type. Refer to the Javadoc for the `org.apache.lucene.analysis.Token` class for more information about token types and payloads.

Factory class: `solr.NumericPayloadTokenFilterFactory`

Arguments:

`payload`: (required) A floating point value that will be added to all matching tokens.

`typeMatch`: (required) A token type name string. Tokens with a matching type name will have their payload set to the above floating point value.

Example:

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solr.NumericPayloadTokenFilterFactory" payload="0.75"
    typeMatch="word"/>
</analyzer>
```

In: "bing bang boom"

Tokenizer to Filter: "bing", "bang", "boom"

Out: "bing"[0.75], "bang"[0.75], "boom"[0.75]

Pattern Replace Filter

This filter applies a regular expression to each token and, for those that match, substitutes the given replacement string in place of the matched pattern. Tokens which do not match are passed through unchanged.

Factory class: `solr.PatternReplaceFilter`

Arguments:

`pattern`: (required) The regular expression to test against each token, as per `java.util.regex.Pattern`.

replacement: (required) A string to substitute in place of the matched pattern. This string may contain references to capture groups in the regex pattern. See the Javadoc for `java.util.regex.Matcher`.

replace: ("all" or "first", default "all") Indicates whether all occurrences of the pattern in the token should be replaced, or only the first.

Example:

Simple string replace:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PatternReplaceFilter" pattern="cat" replacement="dog" />
</analyzer>
```

In: "cat concatenate catycat"

Tokenizer to Filter: "cat", "concatenate", "catycat"

Out: "dog", "condogenate", "dogydog"

Example:

String replacement, first occurrence only:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PatternReplaceFilter" pattern="cat" replacement="dog"
*replace="first" />
</analyzer>
```

In: "cat concatenate catycat"

Tokenizer to Filter: "cat", "concatenate", "catycat"

Out: "dog", "condogenate", "dogycat"

Example:

More complex pattern with capture group reference in the replacement. Tokens that start with non-numeric characters and end with digits will have an underscore inserted before the numbers. Otherwise the token is passed through.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PatternReplaceFilter" pattern="(\D+)(\d+)$" replacement="$1_$2"/>
</analyzer>
```

In: "cat foo1234 9987 blah1234foo"

Tokenizer to Filter: "cat", "foo1234", "9987", "blah1234foo"

Out: "cat", "foo_1234", "9987", "blah1234foo"

Phonetic Filter

This filter creates tokens using one of the phonetic encoding algorithms in the org.apache.commons.codec.language package.

Factory class: solr.PhoneticFilterFactory

Arguments:

`{{encoder}}`: (required) The name of the encoder to use. The encoder name must be one of the following (case insensitive): "[DoubleMetaphone](#)", "[Metaphone](#)", "[Soundex](#)", "[RefinedSoundex](#)", "[Caverphone](#)", or "[ColognePhonetic](#)"

`inject`: (true/false) If true (the default), then new phonetic tokens are added to the stream. Otherwise, tokens are replaced with the phonetic equivalent. Setting this to false will enable phonetic matching, but the exact spelling of the target word may not match.

`maxCodeLength`: (integer) The maximum length of the code to be generated by the Metaphone or Double Metaphone encoders.

Example:

Default behavior for DoubleMetaphone encoding.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PhoneticFilterFactory" encoder="DoubleMetaphone" />
</analyzer>
```

In: "four score and twenty"

Tokenizer to Filter: "four"(1), "score"(2), "and"(3), "twenty"(4)

Out: "four"(1), "FR"(1), "score"(2), "SKR"(2), "and"(3), "ANT"(3), "twenty"(4), "TNT"(4)

The phonetic tokens have a position increment of 0, which indicates that they are at the same position as the token they were derived from (immediately preceding).

Example:

Discard original token.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PhoneticFilterFactory" encoder="DoubleMetaphone" inject="false" />
</analyzer>
```

In: "four score and twenty"

Tokenizer to Filter: "four"(1), "score"(2), "and"(3), "twenty"(4)

Out: "FR"(1), "SKR"(2), "ANT"(3), "TWNT"(4)

Example:

Default Soundex encoder.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PhoneticFilterFactory" encoder="Soundex" />
</analyzer>
```

In: "four score and twenty"

Tokenizer to Filter: "four"(1), "score"(2), "and"(3), "twenty"(4)

Out: "four"(1), "F600"(1), "score"(2), "S600"(2), "and"(3), "A530"(3), "twenty"(4), "T530"(4)

Porter Stem Filter

This filter applies the Porter Stemming Algorithm for English. The results are similar to using the Snowball Porter Stemmer with the `language="English"` argument. But this stemmer is coded directly in Java and is not based on Snowball. Nor does it accept a list of protected words. This stemmer is only appropriate for English language text.

Factory class: solr.PorterStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.PorterStemFilterFactory" />
</analyzer>
```

In: "jump jumping jumped"

Tokenizer to Filter: "jump", "jumping", "jumped"

Out: "jump", "jump", "jump"

Position Filter Factory

This filter sets the position increment values of all tokens in a token stream except the first, which retains its original position increment value.

Factory class: solr.PositionIncrementFilterFactory

Arguments:

positionIncrement: (integer, default = 0) The position increment value to apply to all tokens in a token stream except the first.

Example:

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory" />
  <filter class="solrPositionFilterFactory" positionIncrement="1" />
</analyzer>
```

In: "hello world"

Tokenizer to Filter: "hello", "world"

Out: "hello" (token position 1), "world" (token position 3)

Remove Duplicates Token Filter

The filter removes duplicate tokens in the stream. Tokens are considered to be duplicates if they have the same text and position values.

Factory class: solr.RemoveDuplicatesTokenFilterFactory

Arguments: None

Example:

This is an artificial example that uses the [Synonym Filter](#) to generate duplicate symbols, which are then removed. The file `testsyns.txt` contains the following:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SynonymFilterFactory" synonyms="testsyns.txt" />
  <filter class="solr.RemoveDuplicatesTokenFilterFactory" />
</analyzer>
```

In: "blurt blort"

Tokenizer to Filter: "blurt"(1), "blurt"(2)

Tokenizer to Filter: "foo"(1), "foo"(1), "bar"(2), "bar"(2)

Out: "foo"(1), "bar"(2)

Reversed Wildcard Filter

This filter reverses tokens to provide faster leading wildcard and prefix queries. Tokens without wildcards are not reversed.

Factory class: solr.ReveresedWildcardFilterFactory

Arguments:

`withOriginal` (boolean) If true, the filter produces both original and reversed tokens at the same positions. If false, produces only reversed tokens.

`maxPosAsterisk` (integer, default = 2) The maximum position of the asterisk wildcard ('*') that triggers the reversal of the query term. Terms with asterisks at positions above this value are not reversed.

`maxPosQuestion` (integer, default = 1) The maximum position of the question mark wildcard ('?') that triggers the reversal of query term. To reverse only pure suffix queries (queries with a single leading asterisk), set this to 0 and `maxPosAsterisk` to 1.

`maxFractionAsterisk` (float, default = 0.0) An additional parameter that triggers the reversal if asterisk ('*') position is less than this fraction of the query token length.

`minTrailing` (integer, default = 2) The minimum number of trailing characters in a query token after the last wildcard character. For good performance this should be set to a value larger than 1.

Example:

```
<analyzer type="index">
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrReversedWildcardFilterFactory" withOriginal="true"
    maxPosAsterisk="2" maxPosQuestion="1" minTrailing="2" maxFractionAsterisk="0"/>
</analyzer>
```

In: "*foo *bar"

Tokenizer to Filter: "*foo", "*bar"

Out: "oof*", "rab*"

Shingle Filter

This filter constructs shingles, which are token n-grams, from the token stream. It combines runs of tokens into a single token.

Factory class: solr.ShingleFilterFactory

Arguments:

`maxShingleSize`: (integer, must be ≥ 2 , default 2) The maximum number of tokens per shingle.

`outputUnigrams`: (true/false) If true (the default), then each individual token is also included at its original position.

Example:

Default behavior.

```
<analyzer>
  <tokenizer class="solrStandardTokenizerFactory"/>
  <filter class="solrShingleFilterFactory"/>
</analyzer>
```

In: "To be, or what?"

Tokenizer to Filter: "To"(1), "be"(2), "or"(3), "what"(4)

Out: "To"(1), "To be"(1), "be"(2), "be or"(2), "or"(3), "or what"(3), "what"(4)

Example:

A shingle size of four, do not include original token.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.ShingleFilterFactory" maxShingleSize="4" outputUnigrams="false"/>
</analyzer>
```

In: "To be, or not to be."

Tokenizer to Filter: "To"(1), "be"(2), "or"(3), "not"(4), "to"(5), "be"(6)

Out: "To be"(1), "To be or"(1), "To be or not"(1), "be or"(2), "be or not"(2), "be or not to"(2), "or not"(3), "or not to"(3), "or not to be"(3), "not to"(4), "not to be"(4), "to be"(5)

Snowball Porter Stemmer Filter

This filter factory instantiates a language-specific stemmer generated by Snowball. Snowball is a software package that generates pattern-based word stemmers. This type of stemmer is not as accurate as a table-based stemmer, but is faster and less complex. Table-driven stemmers are labor intensive to create and maintain and so are typically commercial products.

Solr contains Snowball stemmers for Armenian, Basque, Catalan, Danish, Dutch, English, Finnish, French, German, Hungarian, Italian, Norwegian, Portuguese, Romanian, Russian, Spanish, Swedish and Turkish. For more information on Snowball, visit <http://snowball.tartarus.org/>.

`StopFilterFactory`, `CommonGramsFilterFactory`, and `CommonGramsQueryFilterFactory` can optionally read stopwords in Snowball format (specify `format="snowball"` in the configuration of those FilterFactories).

Factory class: `solr.SnowballPorterFilterFactory`

Arguments:

`language`: (default "English") The name of a language, used to select the appropriate Porter stemmer to use. Case is significant. This string is used to select a package name in the "org.tartarus.snowball.ext" class hierarchy.

`protected`: Path of a text file containing a list of protected words, one per line. Protected words will not be stemmed. Blank lines and lines that begin with "#" are ignored. This may be an absolute path, or a simple file name in the Solr config directory.

Example:

Default behavior:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SnowballPorterFilterFactory" />
</analyzer>
```

In: "flip flipped flipping"

Tokenizer to Filter: "flip", "flipped", "flipping"

Out: "flip", "flip", "flip"

Example:

French stemmer, English words:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SnowballPorterFilterFactory" language="French" />
</analyzer>
```

In: "flip flipped flipping"

Tokenizer to Filter: "flip", "flipped", "flipping"

Out: "flip", "flipped", "flipping"

Example:

Spanish stemmer, Spanish words:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SnowballPorterFilterFactory" language="Spanish" />
</analyzer>
```

In: "cante canta"

Tokenizer to Filter: "cante", "canta"

Out: "cant", "cant"

Standard Filter

This filter removes dots from acronyms and the substring "'s" from the end of tokens. This filter depends on the tokens being tagged with the appropriate term-type to recognize acronyms and words with apostrophes.

Factory class: solr.StandardFilterFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.StandardFilterFactory" />
</analyzer>
```

In: "Bob's I.O.U."

Tokenizer to Filter: "Bob's", "I.O.U."

Out: "Bob". "IOU"

Stop Filter

This filter discards, or *stops* analysis of, tokens that are on the given stop words list. A standard stop words list is included in the Solr config directory, named stopwords.txt, which is appropriate for typical English language text.

Factory class: solr.StopFilterFactory

Arguments:

words: (optional) The path to a file that contains a list of stop words, one per line. Blank lines and lines that begin with "#" are ignored. This may be an absolute path, or path relative to the Solr config directory.

format: (optional) If the stopwords list has been formatted for Snowball, you can specify `format="snowball"` so Solr can read the stopwords file.

ignoreCase: (true/false, default false) Ignore case when testing for stop words. If true, the stop list should contain lowercase words.

enablePositionIncrements: (true/false, default false) When true, if a token is stopped (discarded) then the position of the following token is incremented.

Example:

Case-sensitive matching, capitalized words not stopped. Token positions skip stopped words.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.StopFilterFactory" words="stopwords.txt"/>
</analyzer>
```

In: "To be or what?"

Tokenizer to Filter: "To"(1), "be"(2), "or"(3), "what"(4)

Out: "To"(1), "what"(2)

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.StopFilterFactory" words="stopwords.txt" ignoreCase="true"/>
</analyzer>
```

In: "To be or what?"

Tokenizer to Filter: "To"(1), "be"(2), "or"(3), "what"(4)

Out: "what"(1)

Example:

Position increment enabled, original positions retained. No tokens at positions of stopped words.

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.StopFilterFactory" words="stopwords.txt" ignoreCase="true"
enablePositionIncrements="true"/>
</analyzer>
```

In: "You are a star"

Tokenizer to Filter: "You"(1), "are"(2), "a"(3), "star"(4)

Out: "You"(1), "star"(4)

Synonym Filter

This filter does synonym mapping. Each token is looked up in the list of synonyms and if a match is found, then the synonym is emitted in place of the token. The position value of the new tokens are set such they all occur at the same position as the original token.

Factory class: solr.SynonymFilterFactory**Arguments:**

`synonyms`: (required) The path of a file that contains a list of synonyms, one per line. Blank lines and lines that begin with "#" are ignored. This may be an absolute path, or path relative to the Solr config directory. There are two ways to specify synonym mappings:

- A comma-separated list of words. If the token matches any of the words, then all the words in the list are substituted, which will include the original token.
- Two comma-separated lists of words with the symbol ">" between them. If the token matches any word on the left, then the list on the right is substituted. The original token will not be included unless it is also in the list on the right.

For the following examples, assume the following `synonyms.txt` file:

```
couch,sofa,divan
teh => the
huge,ginormous,humungous => large
small => tiny,teeny,weeny
```

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SynonymFilterFactory" synonyms="mysynonyms.txt" />
</analyzer>
```

In: "teh small couch"

Tokenizer to Filter: "teh"(1), "small"(2), "couch"(3)

Out: "the"(1), "tiny"(2), "teeny"(2), "weeny"(2), "couch"(3), "sofa"(3), "divan"(3)

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.SynonymFilterFactory" synonyms="mysynonyms.txt" />
</analyzer>
```

In: "teh ginormous, humungous sofa"

Tokenizer to Filter: "teh"(1), "ginormous"(2), "humungous"(3), "sofa"(4)

Out: "the"(1), "large"(2), "large"(3), "couch"(4), "sofa"(4), "divan"(4)

Token Offset Payload Filter

This filter adds the numeric character offsets of the token as a payload value for that token.

Factory class: solr.TokenOffsetPayloadTokenFilterFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory" />
  <filter class="solr.TokenOffsetPayloadTokenFilterFactory" />
</analyzer>
```

In: "bing bang boom"

Tokenizer to Filter: "bing", "bang", "boom"

Out: "bing"[0,4], "bang"[5,9], "boom"[10,14]

Trim Filter

This filter trims leading and/or trailing whitespace from tokens. Most tokenizers break tokens at whitespace, so this filter is most often used for special situations.

Factory class: solr.TrimFilterFactory

Arguments:

updateOffsets: (true/false, default false) If true, the token's start/end offsets are adjusted to account for any whitespace that was removed.

Example:

The PatternTokenizerFactory configuration used here splits the input on simple commas, it does not remove whitespace.

```
<analyzer>
  <tokenizer class="solr.PatternTokenizerFactory" pattern=" , " />
  <filter class="solr.TrimFilterFactory" />
</analyzer>
```

In: "one, two , three ,four "

Tokenizer to Filter: "one", "two", "three", "four"

Out: "one", "two", "three", "four"

Type As Payload Filter

This filter adds the token's type, as an encoded byte sequence, as its payload.

Factory class: solr.TypeAsPayloadTokenFilterFactory

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solr.TypeAsPayloadTokenFilterFactory"/>
</analyzer>
```

In: "Pay Bob's I.O.U."

Tokenizer to Filter: "Pay", "Bob's", "I.O.U."

Out: "Pay"[<ALPHANUM>], "Bob's"[<APOSTROPHE>], "I.O.U." [<ACRONYM>]

Word Delimiter Filter

This filter splits tokens at word delimiters. The rules for determining delimiters are determined as follows:

- A change in case within a word: "CamelCase" -> "Camel", "Case" This can be disabled by setting splitOnCaseChange="0" (see below).
- A transition from alpha to numeric characters or vice versa: "Gonzo5000" ➤ "Gonzo", "5000""4500XL" > "4500", "XL" This can be disabled by setting splitOnNumerics ="0".
- Non-alphanumeric characters (discarded): "hot-spot" -> "hot", "spot"
- A trailing "'s" is removed: "O'Reilly's" -> "O", "Reilly"
- Any leading or trailing delimiters are discarded: "-hot-spot" > "hot", "spot"

Factory class: solr.WordDelimiterFilterFactory

Arguments:

generateWordParts: (integer, default 1) If non-zero, splits words at delimiters. For example:"CamelCase", "hot-spot" -> "Camel", "Case", "hot", "spot"

generateNumberParts: (integer, default 1) If non-zero, splits numeric strings at delimiters:"1947-32" ->"1947", "32"

splitOnCaseChange: (integer, default 1) If 0, words are not split on camel-case changes:"BugBlaster-XL" -> "BugBlaster", "XL"Example 1 below illustrates the default (non-zero) splitting behavior.

splitOnNumerics: (integer, default 1) If 0, don't split words on transitions from alpha to numeric:"FemBot3000" -> "Fem", "Bot3000"

catenateWords: (integer, default 0) If non-zero, maximal runs of word parts will be joined: "hot-spot-sensor's" -> "hotspotsensor"

catenateNumbers: (integer, default 0) If non-zero, maximal runs of number parts will be joined: 1947-32" -> "194732"

catenateAll: (0/1, default 0) If non-zero, runs of word and number parts will be joined: "Zap-Master-9000" -> "ZapMaster9000"

preserveOriginal: (integer, default 0) If non-zero, the original token is preserved: "Zap-Master-9000" -> "Zap-Master-9000", "Zap", "Master", "9000"

protected: (optional) The pathname of a file that contains a list of protected words that should be passed through without splitting.

stemEnglishPossessive: (integer, default 1) If 1, strips the possessive "'s" from each subword.

Example:

Default behavior. The whitespace tokenizer is used here to preserve non-alphanumeric characters.

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory" />
  <filter class="solrWordDelimiterFilterFactory" />
</analyzer>
```

In: "hot-spot RoboBlaster/9000 100XL"

Tokenizer to Filter: "hot-spot", "RoboBlaster/9000", "100XL"

Out: "hot", "spot", "Robo", "Blaster", "9000", "100", "XL"

Example:

Do not split on case changes, and do not generate number parts. Note that by not generating number parts, tokens containing only numeric parts are ultimately discarded.

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrWordDelimiterFilterFactory" generateNumberParts="0"
splitOnCaseChange="0"/>
</analyzer>
```

In: "hot-spot RoboBlaster/9000 100-42"

Tokenizer to Filter: "hot-spot", "RoboBlaster/9000", "100-42"

Out: "hot", "spot", "RoboBlaster", "9000"

Example:

Concatenate word parts and number parts, but not word and number parts that occur in the same token.

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrWordDelimiterFilterFactory" catenateWords="1"
catenateNumbers="1"/>
</analyzer>
```

In: "hot-spot 100+42 XL40"

Tokenizer to Filter: "hot-spot"(1), "100+42"(2), "XL40"(3)

Out: "hot"(1), "spot"(2), "hotspot"(2), "100"(3), "42"(4), "10042"(4), "XL"(5), "40"(6)

Example:

Concatenate all. Word and/or number parts are joined together.

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrWordDelimiterFilterFactory" catenateAll="1"/>
</analyzer>
```

In: "XL-4000/ES"

Tokenizer to Filter: "XL-4000/ES"(1)

Out: "XL"(1), "4000"(2), "ES"(3), "XL4000ES"(3)

Example:

Using a protected words list that contains "AstroBlaster" and "XL-5000" (among others).

```
<analyzer>
  <tokenizer class="solrWhitespaceTokenizerFactory"/>
  <filter class="solrWordDelimiterFilterFactory" protected="protwords.txt"/>
</analyzer>
```

In: "FooBar AstroBlaster XL-5000 ==ES-34-"

Tokenizer to Filter: "FooBar", "AstroBlaster", "XL-5000", "==ES-34-"

Out: "FooBar", "FooBar", "AstroBlaster", "XL-5000", "ES", "34"

Related Topics

- [TokenFilterFactories](#)



CharFilterFactories

Char Filter is a component that pre-processes input characters. Char Filters can be chained like Token Filters and placed in front of a Tokenizer. **Char Filters** can add, change, or remove characters without worrying about fault of Token offsets.

solr.MappingCharFilterFactory

This filter creates `org.apache.lucene.analysis.MappingCharFilter`, which can be used for changing one character to another (for example, for normalizing é to e.).

This filter requires specifying a `mapping` argument, which is the path and name of a file containing the mappings to perform.

Example:

```
<analyzer>
  <charFilter class="solr.MappingCharFilterFactory"
    mapping="mapping-ISOLatin1Accent.txt"/>
</analyzer>
```

solr.HTMLStripCharFilterFactory

This filter creates `org.apache.solr.analysis.HTMLStripCharFilter`. This Char Filter strips HTML from the input stream and passes the result to another Char Filter or a Tokenizer.



The behavior of this Char Filter changed in v3.6 of Solr. The old behavior has been retained in `LegacyHTMLCharFilterFactory`, although it has been deprecated and may be removed in a future release. What follows is a description of how the filter works in Solr 4; for an overview of the changes made, see Major Changes from Solr 3 to Solr 4.

This filter:

- Removes HTML/XML tags while preserving other content.
- Removes attributes within tags and supports optional attribute quoting.
- Removes XML processing instructions, such as: `<?foo bar?>`
- Removes XML comments.
- Removes XML elements starting with `<!>`.
- Removes contents of `<script>` and `<style>` elements.

- Handles XML comments inside these elements (normal comment processing will not always work).
- Replaces numeric character entities references like `A` or ``.
- The terminating `';` is optional if the entity reference is followed by whitespace.
- Replaces all named character entity references.
- ` s;` is replaced with a space instead of `0xa0`.
- The terminating `';` is mandatory to avoid false matches on something like "Alpha&Omega Corp".
- Newlines are substituted for block-level elements.
- `<CDATA>` sections are recognized.
- Inline tags, such as ``, `<i>`, or `` will be replaced by a space.
- Uppercase character entities like `quot`, `gt`, `lt` and `amp` are recognized and handled as lowercase.

 The input need not be an HTML document. The filter removes only constructs that look like HTML. If the input doesn't include anything that looks like HTML, the filter won't remove any input.

The table below presents examples of HTML stripping.

Input	Output
<code>my link</code>	<code>my link</code>
<code>
hello<!--comment--></code>	<code>hello</code>
<code>hello<script><!-- f('<!--internal--&gt;&lt;/script&gt;'); --&gt;&lt;/script&gt;</code></code>	<code>hello</code>
<code>if a<b then print a;</code>	<code>if a<b then print a;</code>
<code>hello <td height=22 nowrap align="left"></code>	<code>hello</code>
<code>a<b &#65 Alpha&Omega</code>	<code>a<b A Alpha&Omega Ω</code>

solr.LegacyHTMLStripCharFilterFactory

This filter strips HTML from the input stream and passes the result to another Char Filter or a Tokenizer. It has been **deprecated** in favor of improvements introduced in Solr 3.6 to the `HTMLCharFilter`. This filter creates `org.apache.solr.analysis.LegacyHTMLStripCharFilter`.

In Solr versions 3.5 and earlier, this filter had known bugs in the character offsets it provided, triggering (for example) exceptions in highlighting. With version 3.6, `HTMLStripCharFilter` was fixed to address this and other issues. If you depend on the behavior of `HTMLStripCharFilter` in version 3.5 or earlier, the previous implementation, including bugs, is preserved as `LegacyHTMLStripCharFilter`. For more information on the changes, see the section Major Changes from Solr 3 to Solr 4.

solr.PatternReplaceCharFilterFactory

This filter uses [regular expressions](#) to replace or change character patterns.

Arguments:

`pattern`: the regular expression pattern to apply to the incoming text.

`replaceWith`: the text to use to replace matching patterns.

You can configure this filter in `schema.xml` like this:

```
<analyzer>
  <charFilter class="solr.PatternReplaceCharFilterFactory"
    pattern="([nN][oO]\.)\s*(\d+)" replaceWith="$1$2"/>
</analyzer>
```

The table below presents examples of regex-based pattern replacement:

Input	pattern	replaceWith	Output	Description
see-ing looking	(\w+)(ing)	1	see-ing look	Removes "ing" from the end of word.
see-ing looking	(\w+)ing	1	see-ing look	Same as above. 2nd parentheses can be omitted.
No.1 NO. no. 543	[nN][oO] \.\s*(\d+)	{#},1	#1 NO. #543	Example of literal. Do not forget to set a non-period <code>blockDelimiter</code> when using periods in patterns.
abc=1234=5678	(\w+)=(\d+)=(\d+)	3,{=},1,{=},2	5678=abc=1234	Change the order of the groups.

Related Topics

- [CharFilterFactories](#)



Language Analysis

This section contains information about tokenizers and filters related to character set conversion or for use with specific languages. For the European languages, tokenization is fairly straightforward. Tokens are delimited by white space and/or a relatively small set of punctuation characters. In other languages the tokenization rules are often not so simple. Some European languages may require special tokenization rules as well, such as rules for decompounding German words.

For information about language detection at index time, see [Detecting Languages During Indexing](#).

Topics discussed in this section:

- [KeyWordMarkerFilterFactory](#)
- [StemmerOverrideFilterFactory](#)
- [Dictionary Compound Word Token Filter](#)
- [Unicode Collation](#)
- [ISO Latin Accent Filter](#)
- [Language-Specific Factories](#)
- [Related Topics](#)

KeyWordMarkerFilterFactory

Protects words from being modified by stemmers. A customized protected word list may be specified with the "protected" attribute in the schema. Any words in the protected word list will not be modified by any stemmer in Solr.

A sample Solr protwords.txt with comments can be found in the `/solr/conf/` directory:

```
<fieldtype name="myfieldtype" class="solr.TextField">
  <analyzer>
    <tokenizer class="solrWhitespaceTokenizerFactory"/>
    <filter class="solr.KeyWordMarkerFilterFactory" protected="protwords.txt" />
    <filter class="solr.PorterStemFilterFactory" />
  </analyzer>
</fieldtype>
```

StemmerOverrideFilterFactory

Overrides stemming algorithms by applying a custom mapping, then protecting these terms from being modified by stemmers.

A customized mapping of words to stems, in a tab-separated file, can be specified to the "dictionary" attribute in the schema. Words in this mapping will be stemmed to the stems from the file, and will not be further changed by any stemmer.

A sample [stemdict.txt](#) with comments can be found in the Source Repository.

```
<fieldtype name="myfieldtype" class="solr.TextField">
  <analyzer>
    <tokenizer class="solrWhitespaceTokenizerFactory"/>
    <filter class="solr.StemmerOverrideFilterFactory" dictionary="stemdict.txt" />
    <filter class="solr.PorterStemFilterFactory" />
  </analyzer>
</fieldtype>
```

Dictionary Compound Word Token Filter

This filter splits, or *decompounds*, compound words into individual words using a dictionary of the component words. Each input token is passed through unchanged. If it can also be decompounded into subwords, each subword is also added to the stream at the same logical position.

Compound words are most commonly found in Germanic languages.

Factory class: solr.DictionaryCompoundWordTokenFilterFactory

Arguments:

dictionary: (required) The path of a file that contains a list of simple words, one per line. Blank lines and lines that begin with "#" are ignored. This path may be an absolute path, or path relative to the Solr config directory.

minWordSize: (integer, default 5) Any token shorter than this is not decompounded.

minSubwordSize: (integer, default 2) Subwords shorter than this are not emitted as tokens.

maxSubwordSize: (integer, default 15) Subwords longer than this are not emitted as tokens.

onlyLongestMatch: (true/false) If true (the default), only the longest matching subwords will generate new tokens.

Example:

Assume that `germanwords.txt` contains at least the following words:

dummkopfdonaudampfschiff

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.DictionaryCompoundWordTokenFilterFactory"
dictionary="germanwords.txt" />
</analyzer>
```

In: "Donaudampfschiff dummkopf"

Tokenizer to Filter: "Donaudampfschiff"(1), "dummkopf"(2),

Out: "Donaudampfschiff"(1), "Donau"(1), "dampf"(1), "schiff"(1), "dummkopf"(2), "dumm"(2), "kopf"(2)

Unicode Collation

Unicode Collation is a language-sensitive method of sorting text that also be used for advanced search purposes.

Unicode Collation in Solr is fast, because all the work is done at index time. It uses a `KeywordTokenizerFactory` to create a sort field, followed by `CollationKeyFilterFactory`. The `CollationKeyFilterFactory` adds "sort keys" to the `sort` field at index time, so that at query time you can sort on the `sort` field and your results comes back in collated order.

You can also name `CollatedField` and `ICUCollatedField` to hold the results of your collation.

Sorting Text for a Specific Language

In this example, text is sorted according to the default German rules provided by Java. The rules for sorting German in Java are defined in a package called a Java Locale.

Locales are typically defined as a combination of language and country, but you can specify just the language if you want. For example, if you specify "de" as the language, you will get sorting that works well for German language. If you specify "de" as the language and "CH" as the country, you will get German sorting specifically tailored for Switzerland.

You can see a list of supported Locales [here](#).

```
<!-- define a field type for German collation -->
<fieldType name="collatedGERMAN" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.KeywordTokenizerFactory" />
    <filter class="solr.CollationKeyFilterFactory"
      language="de"
      strength="primary"
    />
  </analyzer>
</fieldType>
...
<!-- define a field to store the German collated manufacturer names -->
<field name="manuGERMAN" type="collatedGERMAN" indexed="true" stored="false" />
...
<!-- copy the text to this field. we could create French, English, Spanish versions
too,
      and sort differently for different users! --
<copyField source="manu" dest="manuGERMAN"/>
```

In the example above, we defined the strength as "primary". The strength of the collation determines how strict the sort order will be, but it also depends upon the language. For example, in English, "primary" strength ignores differences in case and accents.

For more information, see the [Collator javadocs](#).

Sorting Text for Multiple Languages

There are two approaches to supporting multiple languages: if there is a small list of languages you wish to support, consider defining collated fields for each language and using `copyField`. However, adding a large number of sort fields can increase disk and indexing costs. An alternative approach is to use the Unicode default collator.

The Unicode default or `ROOT` locale has rules that are designed to work well for most languages. To use the default locale, simply define the language as the empty string. This Unicode default sort is still significantly more advanced than the standard Solr sort.

```
<fieldType name="collatedROOT" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.KeywordTokenizerFactory"/>
    <filter class="solr.CollationKeyFilterFactory"
      language=""
      strength="primary"
    />
  </analyzer>
</fieldType>
```

Sorting Text with Custom Rules

You can define your own set of sorting rules. Its easiest to take existing rules that are close to what you want and customize them.

In the example below, we create a custom rule set for German called DIN 5007-2. This rule set treats umlauts in German differently: it treats ö as equivalent to oe. For more information, see the [RuleBasedCollator javadocs](#).

This example shows how to create a custom rule set and dump it to a file:

```
// get the default rules for Germany
// these are called DIN 5007-1 sorting
RuleBasedCollator baseCollator = (RuleBasedCollator) Collator.getInstance(new
Locale("de", "DE"));

// define some tailorings, to make it DIN 5007-2 sorting.
// For example, this makes ö equivalent to oe
String DIN5007_2_tailorings =
"& ae , a\u0308 & AE , A\u0308"+
"& oe , o\u0308 & OE , O\u0308"+
"& ue , u\u0308 & UE , u\u0308";

// concatenate the default rules to the tailorings, and dump it to a String
RuleBasedCollator tailoredCollator = new RuleBasedCollator(baseCollator.getRules() +
DIN5007_2_tailorings);
String tailoredRules = tailoredCollator.getRules();
// write these to a file, be sure to use UTF-8 encoding!!!
IOUtils.write(tailoredRules, new FileOutputStream("/solr_home/conf/customRules.dat"),
"UTF-8");
```

This rule set can now be used for custom collation in Solr:

```
<fieldType name="collatedCUSTOM" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.KeywordTokenizerFactory"/>
    <filter class="solr.CollationKeyFilterFactory"
      custom="customRules.dat"
      strength="primary"
    />
  </analyzer>
</fieldType>
```

Searching

Collation can also be used to search on a tokenized field.

In this example, we use the same custom German rules defined above on a tokenized field. As with stemmers, although the output tokens are nonsense they are the same values and will match for search purposes.

```
<fieldType name="collatedCUSTOM" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.CollationKeyFilterFactory"
      custom="customRules.dat"
      strength="primary"
    />
  </analyzer>
</fieldType>
```

Collation Key Filter

The filter `solr.CollationKeyFilter` is used at index time, indexing special “sort keys” into the sort field. It lets you choose the collator related to the target country and language. You can also choose the strength of the collation which determines the minimum level of difference considered significant during comparison. For example:

```
<filter class="solr.CollationKeyFilterFactory" language="es" country="ES"
strength="primary" />
```

The example above shows the configuration of the `CollationKeyFilterFactory`, where we want to handle the Spanish language with primary strength.

You can add the filter into field type definitions, as in the example below:

```
<fieldType name="polishLowercase" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.KeywordTokenizerFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.TrimFilterFactory"/>
    <filter class="solr.CollationKeyFilterFactory: language="pl" country="PL"
strength="primary"/>
  </analyzer>
</fieldType>
```

Handling the Polish language has been added to the definition of the currently existing lowercase type. The type will be used for the fields, where the data contains Polish signs. For example, you could also change the type for the `city_sort` field to `polishLowercase`:

```
<field name="city_sort" type="polishLowercase" indexed="true" stored="false" />
```

You can check the test query result:

```
q=*&fl=city&sort=city_sort+asc
```

And the result may look like this:

```
<result name="response" numFound="6" start="0">
  <doc>
    <str name="city">Biaystok</str>
  </doc>
  <doc>
    <str name="city">Koszalin</str>
  </doc>
  <doc>
    <str name="city">owicz</str>
  </doc>
  <doc>
    <str name="city">Szczecin</str>
  </doc>
  <doc>
    <str name="city">widnik</str>
  </doc>
  <doc>
    <str name="city">Warszawa</str>
  </doc>
</result>
```

ICU Collation

For better performance, less memory usage, and support for more locales, you can add the analysis-extras contrib and use `ICUCollationKeyFilterFactory` instead. See the [javadocs](#) for more information.

The principles of ICU Collation are the same as those of Unicode Collation; you just specify an RFC3066 language identifier with the locale parameter instead of specifying language+country+variant.

For example, to get German phonebook sort order:

```
<fieldType name="collatedICU" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.KeywordTokenizerFactory"/>
    <filter class="solr.ICUCollationKeyFilterFactory"
      locale="de@collation=phonebook"
      strength="primary"
    />
  </analyzer>
</fieldType>
```

To use the `ICUCollationKeyFilterFactory` filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `SOLR_HOME/lib`.

ISO Latin Accent Filter

This filter replaces any accented characters in a token with the unaccented equivalent. This can increase recall by causing more matches. On the other hand, it can reduce precision because language-specific character differences may be lost.

Characters in the ISO Latin 1 (ISO-8859-1) character set are recognized and letter case will be preserved, so that "Ã" becomes "A" and "á" becomes "a".

 This filter only looks for accented characters, it does not filter out other non-ASCII characters.

Factory class: `solr.ISOLatin1AccentFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.ISOLatin1AccentFilterFactory" />
</analyzer>
```

In: "Björn Ångström"

Tokenizer to Filter: "Björn", "Ångström"

Out: "Bjorn", "Angstrom"

Language-Specific Factories

These factories are each designed to work with specific languages. The languages covered here are:

- Arabic
- Brazilian Portuguese
- Bulgarian
- Chinese
- Simplified Chinese
- CJK
- Czech
- Dutch
- Finnish
- French
- Galician
- German
- Greek
- Hindi
- Indonesian
- Italian
- Kuromoji (Japanese)
- Lao, Myanmar, Khmer
- Latvian
- Norwegian
- Persian
- Polish
- Portuguese
- Russian
- Spanish
- Swedish
- Thai
- Turkish

Arabic

Solr provides support for the [Light-10](#) (PDF) stemming algorithm, and Lucene includes an example stopword list.

This algorithm defines both character normalization and stemming, so these are split into two filters to provide more flexibility.

Factory classes: solr.ArabicStemFilterFactory, solr.ArabicNormalizationFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.ArabicNormalizationFilterFactory"/>
  <filter class="solr.ArabicStemFilterFactory"/>
</analyzer>
```

In: لوتق م

Tokenizer to Filter: لوتق، لتق

Out: لتق لتق

Brazilian Portuguese

This is a Java filter written specifically for stemming the Brazilian dialect of the Portuguese language. It uses the Lucene class `org.apache.lucene.analysis.br.BrazilianStemmer`. Although that stemmer can be configured to use a list of protected words (which should not be stemmed), this factory does not accept any arguments to specify such a list.

Factory class: `solr.BrazilianStemFilterFactory`

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.BrazilianStemFilterFactory"/>
</analyzer>
```

In: "praia praias"

Tokenizer to Filter: "praia", "praias"

Out: "pra", "pra"

Bulgarian

Solr includes a light stemmer for Bulgarian, following [this algorithm](#) (PDF), and Lucene includes an example stopword list.

Factory class: `solr.BulgarianStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.BulgarianStemFilterFactory"/>
</analyzer>
```

In: "вания ване ването"

Tokenizer to Filter: "вания", "ване", "ването"

Out: "ван", "ван", "ван"

Chinese

Chinese Tokenizer

The Chinese Tokenizer is deprecated as of Solr 3.4. Use the [solr.StandardTokenizerFactory](#) instead.

Factory class: solr.ChineseTokenizerFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.ChineseTokenizerFactory"/>
</analyzer>
```

In: "你好，我不讲中文"

Out: "你", "好", "我", "不", "讲", "中", "文"

Chinese Filter Factory

The Chinese Filter Factory is deprecated as of Solr 3.4. Use the [solr.StopFilterFactory](#) instead.

Factory class: solr.ChineseFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.ChineseFilterFactory"/>
</analyzer>
```

In: "你好, and 我不讲中文"

Tokenizer to Filter: "你", "好", "and", "我", "不", "讲", "中", "文"

Out: "你", "好", "我", "不", "讲", "中", "文"

Simplified Chinese

For Simplified Chinese, Solr provides support for Chinese sentence and word segmentation with the `solr.SmartChineseSentenceTokenFilterFactory` and `solr.SmartChineseWordTokenFilterFactory` in the `analysis-extras contrib` module. This component includes a large dictionary and segments Chinese text into words with the Hidden Markov Model. To use this filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

Factory class: `solr.SmartChineseWordTokenFilterFactory`

Arguments: None

Examples:

To use the default setup with fallback to English Porter stemmer for english words, use:

```
<analyzer class="org.apache.lucene.analysis.cn.smart.SmartChineseAnalyzer"/>
```

Or to configure your own analysis setup, use the `SmartChineseSentenceTokenizerFactory` along with your custom filter setup. The sentence tokenizer tokenizes on sentence boundaries and the `SmartChineseWordTokenFilter` breaks this further up into words.

```
<analyzer>
  <tokenizer class="solr.SmartChineseSentenceTokenizerFactory"/>
  <filter class="solr.SmartChineseWordTokenFilterFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.PositionFilterFactory" />
</analyzer>
```

CJK

This tokenizer breaks Chinese, Japanese and Korean language text into tokens. These are not whitespace delimited languages. The tokens generated by this tokenizer are "doubles", overlapping pairs of CJK characters found in the field text.

Factory class: `solr.CJKTokenizerFactory`

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.CJKTokenizerFactory" />
</analyzer>
```

In: "你好 , 我不讲中文"

Out: "你好", "我不", "不讲", "讲中", "讲", "中文", ""

Czech

Solr includes a light stemmer for Czech, following [this algorithm](#), and Lucene includes an example stopword list.

Factory class: solr.CzechStemFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.LowerCaseFilterFactory" />
  <filter class="solr.CzechStemFilterFactory" />
</analyzer>
```

In: "prezidenští, prezidenta, prezidentského"

Tokenizer to Filter: "prezidenští", "prezidenta", "prezidentského"

Out: "preziden", "preziden", "preziden"

Dutch

This is a Java filter written specifically for stemming the Dutch language. It uses the Lucene class `org.apache.lucene.analysis.nl.DutchStemmer`. Although that stemmer can be configured to use a list of protected words (which should not be stemmed), this factory does not accept any arguments to specify such a list.

Another option for stemming Dutch words is to use the Snowball Porter Stemmer with an argument of `language="Dutch"`.

Factory class: solr.DutchStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.DutchStemFilterFactory" />
</analyzer>
```

In: "kanaal kanalen"

Tokenizer to Filter: "kanaal", "kanalen"

Out: "kanal", "kanal"

Finnish

Solr includes support for stemming Finnish, and Lucene includes an example stopword list.

Factory class: solr.FinnishLightStemFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.FinnishLightStemFilterFactory" />
</analyzer>
```

In: "kala kalat"

Tokenizer to Filter: "kala", "kalat"

Out: "kala", "kala"

French

Elision Filter

Removes article elisions from a token stream. This filter primarily applies to the French language and makes use of the ElisionFilter class in `org.apache.lucene.analysis.fr`.

Factory class: solr.ElisionFilterFactory

Arguments:

`articles`: (required) The pathname of a file that contains a list of articles, one per line, to be stripped. Articles are words such as "le", which are commonly abbreviated, such as *l'avion* (the plane). This file should include the abbreviated form, which precedes the apostrophe. In this case, simply "/".

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.ElisionFilterFactory" />
</analyzer>
```

In: "L'histoire d'art"

Tokenizer to Filter: "L'histoire", "d'art"

Out: "histoire", "art"

French Light Stem Filter

Solr includes three stemmers for French: one in the `solr.SnowballPorterFilterFactory`, a lighter stemmer called `solr.FrenchLightStemFilterFactory`, and an even less aggressive stemmer called `solr.FrenchMinimalStemFilterFactory`. Lucene includes an example stopword list.

Factory classes: `solr.FrenchLightStemFilterFactory`, `solr.FrenchMinimalStemFilterFactory`

Arguments: None

Examples:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.LowerCaseFilterFactory" />
  <filter class="solr.ElisionFilterFactory" />
  <filter class="solr.FrenchLightStemFilterFactory" />
</analyzer>
```

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.LowerCaseFilterFactory" />
  <filter class="solr.ElisionFilterFactory" />
  <filter class="solr.FrenchMinimalStemFilterFactory" />
</analyzer>
```

In: "le chat, les chats"

Tokenizer to Filter: "le", "chat", "les", "chats"

Out: "le", "chat", "le", "chat"

Galician

Solr includes a stemmer for Galician following [this algorithm](#), and Lucene includes an example stopword list.

Factory class: solr.GalicianStemFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.GalicianStemFilterFactory"/>
</analyzer>
```

In: "felizmente Luzes"

Tokenizer to Filter: "felizmente", "luzes"

Out: "feliz", "luz"

German

Solr includes four stemmers for German: one in the `solr.SnowballPorterFilterFactory` language="German", a stemmer called `solr.GermanStemFilterFactory`, a lighter stemmer called `solr.GermanLightStemFilterFactory`, and an even less aggressive stemmer called `solr.GermanMinimalStemFilterFactory`. Lucene includes an example stopword list.

Factory classes: `solr.GermanStemFilterFactory`, `solr.LightGermanStemFilterFactory`, `solr.MinimalGermanStemFilterFactory`

Arguments: None

Examples:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.GermanStemFilterFactory"/>
</analyzer>
```

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.GermanLightStemFilterFactory"/>
</analyzer>
```

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.GermanMinimalStemFilterFactory" />
</analyzer>
```

In: "hund hunden"

Tokenizer to Filter: "hund", "hunden"

Out: "hund", "hund"

Greek

This filter converts uppercase letters in the Greek character set to the equivalent lowercase character.

Factory class: solr.GreekLowerCaseFilterFactory

Arguments:

charset: (optional, default "UnicodeGreek") Specifies the name of the character set to use. Must be "UnicodeGreek", "ISO" or "CP1253".

 Use of custom charsets was deprecated in Solr 1.4 and is unsupported in Solr 3.1. If you need to index text in these encodings, please use Java's character set conversion facilities (InputStreamReader, and so on.) during I/O, so that Lucene can analyze this text as Unicode instead.

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.GreekLowerCaseFilterFactory" />
</analyzer>
```

In: "Ελληνική Δημοκρατία Ellīnikī Dīmokratía"

Tokenizer to Filter: "Ελληνική", "Δημοκρατία", "Ellīnikī", "Dīmokratía"

Out: "ελληνικη", "δημοκρατια", "ellīnikī", "dīmokratía"

Hindi

Solr includes support for stemming Hindi following [this algorithm](#) (PDF), support for common spelling differences through the `solr.HindiNormalizationFilterFactory`, support for encoding differences through the `solr.IndicNormalizationFilterFactory` following [this algorithm](#), and Lucene includes an example stopword list.

Factory classes: `solr.IndicNormalizationFilterFactory`, `solr.HindiNormalizationFilterFactory`, `solr.HindiStemFilterFactory`

Arguments: None

Example:

```
<filter class="solr.IndicNormalizationFilterFactory" />
<filter class="solr.HindiNormalizationFilterFactory" />
<filter class="solr.HindiStemFilterFactory" />
```

In: "बैवहना गुसपेतयोम"

Out: "बैवहन", "गुसपेट"

Indonesian

Solr includes support for stemming Indonesian (Bahasa Indonesia) following [this algorithm](#) (PDF), and Lucene includes an example stopword list.

Factory class: `solr.IndonesianStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.LowerCaseFilterFactory" />
  <filter class="solr.IndonesianStemFilterFactory" stemDerivational="true" />
</analyzer>
```

In: "sebagai sebagainya"

Tokenizer to Filter: "sebagai", "sebagainya"

Out: "bagai", "bagai"

Italian

Solr includes two stemmers for Italian: one in the `solr.SnowballPorterFilterFactory` language="Italian", and a lighter stemmer called `solr.ItalianLightStemFilterFactory`. Lucene includes an example stopword list.

Factory class: `solr.ItalianStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.ItalianLightStemFilterFactory"/>
</analyzer>
```

In: "propaga propagare propagamento"

Tokenizer to Filter: "propaga", "propagare", "propagamento"

Out: "propag", "propag", "propag"

Kuromoji (Japanese)

Solr includes support for stemming Kuromoji (Japanese), and Lucene includes an example stopword list. Kuromoji has a search mode (default) that does segmentation useful for search. A heuristic is used to segment compounds into its parts and the compound itself is kept as a synonym.

With Solr 4, the `JapaneseIterationMarkCharFilterFactory` now is included to normalize Japanese iteration marks.

You can also make discarding punctuation configurable in the `JapaneseTokenizerFactory`, by setting `discardPunctuation` to `false` (to show punctuation) or `true` (to discard punctuation), as in the following example:

Factory class: `solr.KuromojiStemFilterFactory`

Arguments:

`mode`: Use search-mode to get a noun-decompounding effect useful for search. Search mode improves segmentation for search at the expense of part-of-speech accuracy. Valid values for `mode` are:

- `normal`: default segmentation
- `search`: segmentation useful for search (extra compound splitting)
- `extended`: search mode with unigramming of unknown words (experimental)

For some applications it might be good to use search mode for indexing and normal mode for queries to reduce recall and prevent parts of compounds from being matched and highlighted.

Kuromoji also has a convenient user dictionary feature that allows overriding the statistical model with your own entries for segmentation, part-of-speech tags and readings without a need to specify weights. Note that user dictionaries have not been subject to extensive testing. User dictionary attributes are:

`userDictionary: user dictionary filename`

`userDictionaryEncoding: user dictionary encoding (default is UTF-8)`

See `lang/userdict_ja.txt` for a sample user dictionary file.

Punctuation characters are discarded by default. Use `discardPunctuation="false"` to keep them.

Example:

```
<fieldType name="text_ja" positionIncrementGap="100" autoGeneratePhraseQueries="false">
  <analyzer>
    <tokenizer class="solr.JapaneseTokenizerFactory" mode="search"
      userDictionary="lang/userdict_ja.txt"/>
    <filter class="solr.JapaneseBaseFormFilterFactory"/>
    <filter class="solr.JapanesePartOfSpeechStopFilterFactory"
      tags="lang/stoptags_ja.txt" enablePositionIncrements="true"/>
    <filter class="solr.CJKWidthFilterFactory"/>
    <filter class="solr.StopFilterFactory" ignoreCase="true"
      words="lang/stopwords_ja.txt" enablePositionIncrements="true" />
    <filter class="solr.JapaneseKatakanaStemFilterFactory" minimumLength="4"/>
    <filter class="solr.LowerCaseFilterFactory"/>
  </analyzer>
</fieldType>
```

Lao, Myanmar, Khmer

Lucene provides support for segmenting these languages into syllables with the `solr.ICUTokenizerFactory` in the `analysis-extras` contrib module. To use this tokenizer, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

Latvian

Solr includes support for stemming Latvian, and Lucene includes an example stopword list.

Factory class: `solr.LatvianStemFilterFactory`

Arguments: None

Example:

```
<fieldType name="text_lvstem" class="solr.TextField" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.LatvianStemFilterFactory"/>
  </analyzer>
</fieldType>
```

In: "tirgiem tirgus"**Tokenizer to Filter:** "tirgiem", "tirgus"**Out:** "tirg", "tirg"

Norwegian

Solr includes two classes for stemming Norwegian, `NorwegianLightStemFilterFactory` and `NorwegianMinimalStemFilterFactory`. Lucene includes an example stopword list.

Norwegian Light Stemmer

The `NorwegianLightStemFilterFactory` requires a "two-pass" sort for the -dom and -het endings. This means that in the first pass the word "kristendom" is stemmed to "kristen", and then all the general rules apply so it will be further stemmed to "krist". The effect of this is that "kristen," "kristendom," "kristendommen," and "kristendommens" will all be stemmed to "krist."

The second pass is to pick up -dom and -het endings. Consider this example:

One pass		Two passes	
Before	After	Before	After
forlegen	forleg	forlegen	forleg
forlegenhet	forlegen	forlegenhet	forleg
forlegenheten	forlegen	forlegenheten	forleg
forlegenhetens	forlegen	forlegenhetens	forleg
firkantet	firkant	firkantet	firkant
firkantethet	firkantet	firkantethet	firkant
firkantetheten	firkantet	firkantetheten	firkant

Factory class: `solr.NorwegianLightStemFilterFactory`

Arguments: None

Example:

```
<fieldType name="text_no" class="solr.TextField" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.StopFilterFactory" ignoreCase="true"
words="lang/stopwords_no.txt" format="snowball" enablePositionIncrements="true"/>
    <filter class="solr.SnowballPorterFilterFactory" language="Norwegian"/>
    <filter class="solr.NorwegianLightStemFilterFactory"/>
  </analyzer>
</fieldType>
```

In: "Forelskelsen"

Tokenizer to Filter: "forelskelsen"

Out: "forelske"

Norwegian Minimal Stemmer

The NorwegianMinimalStemFilterFactory stems plural forms of Norwegian nouns only.

Factory class: solr.NorwegianMinimalStemFilterFactory

Arguments: None

Example:

```
<fieldType name="text_no" class="solr.TextField" positionIncrementGap="100">
  <analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.StopFilterFactory" ignoreCase="true"
words="lang/stopwords_no.txt" format="snowball" enablePositionIncrements="true"/>
    <filter class="solr.SnowballPorterFilterFactory" language="Norwegian"/>
    <filter class="solr.NorwegianMinimalStemFilterFactory"/>
  </analyzer>
</fieldType>
```

In: "Bilens"

Tokenizer to Filter: "bilens"

Out: "bil"

Persian

Persian Filter Factories

Solr includes support for normalizing Persian, and Lucene includes an example stopword list.

Factory class: solr.PersianNormalizationFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solr.ArabicNormalizationFilterFactory"/>
  <filter class="solr.PersianNormalizationFilterFactory">
</analyzer>
```

In: "اهگرب اه گرب"

Tokenizer to Filter: "گرب", "اه", "اهگرب"

Out: "گرب", "گ"

Polish

Solr provides support for Polish stemming with the `solr.StempelPolishStemFilterFactory` in the `contrib/analysis-extras` module. This component includes an algorithmic stemmer with tables for Polish. To use this filter, see `solr/contrib/analysis-extras/README.txt` for instructions on which jars you need to add to your `solr_home/lib`.

Factory class: solr.StempelPolishStemFilterFactory

Arguments: None

Example:

```
<analyzer>
  <filter class="solrLowerCaseFilterFactory"/>
  <filter class="solr.solr.StempelPolishStemFilterFactory"/>
</analyzer>
```

In: ""studenta studenci"

Tokenizer to Filter: "studenta", "studenci"

Out: "student", "student"

More information about the Stempel stemmer is available in the Lucene javadocs, https://lucene.apache.org/core/4_0_0/analyzers-stempel/index.html.

Portuguese

Solr includes four stemmers for Portuguese: one in the `solr.SnowballPorterFilterFactory`, an alternative stemmer called `solr.PortugueseStemFilterFactory`, a lighter stemmer called `solr.PortugueseLightStemFilterFactory`, and an even less aggressive stemmer called `solr.PortugueseMinimalStemFilterFactory`. Lucene includes an example stopword list.

Factory class: `solr.PortugueseStemFilterFactory`, `solr.PortugueseLightStemFilterFactory`, `solr.PortugueseMinimalStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.PortugueseStemFilterFactory"/>
</analyzer>
```

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.PortugueseLightStemFilterFactory"/>
</analyzer>
```

```
<analyzer>
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.PortugueseMinimalStemFilterFactory"/>
</analyzer>
```

In: "praia praias"

Tokenizer to Filter: "praia", "praias"

Out: "pra", "pra"

Russian

Russian Letter Tokenizer

This tokenizer breaks Russian language text into tokens. It is similar to `LetterTokenizer`, but additionally looks up letters in the appropriate Russian character set.

Factory class: `solr.RussianLetterTokenizerFactory`

Arguments:

`charset`: (optional, default "UnicodeRussian") The name of the character set to use. Must be "UnicodeRussian", "KOI8" or "CP1251".

 Use of custom charsets was deprecated in Solr 1.4 and is unsupported in Solr 3.1. If you need to index text in these encodings, please use Java's character set conversion facilities (`InputStreamReader`, and so on.) during I/O, so that Lucene can analyze this text as Unicode instead.

Example:

```
<analyzer type="index">
  <tokenizer class="solr.RussianLetterTokenizerFactory" />
</analyzer>
```

In: "Здравствуйте!. Я не говорю русского."

Out: "Здравствуйте", "Я", "не", "говорю", "русского"

Russian Lower Case Filter

This filter converts uppercase letters in the Russian character set to the equivalent lowercase character.

Factory class: `solr.RussianLowerCaseFilterFactory`

Arguments:

`charset`: (optional, default "UnicodeRussian") Specifies the name of the character set to use. Must be "UnicodeRussian", "KOI8" or "CP1251".

 Use of custom charsets was deprecated in Solr 1.4 and is unsupported in Solr 3.1. If you need to index text in these encodings, please use Java's character set conversion facilities (`InputStreamReader`, and so on.) during I/O, so that Lucene can analyze this text as Unicode instead.

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.RussianLowerCaseFilterFactory" />
</analyzer>
```

In: "Здравствуйте!. Я не говорю русского."

Tokenizer to Filter: "Здравствуйте", "Я", "не", "говорю", "русского"

Out: "здравствуйте", "я", "не", "говорю", "русского"

Russian Stem Filter

Solr includes two stemmers for Russian: one in the `solr.SnowballPorterFilterFactory` language="Russian", and a lighter stemmer called `solr.RussianLightStemFilterFactory`. Lucene includes an example stopword list.

Factory class: `solr.RussianLightStemFilterFactory`

Arguments:

charset: (optional, default "UnicodeRussian") Specifies the name of the character set to use. Must be "UnicodeRussian", "KOI8" or "CP1251".

 Use of custom charsets was deprecated in Solr 1.4 and is unsupported in Solr 3.4. If you need to index text in these encodings, please use Java's character set conversion facilities (InputStreamReader, and so on.) during I/O, so that Lucene can analyze this text as Unicode instead.

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory" />
  <filter class="solr.RussianLowerCaseFilterFactory" />
  <filter class="solr.RussianLightStemFilterFactory" />
</analyzer>
```

In: "вал валы"

Tokenizer to Filter: "вал", "валы"

Out: "вал", "вал"

Spanish

Solr includes two stemmers for Spanish: one in the `solr.SnowballPorterFilterFactory` `language="Spanish"`, and a lighter stemmer called `solr.SpanishLightStemFilterFactory`. Lucene includes an example stopword list.

Factory class: `solr.SpanishStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <filter class="solrLowerCaseFilterFactory"/>
  <filter class="solr.SpanishLightStemFilterFactory"/>
</analyzer>
```

In: "torear toreada torearlo"

Tokenizer to Filter: "torear", "toreada", "torearlo"

Out: "tor", "tor", "tor"

Swedish

Swedish Stem Filter

Solr includes two stemmers for Swedish: one in the `solr.SnowballPorterFilterFactory` `language="Swedish"`, and a lighter stemmer called `solr.SwedishLightStemFilterFactory`. Lucene includes an example stopword list.

Factory class: `solr.SwedishStemFilterFactory`

Arguments: None

Example:

```
<analyzer>
  <filter class="solrLowerCaseFilterFactory"/>
  <filter class="solr.SwedishLightStemFilterFactory"/>
</analyzer>
```

In: "kloke klokhet klokheten"

Tokenizer to Filter: "kloke", "klokhet", "klokheten"

Out: "klok", "klok", "klok"

Thai

This filter converts sequences of Thai characters into individual Thai words. Unlike European languages, Thai does not use whitespace to delimit words.

Factory class: solr.ThaiWordFilterFactory

Arguments: None

Example:

```
<analyzer type="index">
  <tokenizer class="solr.StandardTokenizerFactory"/>
  <filter class="solr.ThaiWordFilterFactory"/>
</analyzer>
```

In: "ช้างสีบสามเชือก"

Tokenizer to Filter: "ช้างสีบสามเชือก"

Out: "ช้าง", "สีบ", "สาม", "เชือก"

Turkish

Solr includes support for stemming Turkish through the `solr.SnowballPorterFilterFactory`, as well as support for case-insensitive search through the `solr.TurkishLowerCaseFilterFactory`, and Lucene includes an example stopword list.

Factory class: solr.TurkishLowerCaseFilterFactory

Arguments: None

Example:

```
<filter class="solr.TurkishLowerCaseFilterFactory"/>
  <filter class="solr.SnowballPorterFilterFactory" language="Turkish" />
```

Related Topics

- [LanguageAnalysis](#)



Phonetic Matching

Introduced with Solr v3.6, Beider-Morse Phonetic Matching (BMPM) is a "soundalike" tool that lets you search using a new phonetic matching system. BMPM helps you search for personal names (or just surnames) in a Solr/Lucene index, and is far superior to the existing phonetic codecs, such as regular soundex, metaphone, caverphone, etc.

In general, phonetic matching lets you search a name list for names that are phonetically equivalent to the desired name. BMPM is similar to a soundex search in that an exact spelling is not required. Unlike soundex, it does not generate a large quantity of false hits.

From the spelling of the name, BMPM attempts to determine the language. It then applies phonetic rules for that particular language to transliterate the name into a phonetic alphabet. If it is not possible to determine the language with a fair degree of certainty, it uses generic phonetic instead. Finally, it applies language-independent rules regarding such things as voiced and unvoiced consonants and vowels to further insure the reliability of the matches.

For example, assume that the matches found when searching for Stephen in a database are "Stefan", "Steph", "Stephen", "Steve", "Steven", "Stove", and "Stuffin". "Stefan", "Stephen", and "Steven" are probably relevant, and are names that you want to see. "Stuffin", however, is probably not relevant. Also rejected were "Steph", "Steve", and "Stove". Of those, "Stove" is probably not one that we would have wanted. But "Steph" and "Steve" are possibly ones that you might be interested in.

For Solr, BMPM searching is available for the following languages:

- English
- French
- German
- Greek
- Hebrew written in Hebrew letters
- Hungarian
- Italian
- Lithuanian and Latvian
- Polish
- Romanian
- Russian written in Cyrillic letters
- Russian transliterated into English letters
- Spanish
- Turkish

The name matching is also applicable to non-Jewish surnames from the countries in which those languages are spoken.

For more information, see here: <http://stevemorse.org/phoneticinfo.htm> and <http://stevemorse.org/phonetics/bmpm.htm..>

Running Your Analyzer

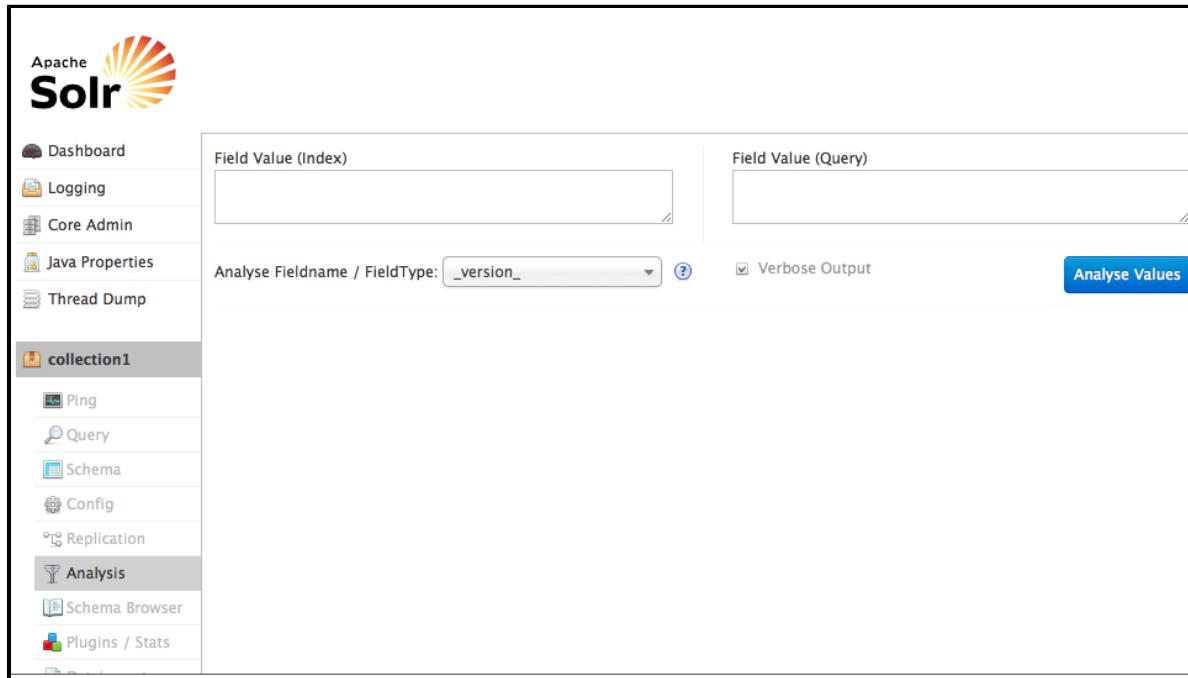
Once you've defined a field type in `schema.xml` and specified the analysis steps that you want applied to it, you should test it out to make sure that it behaves the way you expect it to. Luckily, there is a very handy page in the Solr [admin interface](#) that lets you do just that. You can invoke the analyzer for any text field, provide sample input, and display the resulting token stream.

For example, assume that the following field type definition has been added to `schema.xml`:

```
<fieldType name="mytextfield" class="solr.TextField">
    <analyzer type="index">
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.HyphenatedWordsFilterFactory"/>
        <filter class="solr.LowerCaseFilterFactory"/>
    </analyzer>
    <analyzer type="query">
        <tokenizer class="solr.StandardTokenizerFactory"/>
        <filter class="solr.LowerCaseFilterFactory"/>
    </analyzer>
</fieldType>
```

The objective here (during indexing) is to reconstruct hyphenated words, which may have been split across lines in the text, then to set all words to lowercase. For queries, you want to skip the de-hyphenation step.

To test this out, point your browser at the [Analysis Screen](#) of the Solr Admin Web interface. By default, this will be at the following URL (adjust the hostname and/or port to match your configuration): <http://localhost:8983/solr/#/collection1/analysis>. You should see a page like this.



Empty Analysis screen

We want to test the field type definition for "mytextfield", defined above. The drop-down labeled "Analyse Fieldname/FieldType" allows choosing the field or field type to use for the analysis.

There are two "Field Value" boxes, one for how text will be analyzed during indexing and a second for how text will be analyzed for query processing. In the "Field Value (Index)" box enter some sample text "Super-computer" (in this example) to be processed by the analyzer. We will leave the query field value empty for now.

The result we expect is that `HyphenatedWordsFilter` will join the hyphenated pair "Super-" and "computer" into the single word "Supercomputer", and then `LowerCaseFilter` will set it to "supercomputer". Let's see what happens:

	text	raw_bytes	start	end	position	type				
ST	text	Super	Computer	[53 75 70 65 72]	43 6f 6d 70 75 74 65 72]	0	7	15	2	<ALPHANUM>
HWF	text	Super	Computer	[53 75 70 65 72]	43 6f 6d 70 75 74 65 72]	0	7	15	2	<ALPHANUM>
LCF	text	super	computer	[73 75 70 65 72]	63 6f 6d 70 75 74 65 72]	1	2	7	15	<ALPHANUM>

Running index-time analyzer, verbose output.

The result is two distinct tokens rather than the one we expected. What went wrong? Looking at the first token that came out of StandardTokenizer, we can see the trailing hyphen has been stripped off of "Super-". Checking the documentation for StandardTokenizer, we see that it treats all punctuation characters as delimiters and discards them. What we really want in this case is a whitespace tokenizer that will preserve the hyphen character when it breaks the text into tokens.

Let's make this change and try again:

```

<fieldType name="mytextfield" class="solr.TextField">
    <analyzer type="index">
        <tokenizer class="solrWhitespaceTokenizerFactory"/>
        <filter class="solrHyphenatedWordsFilterFactory"/>
        <filter class="solrLowerCaseFilterFactory"/>
    </analyzer>
    <analyzer type="query">
        <tokenizer class="solrStandardTokenizerFactory"/>
        <filter class="solrLowerCaseFilterFactory"/>
    </analyzer>
</fieldType>

```

Re-submitting the form by clicking "Analyse Values" again, we see the result in the screen shot below.

The screenshot shows the Apache Solr Analysis screen. On the left, there's a sidebar with various collection and analysis-related links. The main area has two input fields: "Field Value (Index)" containing "Super-Computer" and "Field Value (Query)" which is currently empty. Below these is a dropdown menu set to "mytextfield". To the right of the dropdown is a checkbox for "Verbose Output" and a blue "Analyse Values" button. The central part of the screen displays three tables representing different analyzers:

	WT	text	raw_bytes	start	end	position	type	
Super-	text	Super-	[53 75 70 65 72 2d]	0	6	1	word	Computer
								[43 6f 6d 70 75 74 65 72]
				7		15		
					2			

	HWF	text	raw_bytes	start	end	position	type	
SuperComputer	text	SuperComputer	[53 75 70 65 72 43 6f 6d 70 75 74 65 72]	0	15	1	word	

	LCF	text	raw_bytes	position	start	end	type	
supercomputer	text	supercomputer	[73 75 70 65 72 63 6f 6d 70 75 74 65 72]	1	0	15	word	

Using WhitespaceTokenizer, expected results.

That's more like it. Because the whitespace tokenizer preserved the trailing hyphen on the first token, `HyphenatedWordsFilter` was able to reconstruct the hyphenated word, which then passed it on to `LowerCaseFilter`, where capital letters are set to lowercase.

Now let's see what happens when invoking the analyzer for query processing. For query terms, we don't want to do de-hyphenation and we *do* want to discard punctuation, so let's try the same input on it. We'll copy the same text to the "Field Value (Query)" box and clear the one for index analysis. We'll also include the full, unhyphenated word as another term to make sure it is processed to lower case as we expect. Submitting again yields these results:

The screenshot shows the Apache Solr Admin interface for a collection named "collection1". On the left, there's a sidebar with links like Dashboard, Logging, Core Admin, Java Properties, Thread Dump, collection1 (which is selected), Ping, Query, Schema, Config, Replication, Analysis (which is also selected), and Schema Browser. The main area has two text input fields: "Field Value (Index)" containing "Super-Computer Supercomputer" and "Field Value (Query)" containing "Super-Computer Supercomputer". Below these is a dropdown menu "Analyse Fieldname / FieldType:" set to "mytextfield" and a "Analyst Values" button. There are also "Verbose Output" and "Analyse Values" checkboxes. Two tables below show the tokenization details for the analyzed field:

ST	text	raw_bytes	Super	Computer	Supercomputer
		[53 75 70 65 72]	[43 6f 6d 70 75 74 65 72]	[53 75 70 65 72 6f 6d 70 75 74 65 72]	
0			6	15	
end			14	28	
type	<ALPHANUM>		<ALPHANUM>	<ALPHANUM>	
position	1		2	3	

LCF	text	raw_bytes	super	computer	supercomputer
		[73 75 70 65 72]	[63 6f 6d 70 75 74 65 72]	[73 75 70 65 72 6f 6d 70 75 74 65 72]	
			1	2	3
position			0	6	15
start			5	14	28
end				<ALPHANUM>	<ALPHANUM>
type	<ALPHANUM>				

Query-time analyzer, good results.

We can see that for queries the analyzer behaves the way we want it to. Punctuation is stripped out, `HyphenatedWordsFilter` doesn't run, and we wind up with the three tokens we expected.

Refer to the section [Running Field Analysis to Test Analyzers, Tokenizers, and TokenFilters](#) for more information about conducting field analysis through the Admin Web interface.

Indexing and Basic Data Operations

This section describes how Solr adds data to its index. It covers the following topics:

[What Is Indexing?](#): An overview of Solr's indexing process.

[Uploading Data with Solr Cell using Apache Tika](#): Information about using the Solr Cell framework to upload data for indexing.

[Uploading Data with Index Handlers](#): Information about using Solr's Index Handlers to upload XML and CSV data.

[Uploading Structured Data Store Data with the Data Import Handler](#): Information about uploading and indexing data from a structured data store.

[Detecting Languages During Indexing](#): Information about using language identification during the indexing process.

[UIMA Integration](#): Information about integrating Solr with Apache's Unstructured Information Management Architecture (UIMA). UIMA lets you define custom pipelines of Analysis Engines that incrementally add metadata to your documents as annotations.

[Content Streams](#): Information about streaming content to Solr Request Handlers.

What Is Indexing?

This section describes the process of indexing: adding content to a Solr index and, if necessary, modifying that content or deleting it. By adding content to an index, we make it searchable by Solr.

A Solr index can accept data from many different sources, including XML files, comma-separated value (CSV) files, data extracted from tables in a database, and files in common file formats such as Microsoft Word or PDF.

Here are the three most common ways of loading data into a Solr index:

- Using the [Solr Cell](#) framework built on Apache Tika for ingesting binary files or structured files such as Office, Word, PDF, and other proprietary formats.
- Uploading XML files by sending HTTP requests to the Solr server from any environment where such requests can be generated.
- Writing a custom Java application to ingest data through Solr's Java Client API (which is described in more detail in [Client APIs](#). Using the Java API may be the best choice if you're working with an application, such as a Content Management System (CMS), that offers a Java API.

Regardless of the method used to ingest data, there is a common basic data structure for data being fed into a Solr index: a *document* containing multiple *fields*, each with a *name* and containing *content*, which may be empty. One of the fields is usually designated as a unique ID field (analogous to a primary key in a database), although the use of a unique ID field is not strictly required by Solr.

If the field name is defined in the `schema.xml` file that is associated with the index, then the analysis steps associated with that field will be applied to its content when the content is tokenized. Fields that are not explicitly defined in the schema will either be ignored or mapped to a dynamic field definition (see [Documents, Fields, and Schema Design](#)), if one matching the field name exists.

For more information on indexing in Solr, see the [Solr Wiki](#).

The Solr Example Directory

The `example/` directory includes a sample Solr implementation, along with sample documents for uploading into an index. You will find the example docs in `$SOLR_HOME/example/exampledocs`.

The curl Utility for Transferring Files

Many of the instructions and examples in this section make use of the `curl` utility for transferring content through a URL. `curl` posts and retrieves data over HTTP, FTP, and many other protocols. Most Linux distributions include a copy of `curl`. You'll find `curl` downloads for Linux, Windows, and many other operating systems at <http://curl.haxx.se/download.html>. Documentation for `curl` is available here: <http://curl.haxx.se/docs/manpage.html>.



Using `curl` or other command line tools for posting data is just fine for examples or tests, but it's not the recommended method for achieving the best performance for updates in production environments. You will achieve better performance with Solr Cell or the other methods described in this section.

Instead of `curl`, you can use utilities such as GNU `wget` (<http://www.gnu.org/software/wget/>) or manage GETs and POSTS with Perl, although the command line options will differ.

Uploading Data with Solr Cell using Apache Tika

Solr uses code from the [Apache Tika](#) project to provide a framework for incorporating many different file-format parsers such as [Apache PDFBox](#) and [Apache POI](#) into Solr itself. Working with this framework, Solr's `ExtractingRequestHandler` can use Tika to support uploading binary files, including files in popular formats such as Word and PDF, for data extraction and indexing.

 In version 4, Solr uses Apache Tika v1.2.

When this framework was under development, it was called the Solr Content Extraction Library or CEL; from that abbreviation came this framework's name: Solr Cell.

If you want to supply your own [ContentHandler](#) for Solr to use, you can extend the `ExtractingRequestHandler` and override the `createFactory()` method. This factory is responsible for constructing the [SolrContentHandler](#) that interacts with Tika, and allows literals to override Tika-parsed values. Set the parameter `literalsOverride`, which normally defaults to `*true`, to `*false` to append Tika-parsed values to literal values.

For more information on Solr's Extracting Request Handler, see
<https://wiki.apache.org/solr/ExtractingRequestHandler>

Topics covered in this section:

- [Key Concepts](#)
- [Trying out Tika with the Solr Example Directory](#)
- [Input Parameters](#)
- [Order of Operations](#)
- [Configuring the Solr ExtractingRequestHandler](#)
- [Indexing Encrypted Documents with the ExtractingUpdateRequestHandler](#)
- [Examples](#)
- [Sending Documents to Solr with a POST](#)
- [Sending Documents to Solr with Solr Cell and SolrJ](#)
- [Related Topics](#)

Key Concepts

When using the Solr Cell framework, it is helpful to keep the following in mind:

- Tika will automatically attempt to determine the input document type (Word, PDF, HTML) and extract the content appropriately. If you like, you can explicitly specify a MIME type for Tika with the `stream.type` parameter.

- Tika works by producing an XHTML stream that it feeds to a SAX ContentHandler. SAX is a common interface implemented for many different XML parsers. For more information, see <http://www.saxproject.org/quickstart.html>.
- Solr then responds to Tika's SAX events and creates the fields to index.
- Tika produces metadata such as Title, Subject, and Author according to specifications such as the DublinCore. See <http://tika.apache.org/1.0/formats.html> for the file types supported.
- Tika adds all the extracted text to the `content` field. This field is defined as "stored" in `schema.xml`. It is also copied to the `text` field with a `copyField` rule.
- You can map Tika's metadata fields to Solr fields. You can also boost these fields.
- You can pass in literals for field values. Literals will override Tika-parsed values, including fields in the Tika metadata object, the Tika content field, and any "captured content" fields.
- You can apply an XPath expression to the Tika XHTML to restrict the content that is produced.

💡 While Apache Tika is quite powerful, it is not perfect and fails on some files. PDF files are particularly problematic, mostly due to the PDF format itself. In case of a failure processing any file, the `ExtractingRequestHandler` does not have a secondary mechanism to try to extract some text from the file; it will throw an exception and fail.

Trying out Tika with the Solr Example Directory

You can try out the Tika framework using the example application included in Solr.

Start the Solr example server:

```
cd example -jar start.jar
```

In a separate window go to the `docs/` directory (which contains some nice example docs), or the `site` directory if you built Solr from source, and send Solr a file via HTTP POST:

```
curl 'http://localhost:8983/solr/update/extract?literal.id=doc1&commit=true'  
-F "myfile=@tutorial.html"
```

The URL above calls the Extraction Request Handler, uploads the file `tutorial.html` and assigns it the unique ID `doc1`. Here's a closer look at the components of this command:

- The `literal.id=doc1` parameter provides the necessary unique ID for the document being indexed.

- The `commit=true` parameter causes Solr to perform a commit after indexing the document, making it immediately searchable. For optimum performance when loading many documents, don't call the commit command until you are done.
- The `-F` flag instructs curl to POST data using the Content-Type `multipart/form-data` and supports the uploading of binary files. The `@` symbol instructs curl to upload the attached file.
- The argument `myfile=@tutorial.html` needs a valid path, which can be absolute or relative (for example, `myfile=@../../site/tutorial.html` if you are still in `exampledocs` directory).

Now you should be able to execute a query and find that document (open the following link in your browser): <http://localhost:8983/solr/select?q=tutorial>.

You may notice that although you can search on any of the text in the sample document, you may not be able to see that text when the document is retrieved. This is simply because the "content" field generated by Tika is mapped to the Solr field called `text`, which is indexed but not stored. This operation is controlled by default map rule in the `/update/extract` handler in `solrconfig.xml`, and its behavior can be easily changed or overridden. For example, to store and see all metadata and content, execute the following:

```
curl  
'http://localhost:8983/solr/update/extract?literal.id=doc1&uprefix=attr_&fmap.content=attr_content'  
-F "myfile=@tutorial.html"
```

In this command, the `uprefix=attr_` parameter causes all generated fields that aren't defined in the schema to be prefixed with `attr_`, which is a dynamic field that is stored.

The `fmap.content=attr_content` parameter overrides the default `fmap.content=text` causing the content to be added to the `attr_content` field instead.

Then run this command to query the document:

http://localhost:8983/solr/select?q=attr_content:tutorial

Input Parameters

The table below describes the parameters accepted by the Extraction Request Handler.

Parameter	Description
<code>boost.<fieldname></code>	Boosts the specified field by the defined float amount. (Boosting a field alters its importance in a query response. To learn about boosting fields, see Searching .)

capture	Captures XHTML elements with the specified name for a supplementary addition to the Solr document. This parameter can be useful for copying chunks of the XHTML into a separate field. For instance, it could be used to grab paragraphs (<code><p></code>) and index them into a separate field. Note that content is still also captured into the overall "content" field.
captureAttr	Indexes attributes of the Tika XHTML elements into separate fields, named after the element. If set to true, for example, when extracting from HTML, Tika can return the href attributes in <code><a></code> tags as fields named "a". See the examples below.
commitWithin	Add the document within the specified number of milliseconds.
date.formats	Defines the date format patterns to identify in the documents.
defaultField	If the uprefix parameter (see below) is not specified and a field cannot be determined, the default field will be used.
extractOnly	Default is false. If true, returns the extracted content from Tika without indexing the document. This literally includes the extracted XHTML as a string in the response. When viewing manually, it may be useful to use a response format other than XML to aid in viewing the embedded XHTML tags. For an example, see http://wiki.apache.org/solr/TikaExtractOnlyExampleOutput .
extractFormat	Default is "xml", but the other option is "text". Controls the serialization format of the extract content. The xml format is actually XHTML, the same format that results from passing the <code>-x</code> command to the Tika command line application, while the text format is like that produced by Tika's <code>-t</code> command. This parameter is valid only if <code>extractOnly</code> is set to true.
fmap.<source_field>	Maps (moves) one field name to another. The <code>source_field</code> must be a field in incoming documents, and the value is the Solr field to map to. Example: <code>fmap.content=text</code> causes the data in the <code>content</code> field generated by Tika to be moved to the Solr's <code>text</code> field.
literal.<fieldname>	Populates a field with the name supplied with the specified value for each document. The data can be multivalued if the field is multivalued.

literalsOverride	If true (the default), literal field values will override other values with the same field name. If false, literal values defined with literal.<fieldname> will be appended to data already in the fields extracted from Tika. If setting literalsOverride to "false", the field must be multivalued.
lowernames	Values are "true" or "false". If true, all field names will be mapped to lowercase with underscores, if needed. For example, "Content-Type" would be mapped to "content_type."
multipartUploadLimitInKB	Useful if uploading very large documents, this defines the KB size of documents to allow.
passwordsFile	Defines a file path and name for a file of file name to password mappings.
resource.name	Specifies the optional name of the file. Tika can use it as a hint for detecting a file's MIME type.
resource.password	Defines a password to use for a password-protected PDF or OOXML file
tika.config	Defines a file path and name to a customized Tika configuration file. This is only required if you have customized your Tika implementation.
uprefix	Prefixes all fields that are not defined in the schema with the given prefix. This is very useful when combined with dynamic field definitions. Example: uprefix=ignored_ would effectively ignore all unknown fields generated by Tika given the example schema contains <dynamicField name="ignored_*" type="ignored"/>
xpath	When extracting, only return Tika XHTML content that satisfies the given XPath expression. See http://tika.apache.org/1.0/index.html for details on the format of Tika XHTML. See also http://wiki.apache.org/solr/TikaExtractOnlyExampleOutput .

Order of Operations

Here is the order in which the Solr Cell framework, using the Extraction Request Handler and Tika, processes its input.

1. Tika generates fields or passes them in as literals specified by literal.<fieldname>=<value>. If literalsOverride=false, literals will be appended as multi-value to the Tika-generated field.
2. If lowernames=true, Tika maps fields to lowercase.

3. Tika applies the mapping rules specified by `fmap.source=target` parameters.
4. If `uprefix` is specified, any unknown field names are prefixed with that value, else if `defaultField` is specified, any unknown fields are copied to the default field.

Configuring the Solr ExtractingRequestHandler

If you are not working in the supplied `example/solr` directory, you must copy all libraries from `example/solr/libs` into a `libs` directory within your own `solr` directory or to a directory you've specified in `solrconfig.xml` using the new `libs` directive. The `ExtractingRequestHandler` is not incorporated into the Solr WAR file, so you have to install it separately.

Here is an example of configuring the `ExtractingRequestHandler` in `solrconfig.xml`.

```
<requestHandler name="/update/extract"
class="org.apache.solr.handler.extraction.ExtractingRequestHandler">
  <lst name="defaults">
    <str name="fmap.Last-Modified">last_modified</str>
    <str name="uprefix">ignored_</str>
  </lst>
  <!--Optional. Specify a path to a tika configuration file. See the Tika docs for
details.-->
  <str name="tika.config">/my/path/to/tika.config</str>
  <!-- Optional. Specify one or more date formats to parse. See
DateUtil.DEFAULT_DATE_FORMATS
    for default date formats -->
  <lst name="date.formats">
    <str>yyyy-MM-dd</str>
  </lst>
</requestHandler>
```

In the defaults section, we are mapping Tika's Last-Modified Metadata attribute to a field named `last_modified`. We are also telling it to ignore undeclared fields. These are all overridden parameters.

The `tika.config` entry points to a file containing a Tika configuration. The `date.formats` allows you to specify various `java.text.SimpleDateFormat` date formats for working with transforming extracted input to a Date. Solr comes configured with the following date formats (see the `DateUtil` in Solr):

```
YYYY-MM-dd 'T' HH:mm:ss 'Z'  
YYYY-MM-dd 'T' HH:mm:ss  
YYYY-MM-dd  
YYYY-MM-dd hh:mm:ss  
YYYY-MM-dd HH:mm:ss  
EEE MMM d hh:mm:ss z YYYY  
EEE, dd MMM yyyy HH:mm:ss zzz  
EEEE, dd-MMM-yy HH:mm:ss zzz  
EEE MMM d HH:mm:ss YYYY
```

You may also need to adjust the `multipartUploadLimitInKB` attribute as follows if you are submitting very large documents.

```
<requestDispatcher handleSelect="true" >  
  <requestParsers enableRemoteStreaming="false" multipartUploadLimitInKB="20480" />  
  ...
```

Multi-Core Configuration

For a multi-core configuration, specify `sharedLib='lib'` in the `<solr/>` section of `solr.xml` in order for Solr to find the JAR files in `example/solr/lib`.

For more information about Solr cores, see [The Well-Configured Solr Instance](#).

Indexing Encrypted Documents with the ExtractingUpdateRequestHandler

The `ExtractingRequestHandler` will decrypt encrypted files and index their content if you supply a password in either `resource.password` on the request, or in a `passwordsFile` file.

In the case of `passwordsFile`, the file supplied must be formatted so there is one line per rule. Each rule contains a file name regular expression, followed by "=", then the password in clear-text. Because the passwords are in clear-text, the file should have strict access restrictions.

```
# This is a comment  
myFileName = myPassword  
.*\.\.docx\$ = myWordPassword  
.*\.\.pdf\$ = myPdfPassword
```

Examples

Metadata

As mentioned before, Tika produces metadata about the document. Metadata describes different aspects of a document, such as the author's name, the number of pages, the file size, and so on. The metadata produced depends on the type of document submitted. For instance, PDFs have different metadata than Word documents do.

In addition to Tika's metadata, Solr adds the following metadata (defined in `ExtractingMetadataConstants`):

Solr Metadata	Description
<code>stream_name</code>	The name of the Content Stream as uploaded to Solr. Depending on how the file is uploaded, this may or may not be set
<code>stream_source_info</code>	Any source info about the stream. (See the section on Content Streams later in this section.)
<code>stream_size</code>	The size of the stream in bytes.
<code>stream_content_type</code>	The content type of the stream, if available.



We recommend that you try using the `extractOnly` option to discover which values Solr is setting for these metadata elements.

Examples of Uploads Using the Extraction Request Handler

Capture and Mapping

The command below captures `<div>` tags separately, and then maps all the instances of that field to a dynamic field named `foo_t`.

```
curl "http://localhost:8983/solr/update/extract?literal.id=doc2&captureAttr=true&defaultField=text&fmap.div=foo_t&capture=div" -F "tutorial=@tutorial.pdf"
```

Capture, Mapping, and Boosting

The command below captures `<div>` tags separately, maps the field to a dynamic field named `foo_t`, then boosts `foo_t` by 3.

```
curl "http://localhost:8983/solr/update/extract?literal.id=doc3&captureAttr=true&defaultField=text&capture=div&fmap.div=foo_t&boost.foo_t=3" -F "tutorial=@tutorial.pdf"
```

Using Literals to Define Your Own Metadata

To add in your own metadata, pass in the literal parameter along with the file:

```
curl "http://localhost:8983/solr/update/extract?literal.id=doc4&captureAttr=true&defaultField=text&capture=div&fmap.div=foo_t&boost.foo_t=3&literal.blah_s=Bah"-F "tutorial=@tutorial.pdf"
```

XPath

The example below passes in an XPath expression to restrict the XHTML returned by Tika:

```
curl "http://localhost:8983/solr/update/extract?literal.id=doc5&captureAttr=true&defaultField=text&capture=div&fmap.div=foo_t&boost.foo_t=3&literal.id=id&xpath=/xhtml:html/xhtml:body/xhtml:div/ancestor::node()"-F "tutorial=@tutorial.pdf"
```

Extracting Data without Indexing It

Solr allows you to extract data without indexing. You might want to do this if you're using Solr solely as an extraction server or if you're interested in testing Solr extraction.

The example below sets the `extractOnly=true` parameter to extract data without indexing it.

```
curl "http://localhost:8983/solr/update/extract?&extractOnly=true"--data-binary @tutorial.html-H 'Content-type:text/html'
```

The output includes XML generated by Tika (and further escaped by Solr's XML) using a different output format to make it more readable:

```
curl "http://localhost:8983/solr/update/extract?&extractOnly=true&wt=ruby&indent=true"--data-binary @tutorial.html-H 'Content-type:text/html'
```

Sending Documents to Solr with a POST

The example below streams the file as the body of the POST, which does not, then, provide information to Solr about the name of the file.

```
curl "http://localhost:8983/solr/update/extract?literal.id=doc5&defaultField=text"
--data-binary @tutorial.html
-H 'Content-type:text/html'
```

Sending Documents to Solr with Solr Cell and SolrJ

SolrJ is a Java client that you can use to add documents to the index, update the index, or query the index. You'll find more information on SolrJ in [Client APIs](#).

Here's an example of using Solr Cell and SolrJ to add documents to a Solr index.

First, let's use SolrJ to create a new SolrServer, then we'll construct a request containing a ContentStream (essentially a wrapper around a file) and sent it to Solr:

```
public class SolrCellRequestDemo {
    public static void main (String[] args){color} throws IOException,
    SolrServerException {
        SolrServer server = new HttpSolrServer("http://localhost:8983/solr");
        ContentStreamUpdateRequest req = new ContentStreamUpdateRequest("/update/extract");
        req.addFile(new File("apache-solr/site/features.pdf"));
        req.setParam(ExtractingParams.EXTRACT_ONLY, "true");
        NamedList<Object> result = server.request(req);
        System.out.println("Result: " + result);
    }
}
```

This operation streams the file features.pdf into the Solr index.

The sample code above calls the extract command, but you can easily substitute other commands that are supported by Solr Cell. The key class to use is the ContentStreamUpdateRequest, which makes sure the ContentStreams are set properly. SolrJ takes care of the rest.

Note that the ContentStreamUpdateRequest is not just specific to Solr Cell. You can send CSV to the CSV Update handler and to any other Request Handler that works with Content Streams for updates.

Related Topics

- [ExtractingRequestHandler](#)



Uploading Data with Index Handlers

Index Handlers are Update Handlers designed to add, delete and update documents to the index. Solr includes several of these to allow indexing documents in XML, CSV and JSON.

The example URLs given here reflect the handler configuration in the supplied `solrconfig.xml`. If the name associated with the handler is changed then the URLs will need to be modified. It is quite possible to access the same handler using more than one name, which can be useful if you wish to specify different sets of default options.

New `UpdateProcessors` now default to the `uniqueKey` field if it is the appropriate type for configured fields. The processors automatically add fields with new UUIDs and Timestamps to `SolrInputDocuments`. These work similarly to the `<field default="..." />` option in `schema.xml`, but are applied in the `UpdateProcessorChain`. They may be used prior to other `UpdateProcessors`, or to generate a `uniqueKey` field value when using the `DistributedUpdateProcessor` (i.e., `SolrCloud`), `TimestampUpdateProcessorFactory`, `UUIDUpdateProcessorFactory`, and `DefaultValueUpdateProcessorFactory`.

Index Handlers covered in this section:

- [Combined UpdateRequestHandlers](#)
- [XMLUpdateRequestHandler for XML-formatted Data](#)
- [XSLTRequestHandler to Transform XML Content](#)
- [CSVRequestHandler for CSV Content](#)
- [Using the JsonRequestHandler for JSON Content](#)
- [Updating Only Part of a Document](#)
- [Using SimplePostTool](#)
- [Indexing Using SolrJ](#)

Combined UpdateRequestHandlers

For the separate XML, CSV, JSON, and javabin update request handlers explained below, Solr provides a single `RequestHandler`, and chooses the appropriate `ContentStreamLoader` based on the the `Content-Type` header, entered as the `qt` (query type) parameter matching the name of registered handlers. The "standard" request handler is the default and will be used if `qt` is not specified in the request.

```
<requestHandler name="standard" />
<requestHandler name="custom" />
```

Configuring Shard Handlers for Distributed Searches

Inside the RequestHandler, you can configure and specify the shard handler used for distributed search. You can also plug in custom shard handlers as well.

Configuring the standard handler, set up the configuration as in this example:

```
<requestHandler name="standard" default="true">
    <!-- other params go here -->
    <shardHandlerFactory>
        <int name="socketTimeOut">1000</int>
        <int name="connTimeOut">5000</int>
    </shardHandler>
</requestHandler>
```

The parameters that can be specified are as follows:

Parameter	Default	Explanation
socketTimeout	default: 0 (use OS default)	The amount of time in ms that a socket is allowed to wait
connTimeout	default: 0 (use OS default)	The amount of time in ms that is accepted for binding / connection a socket
maxConnectionsPerHost	default: 20	The maximum number of connections that is made to each individual shard in a distributed search
corePoolSize	default: 0	The retained lowest limit on the number of threads used in coordinating distributed search
maximumPoolSize	default: Integer.MAX_VALUE	The maximum number of threads used for coordinating distributed search
maxThreadIdleTime	default: 5 seconds	The amount of time to wait for before threads are scaled back in response to a reduction in load
sizeOfQueue	default: -1	If specified, the thread pool will use a backing queue instead of a direct handoff buffer. High throughput systems will want to configure this to be a direct hand off (with -1). Systems that desire better latency will want to configure a reasonable size of queue to handle variations in requests.

fairnessPolicy	default: false	Chooses the JVM specifics dealing with fair policy queuing. If enabled, distributed searches will be handled in a First in - First out method at a cost to throughput. If disabled, throughput will be favored over latency.
----------------	----------------	--

XMLUpdateRequestHandler for XML-formatted Data

Configuration

The default configuration file has the update request handler configured by default.

```
<requestHandler name="/update" class="solr.XmlUpdateRequestHandler" />
```

Adding Documents

Documents are added to the index by sending an XML message to the update handler.

The XML schema recognized by the update handler is very straightforward:

- The `<add>` element introduces one or more documents to be added.
- The `<doc>` element introduces the fields making up a document.
- The `<field>` element presents the content for a specific field.

For example:

```

<add>
  <doc>
    <field name="authors">Patrick Eagar</field>
    <field name="subject">Sports</field>
    <field name="dd">796.35</field>
    <field name="numpages">128</field>
    <field name="desc"></field>
    <field name="price">12.40</field>
    <field name="title" boost="2.0">Summer of the all-rounder: Test and championship
cricket in England 1982</field>
    <field name="isbn">0002166313</field>
    <field name="yearpub">1982</field>
    <field name="publisher">Collins</field>
  </doc>
  <doc boost="2.5">
    ...
  </doc>
</add>

```

If the document schema defines a unique key, then an `/update` operation silently replaces a document in the index with the same unique key, unless the `<add>` element sets the `allowDups` attribute to `true`. If no unique key has been defined, indexing performance is somewhat faster, as no search has to be made for an existing document.

Each element has certain optional attributes which may be specified.

Command	Command Description	Optional Parameter	Parameter Description
<code><add></code>	Introduces one or more documents to be added to the index.	<code>commitWithin=number</code>	Add the document within the specified number of milliseconds
<code><doc></code>	Introduces the definition of a specific document.	<code>boost=float</code>	Default is 1.0. Sets a boost value for the document. To learn more about boosting, see Searching .
<code><field></code>	Defines a field within a document.	<code>boost=float</code>	Default is 1.0. Sets a boost value for the field.

⚠ Other optional parameters for `<add>`, including `allowDups`, `overwritePending`, and `overwriteCommitted`, are now deprecated. However, you can specify `overwrite=false` for XML updates to avoid overwriting.

Commit and Optimize Operations

The `<commit>` operation writes all documents loaded since the last commit to one or more segment files on the disk. Before a commit has been issued, newly indexed content is not visible to searches. The commit operation opens a new searcher, and triggers any event listeners that have been configured.

Commits may be issued explicitly with a `<commit/>` message, and can also be triggered from `<autocommit>` parameters in `solrconfig.xml`.

The `<optimize>` operation requests Solr to merge internal data structures in order to improve search performance. For a large index, optimization will take some time to complete, but by merging many small segment files into a larger one, search performance will improve. If you are using Solr's replication mechanism to distribute searches across many systems, be aware that after an optimize, a complete index will need to be transferred. In contrast, post-commit transfers are usually much smaller.

The `<commit>` and `<optimize>` elements accept these optional attributes:

Optional Attribute	Description
maxSegments	Default is 1. Optimizes the index to include no more than this number of segments.
waitFlush	Default is true. Blocks until index changes are flushed to disk.
waitSearcher	Default is true. Blocks until a new searcher is opened and registered as the main query searcher, making the changes visible.
expungeDeletes	Default is false. Merges segments and removes deleted documents.

Here are examples of `<commit>` and `<optimize>` using optional attributes:

```
<commit waitFlush="false" waitSearcher="false"/>
<commit waitFlush="false" waitSearcher="false" expungeDeletes="true"/>
<optimize waitFlush="false" waitSearcher="false"/>
```

Delete Operations

Documents can be deleted from the index in two ways. "Delete by ID" deletes the document with the specified ID, and can be used only if a UniqueID field has been defined in the schema. "Delete by Query" deletes all documents matching a specified query, although `commitWithin` is ignored for a Delete by Query. A single delete message can contain multiple delete operations.

```
<delete>
<id>0002166313</id>
<id>0031745983</id>
<query>subject:sport</query>
<query>publisher:penguin</query>
</delete>
```

Rollback Operations

The rollback command rolls back all add and deletes made to the index since the last commit. It neither calls any event listeners nor creates a new searcher. Its syntax is simple: <rollback/>.

Using curl to Perform Updates with the Update Request Handler.

You can use the `curl` utility to perform any of the above commands, using its `--data-binary` option to append the XML message to the `curl` command, and generating a HTTP POST request. For example:

```
curl http://localhost:8983/update -H "Content-Type: text/xml" --data-binary '
<add>
<doc>
<field name="authors">Patrick Eagar</field>
<field name="subject">Sports</field>
<field name="dd">796.35</field>
<field name="isbn">0002166313</field>
<field name="yearpub">1982</field>
<field name="publisher">Collins</field>
</doc>
</add>'
```

For posting XML messages contained in a file, you can use the alternative form:

```
curl http://localhost:8983/update -H "Content-Type: text/xml"
--data-binary @myfile.xml
```

Short requests can also be sent using a HTTP GET command, URL-encoding the request, as in the following. Note the escaping of "<" and ">":

```
curl http://localhost:8983/update?stream.body=%3Ccommit%3E
```

Responses from Solr take the form shown here:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<lst name="responseHeader">
<int name="status">0</int>
<int name="QTime">127</int>
</lst>
</response>
```

The status field will be non-zero in case of failure. The servlet container will generate an appropriate HTML-formatted message in the case of an error at the HTTP layer.

A Simple Cross-Platform Posting Tool

For demo purposes, the file `$SOLR/example/exampledocs/post.jar` includes a cross-platform Java tool for POST-ing XML documents. Open a window and run:

```
java -jar post.jar <list of files with messages>
```

By default, this will contact the server at `localhost:8983`. The "-help" option outputs the following information on its usage:

```
SimplePostTool: version 1.2
```

This is a simple command line tool for POSTing raw XML to a Solr port. XML data can be read from files specified as command line args; as raw commandline arg strings; or via STDIN.

Examples:

```
java -Ddata=files -jar post.jar *.xml
java -Ddata=args -jar post.jar '<delete><id>42</id></delete>'
java -Ddata=stdin -jar post.jar < hd.xml
```

Other options controlled by System Properties include the Solr URL to POST to, and whether a commit should be executed. These are the defaults for all System Properties.

```
-Ddata=files
-Durl=[http://localhost:8983/solr/update|http://localhost:8983/solr/update]
-Dcommit=yes
```

For more information about the XML Update Request Handler, see
<https://wiki.apache.org/solr/UpdateXmlMessages>.

XSLTRequestHandler to Transform XML Content

Configuration

The default configuration file has the update request handler configured by default, although the "lazy load" flag is set.

The XSLTRequestHandler allows you to index any XML data with the [XML <tr> command](#). You must have an XSLT stylesheet in the solr/conf/xslt directory that can transform the incoming data to the expected <add><doc/></add> format.

```
<requestHandler name="/update/xslt" startup="lazy"
class="solr.XsltUpdateRequestHandler"/>
```

Here is an example XSLT stylesheet:

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
<xsl:template match="/">
<add>
<xsl:apply-templates select="/random/document"/>
</add>
</xsl:template>

<xsl:template match="document">

<doc boost="5.5">
<xsl:apply-templates select="*"/>
</doc>
</xsl:template>

<xsl:template match="node">
<field name="{@name}">
<xsl:if test="@enhance != ''">
<xsl:attribute name="boost"><xsl:value-of select="@enhance"/></xsl:attribute>
</xsl:if>
<xsl:value-of select="@value"/>
</field>
</xsl:template>

</xsl:stylesheet>
```

Attaching the stylesheet "updateXml.xsl" transforms a search result to Solr's `UpdateXml` syntax. One example is to copy a Solr1.3 index (which does not have CSV response writer) into a format which can be indexed into another Solr file (provided that all fields are stored):

```
http://localhost:8983/solr/select?q=*:*&wt=xslt&tr=updateXml.xsl&rows=1000
```

You can also use the stylesheet in `xsltUpdateRequestHandler` to transform an index when updating:

```
curl "http://localhost:8983/solr/update/xslt?commit=true&tr=updateX
```

CSVRequestHandler for CSV Content

Configuration

The default configuration file has the update request handler configured by default, although the "lazy load" flag is set.

```
<requestHandler name="/update/csv" class="solr.CSVRequestHandler" startup="lazy" />
```

Parameters

The CSV handler allows the specification of many parameters in the URL in the form: `f.parameter.optional_fieldname=value`.

The table below describes the parameters for the update handler.

Parameter	Usage	Global (g) or Per Field (f)	Example
separator	Character used as field separator; default is ","	g,(f: see split)	separator=%
trim	If true, remove leading and trailing whitespace from values. Default=false.	g,f	f.isbn.trim=true trim=false
header	Set to true if first line of input contains field names. These will be used if the field_name parameter is absent.	g	

field_name	Comma separated list of field names to use when adding documents.	g	field_name=isbn,price,title
literal.<field_name>	Comma separated list of field names to use when processing literal values.	g	literal.color=red,blue,black
skip	Comma separated list of field names to skip.	g	skip=uninteresting,shoesize
skipLines	Number of lines to discard in the input stream before the CSV data starts, including the header, if present. Default=0.	g	skipLines=5
encapsulator	The character optionally used to surround values to preserve characters such as the CSV separator or whitespace. This standard CSV format handles the encapsulator itself appearing in an encapsulated value by doubling the encapsulator.	g,(f: see split)	encapsulator="
escape	The character used for escaping CSV separators or other reserved characters. If an escape is specified, the encapsulator is not used unless also explicitly specified since most formats use either encapsulation or escaping, not both	g	escape=\
keepEmpty	Keep and index zero length (empty) fields. Default=false.	g,f	f.price.keepEmpty=true
map	Map one value to another. Format is value:replacement (which can be empty.)	g,f	map=left:right f.subject.map=history:bunk
split	If true, split a field into multiple values by a separate parser.	f	

overwrite	If true (the default), check for and overwrite duplicate documents, based on the uniqueKey field declared in the Solr schema. If you know the documents you are indexing do not contain any duplicates then you may see a considerable speed up setting this to false.	g	
commit	Issues a commit after the data has been ingested.	g	
commitWithin	Add the document within the specified number of milliseconds.	g	commitWithin=10000

For more information on the CSV Update Request Handler, see <https://wiki.apache.org/solr/UpdateCSV>.

Using the **JSONRequestHandler** for JSON Content

JSON formatted update requests may be sent to Solr using the /solr/update/json URL. All of the normal methods for uploading content are supported.

Configuration

The default configuration file has the update request handler configured by default, although the "lazy load" flag is set.

```
<requestHandler name="/update/json" class="solr.JsonUpdateRequestHandler"
  startup="lazy" />
```

Examples

There is a sample JSON file at example/exampledocs/books.json that you can use to add documents to the Solr example server.

```
cd example/exampledocs
curl 'http://localhost:8983/solr/update/json?commit=true'
--data-binary @books.json -H 'Content-type:application/json'
```

Adding commit=true to the URL makes the documents immediately searchable.

You should now be able to query for the newly added documents:

`http://localhost:8983/solr/select?q=title:monsters&wt=json&indent=true` returns:

```
{  
  "responseHeader":{  
    "status":0,  
    "QTime":2,  
    "params":{  
      "indent":"true",  
      "wt":"json",  
      "q":"title:monsters"}},  
  "response":{ "numFound":1, "start":0, "docs": [  
    {  
      "id": "978-1423103349",  
      "author": "Rick Riordan",  
      "series_t": "Percy Jackson and the Olympians",  
      "sequence_i": 2,  
      "genre_s": "fantasy",  
      "inStock": true,  
      "price": 6.49,  
      "pages_i": 304,  
      "title": [  
        "The Sea of Monsters"],  
      "cat": ["book", "paperback"] } ]  
  }  
}
```

Update Commands

The JSON update handler accepts all of the update commands that the XML update handler supports, through a straightforward mapping. Multiple commands may be contained in one message:

```
{  
  "add": {  
    "doc": {  
      "id": "DOC1",  
      "my_boosted_field": { /* use a map with boost/value for a boosted field */  
        "boost": 2.3,  
        "value": "test"  
      },  
      "my_multivalued_field": [ "aaa", "bbb" ] /* use an array for a multi-valued field */  
    }/  
  }/  
},  
"add": {  
  "commitWithin": 5000, /* commit this document within 5 seconds */  
  "overwrite": false, /* don't check for existing documents with the same  
uniqueKey */  
  "boost": 3.45, /* a document boost */  
  "doc": {  
    "f1": "v1",  
    "f1": "v2"  
  }/  
},  
  
"commit": {},  
"optimize": { "waitFlush":false, "waitSearcher":false },  
  
"delete": { "id":"ID" }, /* delete by ID */  
"delete": { "query":"QUERY" } /* delete by query */  
}
```

 Comments are not allowed JSON, but duplicate names are.

As with other update handlers, parameters such as `commit`, `commitWithin`, `optimize`, and `overwrite` may be specified in the URL instead of in the body of the message.

The JSON update format allows for a simple delete-by-id. The value of a `delete` can be an array which contains a list of zero or more specific document id's (not a range) to be deleted. For example:

```
"delete":"myid"
```

```
"delete": [ "id1", "id2" ]
```

The value of a "delete" can be an array which contains a list of zero or more id's to be deleted. It is not a range (start and end).

You can also specify *version* with each "delete":

```
String str = "
```

```
"
```

You can specify the version of deletes in the body of the update request as well.

For more information about the JSON Update Request Handler, see

<https://wiki.apache.org/solr/UpdateJSON>.

Updating Only Part of a Document

Solr supports several modifiers that atomically update values of a document.

Modifier	Usage
set	Set or replace a particular value, or remove the value if null is specified as the new value.
add	Adds an additional value to a list.
inc	Increments a numeric value by a specific amount.



All original source fields must be stored for field modifiers to work correctly, which is the Solr default.

For example:

```
{"id": "mydoc", "f1": {"set": 10}, "f2": {"add": 20}}
```

This example results in field f1 being set to "10", and field f2 having an additional value of "20" added. All other existing fields from the original document remain unchanged.

Using SimplePostTool

This is a simple command line tool for POSTing raw data to a Solr port. Data can be read from files specified as command line arguments, as raw command line argument strings, or via STDIN. Options controlled by System Properties include the Solr URL to post to, the Content-Type of the data, whether a commit or optimize should be executed, and whether the response should be written to STDOUT. If auto=yes the tool will try to guess the type and set Content-Type and the URL automatically. When posting rich documents the file name will be propagated as resource.name and also used as literal.id. You may override these or any other request parameter through the -Dparams property

Supported System Properties and their defaults:

Parameter	Values	Default
-Ddata	yes, no	default=files
-Dtype	<content-type>	default=application/xml
-Durl	<solr-update-url>	default= http://localhost:8983/solr/update
-Dauto	yes, no	default=no
-Drecursive	yes, no	default=no
-Dfiletypes	<type>[,<type>,...]	default=xml, json, csv, pdf, doc, docx, ppt, pptx, xls, xlsx, odt, odp, ods, rtf, htm, html
-Dparams	"<key>=<value>[&<key>=<value>...]"	values must be URL-encoded
-Dcommit	yes, no	default=yes
-Doptimize	yes, no	default=no
-Dout	yes,no	default=no

Examples:

```
java -jar post.jar *.xml
java -Ddata=args -jar post.jar '<delete><id>42</id></delete>'
java -Ddata=stdin -jar post.jar < hd.xml
java -Dtype=text/csv -jar post.jar *.csv
java -Dtype=application/json -jar post.jar *.json
java -Durl=[http://localhost:8983/solr/update/extract] -Dparams=literal.id=a
-Dtype=application/pdf -jar post.jar a.pdf
java -Dauto=yes -jar post.jar a.pdf
java -Dauto=yes -Drecursive=yes -jar post.jar afolder
java -Dauto=yes -Dfiletypes=ppt,html -jar post.jar afolder
```

In the above example:

-Dauto=yes	Will guess file type from file name suffix, and set type and url accordingly. It also sets the ID and file name automatically.
-Drecursive=yes	Will recurse into sub-folders and index all files.
-Dfiletypes	Specifies the file types to consider when indexing folders.
-Dparams	HTTP GET params to add to the request, so you don't need to write the whole URL again.

Indexing Using SolrJ

Use of the the SolrJ client library is covered in the section on [Using SolrJ](#).

Uploading Structured Data Store Data with the Data Import Handler

Many search applications store the content to be indexed in a structured data store, such as a relational database. The Data Import Handler (DIH) provides a mechanism for importing content from a data store and indexing it. In addition to relational databases, DIH can index content from HTTP based data sources such as RSS and ATOM feeds, e-mail repositories, and structured XML where an XPath processor is used to generate fields.



The DataImportHandler jars are no longer included in the Solr WAR. You should add them to Solr's lib directory, or reference them via the <lib> directive in `solrconfig.xml`.

For more information about the Data Import Handler, see <https://wiki.apache.org/solr/DataImportHandler>.

Topics covered in this section:

- Concepts and Terminology
- Configuration
- Data Import Handler Commands
- Property Writer
- Data Sources
- Entity Processors
- Transformers
- Special Commands for the Data Import Handler

Concepts and Terminology

Descriptions of the Data Import Handler use several familiar terms, such as entity and processor, in specific ways, as explained in the table below.

Term	Definition
Datasource	As its name suggests, a datasource defines the location of the data of interest. For a database, it's a DSN. For an HTTP datasource, it's the base URL.

Entity	Conceptually, an entity is processed to generate a set of documents, containing multiple fields, which (after optionally being transformed in various ways) are sent to Solr for indexing. For a RDBMS data source, an entity is a view or table, which would be processed by one or more SQL statements to generate a set of rows (documents) with one or more columns (fields).
Processor	An entity processor does the work of extracting content from a data source, transforming it, and adding it to the index. Custom entity processors can be written to extend or replace the ones supplied.
Transformer	Each set of fields fetched by the entity may optionally be transformed. This process can modify the fields, create new fields, or generate multiple rows/documents from a single row. There are several built-in transformers in the DIH, which perform functions such as modifying dates and stripping HTML. It is possible to write custom transformers using the publicly available interface.

Configuration

Configuring solrconfig.xml

The Data Import Handler has to be registered in `solrconfig.xml`. For example:

```
<requestHandler name="/dataimport"
class="org.apache.solr.handler.dataimport.DataImportHandler">
<lst name="defaults">
<str name="config"/>/path/to/my/DIHconfigfile.xml</str>
</lst>
</requestHandler>
```

The only required parameter is the `config` parameter, which specifies the location of the DIH configuration file that contains specifications for the data source, how to fetch data, what data to fetch, and how to process it to generate the Solr documents to be posted to the index.

You can have multiple DIH configuration files. Each file would require a separate definition in the `solrconfig.xml` file, specifying a path to the file.

Configuring the DIH Config File

There is a sample DIH application distributed with Solr in the directory `example/example-DIH`. This accesses a small hsqldb database. Details of how to run this example can be found in the `README.txt` file. The sample DIH configuration can be found in `example/example-DIH/solr/db/conf/db-data-config.xml`.

An annotated configuration file, based on the sample, is shown below. It extracts fields from the four tables defining a simple product database, with this schema. More information about the parameters and options shown here are described in the sections following.

```
<dataConfig>

<!-- The first element is the dataSource, in this case an HSQLDB database.
     The path to the JDBC driver and the JDBC URL and login credentials are all
     specified here.

     Other permissible attributes include whether or not to autocommit to Solr, the
     batchSize

     used in the JDBC connection, a 'readOnly' flag -->

    <dataSource driver="org.hsqldb.jdbcDriver"
url="jdbc:hsqldb:./example-DIH/hsqldb/ex" user="sa" />

<!-- a 'document' element follows, containing multiple 'entity' elements.
     Note that 'entity' elements can be nested, and this allows the entity
     relationships in the sample database to be mirrored here, so that we can
     generate a denormalized Solr record which may include multiple features
     for one item, for instance -->

<document>

<!--The possible attributes for the entity element are described below.
     Entity elements may contain one or more 'field' elements, which map
     the data source field names to Solr fields, and optionally specify
     per-field transformations -->

<!-- this entity is the 'root' entity. -->

    <entity name="item" query="select * from item"
           deltaQuery="select id from item where last_modified >
 '${dataimporter.last_index_time}'>
        <field column="NAME" name="name" />

    <!-- This entity is nested and reflects the one-to-many relationship between an item
        and its multiple features.

        Note the use of variables; ${item.ID} is the value of the column 'ID' for the
        current item

        ('item' referring to the entity name) -->

        <entity name="feature"
               query="select DESCRIPTION from FEATURE where ITEM_ID='${item.ID}'"
               deltaQuery="select ITEM_ID from FEATURE where last_modified >
 '${dataimporter.last_index_time}'"
```

```
        parentDeltaQuery="select ID from item where ID=${feature.ITEM_ID}">
            <field name="features" column="DESCRIPTION" />
        </entity>

        <entity name="item_category"
            query="select CATEGORY_ID from item_category where
ITEM_ID='${item.ID}'"
            deltaQuery="select ITEM_ID, CATEGORY_ID from item_category where
last_modified > '${dataimporter.last_index_time}'"
            parentDeltaQuery="select ID from item where
ID=${item_category.ITEM_ID}">
            <entity name="category"
                query="select DESCRIPTION from category where ID =
`${item_category.CATEGORY_ID}'"
                deltaQuery="select ID from category where last_modified >
`${dataimporter.last_index_time}'"
                parentDeltaQuery="select ITEM_ID, CATEGORY_ID from
item_category where CATEGORY_ID=${category.ID}">
                <field column="description" name="cat" />
            </entity>
        </entity>
    </entity>
</document>
</dataConfig>
```

Datasources can still be specified in `solrconfig.xml`. These must be specified in the defaults section of the handler in `solrconfig.xml`. However, these are not parsed until the main configuration is loaded.

The entire configuration itself can be passed as a request parameter using the `dataConfig` parameter rather than using a file. When configuration errors are encountered, the error message is returned in XML format.

In Solr 4.1, a new property has been added, the `propertyWriter` element, which allows defining the date format and locale for use with delta queries. It also allows customizing the name and location of the properties file.

The `reload-config` command is still supported, which is useful for validating a new configuration file, or if you want to specify a file, load it, and not have it reloaded again on import. If there is an `xml` mistake in the configuration a user-friendly message is returned in `xml` format. You can then fix the problem and do a `reload-config`.

- ✓ You can also view the DIH configuration in the Solr Admin UI. There is also an interface to import content.

Data Import Handler Commands

DIH commands are sent to Solr via an HTTP request. The following operations are supported.

Command	Description
abort	Aborts an ongoing operation. The URL is <code>http://<host>:<port>/solr/dataimport?command=abort</code> .
delta-import	For incremental imports and change detection. The command is of the form <code>http://<host>:<port>/solr/dataimport?command=delta-import</code> . It supports the same clean, commit, optimize and debug parameters as full-import command.
full-import	A Full Import operation can be started with a URL of the form <code>http://<host>:<port>/solr/dataimport?command=full-import</code> . The command returns immediately. The operation will be started in a new thread and the <i>status</i> attribute in the response should be shown as <i>busy</i> . The operation may take some time depending on the size of dataset. Queries to Solr are not blocked during full-imports. When a full-import command is executed, it stores the start time of the operation in a file located at <code>conf/dataimport.properties</code> . This stored timestamp is used when a delta-import operation is executed. For a list of parameters that can be passed to this command, see below.
reload-config	If the configuration file has been changed and you wish to reload it without restarting Solr, run the command <code>http://<host>:<port>/solr/dataimport?command=reload-config</code> .
status	The URL is <code>http://<host>:<port>/solr/dataimport?command=status</code> . It returns statistics on the number of documents created, deleted, queries run, rows fetched, status, and so on.

Parameters for the full-import Command

The `full-import` command accepts the following parameters:

Parameter	Description
clean	Default is true. Tells whether to clean up the index before the indexing is started.
commit	Default is true. Tells whether to commit after the operation.

debug	Default is false. Runs the command in debug mode. It is used by the interactive development mode. Note that in debug mode, documents are never committed automatically. If you want to run debug mode and commit the results too, add <code>commit=true</code> as a request parameter.
entity	The name of an entity directly under the <code><document></code> tag in the configuration file. Use this to execute one or more entities selectively. Multiple "entity" parameters can be passed on to run multiple entities at once. If nothing is passed, all entities are executed.
optimize	Default is true. Tells Solr whether to optimize after the operation.

Property Writer

The `propertyWriter` element defines the date format and locale for use with delta queries. It is an optional configuration. Add the element to the DIH configuration file, directly under the `dataConfig` element.

```
<propertyWriter dateFormat="yyyy-MM-dd HH:mm:ss" type="SimplePropertiesWriter"
directory="data" filename="my_dih.properties" locale="en_US" />
```

The parameters available are:

Parameter	Description
dateFormat	A <code>java.text.SimpleDateFormat</code> to use when converting the date to text. The default is "yyyy-MM-dd HH:mm:ss".
type	The implementation class. Use <code>SimplePropertiesWriter</code> for non-SolrCloud installations. If using SolrCloud, use <code>ZKPropertiesWriter</code> . If this is not specified, it will default to the appropriate class depending on if SolrCloud mode is enabled.
directory	Used with the <code>SimplePropertiesWriter</code> only). The directory for the properties file. If not specified, the default is "conf".
filename	Used with the <code>SimplePropertiesWriter</code> only). The name of the properties file. If not specified, the default is the <code>requestHandler</code> name (as defined in <code>solrconfig.xml</code> , appended by ".properties" (i.e., "dataimport.properties").
locale	The locale. If not defined, the ROOT locale is used. It must be specified as language-country. For example, en-US.

Data Sources

A data source specifies the origin of data and its type. Somewhat confusingly, some data sources are configured within the associated entity processor. Data sources can also be specified in `solrconfig.xml`, which is useful when you have multiple environments (for example, development, QA, and production) differing only in their data sources.

You can create a custom data source by writing a class that extends `org.apache.solr.handler.dataimport.DataSource`.

The mandatory attributes for a data source definition are its name and type. The name identifies the data source to an Entity element.

The types of data sources available are described below.

ContentStreamDataSource

This takes the POST data as the data source. This can be used with any EntityProcessor that uses a `DataSource<Reader>`.

FieldReaderDataSource

This can be used where a database field contains XML which you wish to process using the `XpathEntityProcessor`. You would set up a configuration with both JDBC and FieldReader data sources, and two entities, as follows:

```
<dataSource name="a1" driver="org.hsqldb.jdbcDriver" ... />
<dataSource name="a2" type=FieldReaderDataSource" />

<!-- processor for database -->

<entity name ="e1" dataSource="a1" processor="SQLEntityProcessor" pk="docid"
query="select * from t1 ..." >

<!-- nested XpathEntity; the field in the parent which is to be used for
Xpath is set in the 'datafield' attribute inplace of the "url" attribute -->

<entity name="e2"
dataSource="a2"
processor="XPathEntityProcessor"
dataField="e1.fieldToUseForXPath"

<!-- Xpath configuration follows -->
...
</entity>
</entity>
```

The FieldReaderDataSource can take an `encoding` parameter, which will default to "UTF-8" if not specified. It must be specified as language-country. For example, en-US.

FileDataSource

This can be used like an [URLDataSource](#), but is used to fetch content from files on disk. The only difference from URLDataSource, when accessing disk files, is how a pathname is specified.

This data source accepts these optional attributes.

Optional Attribute	Description
basePath	The base path relative to which the value is evaluated if it is not absolute.
encoding	Defines the character encoding to use. If not defined, UTF-8 is used.

JdbcDataSource

This is the default datasource. It's used with the [SQLEntityProcessor](#). See the example in the [FieldReaderDataSource](#) section for details on configuration.

URLDataSource

This data source is often used with XPathEntityProcessor to fetch content from an underlying file:// or http:// location. Here's an example:

```
<dataSource name="a"
  type="URLDataSource"
  baseUrl="http://host:port/"
  encoding="UTF-8"
  connectionTimeout="5000"
  readTimeout="10000"/>
```

The URLDataSource type accepts these optional parameters:

Optional Parameter	Description
baseUrl	Specifies a new baseURL for pathnames. You can use this to specify host/port changes between Dev/QA/Prod environments. Using this attribute isolates the changes to be made to the solrconfig.xml
connectionTimeout	Specifies the length of time in milliseconds after which the connection should time out. The default value is 5000ms.

encoding	By default the encoding in the response header is used. You can use this property to override the default encoding.
readTimeout	Specifies the length of time in milliseconds after which a read operation should time out. The default value is 10000ms.

Entity Processors

Entity processors extract data, transform it, and add it to a Solr index. Examples of entities include views or tables in a data store.

Each processor has its own set of attributes, described in its own section below. In addition, there are non-specific attributes common to all entities which may be specified.

Attribute	Use
datasource	The name of a dataSource. Used if there are multiple datasources, specified, in which case each one must have a name.
name	Required. The unique name used to identify an entity.
pk	The primary key for the entity. It is optional, and required only when using delta-imports. It has no relation to the uniqueKey defined in schema.xml but they can both be the same. It is mandatory if you do delta-imports and then refers to the column name in \${dataimporter.delta.<column-name>} which is used as the primary key.
processor	Default is SQLEntityProcessor. Required only if the datasource is not RDBMS.
onError	Permissible values are (abort skip continue) . The default value is 'abort'. 'Skip' skips the current document. 'Continue' ignores the error and processing continues.
preImportDeleteQuery	Before a full-import command, use this query this to cleanup the index instead of using ':'. This is honored only on an entity that is an immediate sub-child of <document>.
postImportDeleteQuery	Similar to the above, but executed after the import has completed.
rootEntity	By default the entities immediately under the <document> are root entities. If this attribute is set to false, the entity directly falling under that entity will be treated as the root entity (and so on). For every row returned by the root entity, a document is created in Solr.

transformer	Optional. One or more transformers to be applied on this entity.
-------------	--

The SQL Entity Processor

The SqlEntityProcessor is the default processor. The associated [data source](#) should be a JDBC URL.

The entity attributes specific to this processor are shown in the table below.

Attribute	Use
query	Required. The SQL query used to select rows.
deltaQuery	SQL query used if the operation is delta-import. This query selects the primary keys of the rows which will be parts of the delta-update. The pks will be available to the deltaImportQuery through the variable \${dataimporter.delta.<column-name>}.
parentDeltaQuery	SQL query used if the operation is delta-import.
deletedPkQuery	SQL query used if the operation is delta-import.
deltaImportQuery	SQL query used if the operation is delta-import. If this is not present, DIH tries to construct the import query by(after identifying the delta) modifying the 'query' (this is error prone). There is a namespace \${dataimporter.delta.<column-name>} which can be used in this query. For example, select * from tbl where id=\${dataimporter.delta.id}.

The XPathEntityProcessor

This processor is used when indexing XML formatted data. The data source is typically [URLDataSource](#) or [FileDataSource](#). Xpath can also be used with the [FileListEntityProcessor](#) described below, to generate a document from each file.

The entity attributes unique to this processor are shown below.

Attribute	Use
Processor	Required. Must be set to "XpathEntityProcessor".
url	Required. HTTP URL or file location.
stream	Optional: Set to true for a large file or download.
forEach	Required unless you define useSolrAddSchema. The Xpath expression which demarcates each record. This will be used to set up the processing loop.

xsl	Optional: Its value (a URL or filesystem path) is the name of a resource used as a preprocessor for applying the XSL transformation.
useSolrAddSchema	Set this to true if the content is in the form of the standard Solr update XML schema.
flatten	Optional: If set true, then text from under all the tags is extracted into one field.

Each field element in the entity can have the following attributes as well as the default ones.

Attribute	Use
xpath	Required. The XPath expression which will extract the content from the record for this field. Only a subset of Xpath syntax is supported.
commonField	Optional. If true, then when this field is encountered in a record it will be copied to future records when creating a Solr document.

Example:

```
<!-- slashdot RSS Feed -->
<dataConfig>
  <dataSource type="HttpDataSource" />
  <document>
<entity name="slashdot"
  pk="link"
  url="http://rss.slashdot.org/Slashdot/slashdot"
  processor="XPathEntityProcessor"

  <!-- forEach sets up a processing loop ; here there are two expressions-->

  forEach="/RDF/channel | /RDF/item"
  transformer="DateFormatTransformer">
    <field column="source"
      xpath="/RDF/channel/title"
      commonField="true" />
    <field column="source-link"
      xpath="/RDF/channel/link"
      commonField="true"/>
    <field column="subject"
      xpath="/RDF/channel/subject"
      commonField="true" />
    <field column="title"
      xpath="/RDF/item/title" />
    <field column="link"
      xpath="/RDF/item/link" />
    <field column="description"
      xpath="/RDF/item/description" />
    <field column="creator"
      xpath="/RDF/item/creator" />
    <field column="item-subject"
      xpath="/RDF/item/subject" />
    <field column="date"
      xpath="/RDF/item/date"
        dateTimeFormat="yyyy-MM-dd'T'hh:mm:ss" />
    <field column="slash-department"
      xpath="/RDF/item/department" />
    <field column="slash-section"
      xpath="/RDF/item/section" />
    <field column="slash-comments"
      xpath="/RDF/item/comments" />
  </entity>
</document>
</dataConfig>
```

The **FileListEntityProcessor**

This processor is basically a wrapper, and is designed to generate a set of files satisfying conditions specified in the attributes which can then be passed to another processor, such as the [XPathEntityProcessor](#). The entity information for this processor would be nested within the `FileListEntity` entry. It generates four implicit fields: `fileAbsolutePath`, `fileSize`, `fileLastModified`, `fileName` which can be used in the nested processor. This processor does not use a data source.

The attributes specific to this processor are described in the table below:

Attribute	Use
<code>fileName</code>	Required. A regular expression pattern to identify files to be included.
<code>basedir</code>	Required. The base directory (absolute path).
<code>recursive</code>	Whether to search directories recursively. Default is 'false'.
<code>excludes</code>	A regular expression pattern to identify files which will be excluded.
<code>newerThan</code>	A date in the format <code>yyyy-MM-ddHH:mm:ss</code> or a date math expression (<code>NOW - 2YEARS</code>).
<code>olderThan</code>	A date, using the same formats as <code>newerThan</code> .
<code>rootEntity</code>	This should be set to false. This ensures that each row (filepath) emitted by this processor is considered to be a document.
<code>dataSource</code>	Must be set to null.

The example below shows the combination of the `FileListEntityProcessor` with another processor which will generate a set of fields from each file found.

```

<dataConfig>
<dataSource type="FileDataSource"/><document>
    <!-- this outer processor generates a list of files satisfying the conditions
        specified in the attributes -->
    <entity name="f" processor="FileListEntityProcessor"
        fileName=".xml"
        newerThan=" 'NOW-30DAYS' "
        recursive="true"
        rootEntity="false"
        dataSource="null"
        baseDir="/my/document/directory">

        <!-- this processor extracts content using Xpath from each file found -->

        <entity name="nested" processor="XPathEntityProcessor"
            forEach="/rootelement" url="${f.fileAbsolutePath}" >
            <field column="name" xpath="/rootelement/name"/>
            <field column="number" xpath="/rootelement/number"/>
        </entity>
    </entity>
</document>
</dataConfig>

```

LineEntityProcessor

This EntityProcessor reads all content from the data source on a line by line basis and returns a field called `rawLine` for each line read. The content is not parsed in any way; however, you may add transformers to manipulate the data within the `rawLine` field, or to create other additional fields.

The lines read can be filtered by two regular expressions specified with the `acceptLineRegex` and `omitLineRegex` attributes. The table below describes the LineEntityProcessor's attributes:

Attribute	Description
url	A required attribute that specifies the location of the input file in a way that is compatible with the configured data source. If this value is relative and you are using FileDataSource or URLDataSource, it assumed to be relative to <code>baseLoc</code> .
acceptLineRegex	An optional attribute that if present discards any line which does not match the <code>regExp</code> .
omitLineRegex	An optional attribute that is applied after any <code>acceptLineRegex</code> and that discards any line which matches this <code>regExp</code> .

For example:

```
<entity name="jc"
    processor="LineEntityProcessor"
    acceptLineRegex="^.*\.\xml$"
    omitLineRegex="/obsolete"
    url="file:///Volumes/ts/files.lis"
    rootEntity="false"
    dataSource="myURIreader1"
    transformer="RegexTransformer,DateFormatTransformer"
    >
...
...
```

While there are use cases where you might need to create a Solr document for each line read from a file, it is expected that in most cases that the lines read by this processor will consist of a pathname, which in turn will be consumed by another EntityProcessor, such as [XPathEntityProcessor](#).

PlainTextEntityProcessor

This EntityProcessor reads all content from the data source into a single implicit field called `plainText`. The content is not parsed in any way, however you may add transformers to manipulate the data within the `plainText` as needed, or to create other additional fields.

For example:

```
<entity processor="PlainTextEntityProcessor" name="x" url="http://abc.com/a.txt"
dataSource="data-source-name">
    <!-- copies the text to a field called 'text' in Solr-->
    <field column="plainText" name="text"/>
</entity>
```

Ensure that the `dataSource` is of type `DataSource<Reader>` (`FileDataSource`, `URLDataSource`).

Transformers

Transformers manipulate the fields in a document returned by an entity. A transformer can create new fields or modify existing ones. You must tell the entity which transformers your import operation will be using, by adding an attribute containing a comma separated list to the `<entity>` element.

```
<entity name="abcde"
transformer="org.apache.solr....,my.own.transformer,..." />
```

Specific transformation rules are then added to the attributes of a <field> element, as shown in the examples below. The transformers are applied in the order in which they are specified in the transformer attribute.

The Data Import Handler contains several built-in transformers. You can also write your own custom transformers, as described in the Solr Wiki (see <http://wiki.apache.org/solr/DIHCustomTransformer>). The ScriptTransformer (described below) offers an alternative method for writing your own transformers.

Solr includes the following built-in transformers:

Transformer Name	Use
ClobTransformer	Used to create a String out of a Clob type in database.
DateFormatTransformer	Parse date/time instances.
HTMLStripTransformer	Strip HTML from a field.
LogTransformer	Used to log data to log files or a console.
NumberFormatTransformer	Uses the NumberFormat class in java to parse a string into a number.
RegexTransformer	Use regular expressions to manipulate fields.
ScriptTransformer	Write transformers in Javascript or any other scripting language supported by Java. Requires Java 6.
TemplateTransformer	Transform a field using a template.

These transformers are described below.

ClobTransformer

You can use the ClobTransformer to create a string out of a CLOB in a database. A CLOB is a character large object: a collection of character data typically stored in a separate location that is referenced in the database. See http://en.wikipedia.org/wiki/Character_large_object. Here's an example of invoking the ClobTransformer.

```
<entity name="e" transformer="ClobTransformer" ...>
<field column="hugeTextField" clob="true" />
...
</entity>
```

The ClobTransformer accepts these attributes:

Attribute	Description
clob	Boolean value to signal if ClobTransformer should process this field or not. If this attribute is omitted, then the corresponding field is not transformed.
sourceColName	The source column to be used as input. If this is absent source and target are same

The DateFormatTransformer

This transformer converts dates from one format to another. This would be useful, for example, in a situation where you wanted to convert a field with a fully specified date/time into a less precise date format, for use in faceting.

DateFormatTransformer applies only on the fields with an attribute `dateTimeFormat`. Other fields are not modified.

This transformer recognizes the following attributes:

Attribute	Description
dateTimeFormat	The format used for parsing this field. This must comply with the syntax of the JavaSimpleDateFormat class.
sourceColName	The column on which the dateFormat is to be applied. If this is absent source and target are same.
locale	The locale to use for date transformations. If not specified, the ROOT locale will be used. It must be specified as language-country. For example, <code>en-US</code> .

Here is example code that returns the date rounded up to the month "2007-JUL":

```
<entity name="en" pk="id" transformer="DateTimeTransformer" ... >
...
<field column="date"
  sourceColName="fulldate"
  dateTimeFormat="yyyy-MMM" />
</entity>
```

The HTMLStripTransformer

You can use this transformer to strip HTML out of a field. For example:

```
<entity name="e" transformer="HTMLStripTransformer" ...>
<field column="htmlText" stripHTML="true" />
...
</entity>
```

There is one attribute for this transformer, `stripHTML`, which is a boolean value (true/false) to signal if the `HTMLStripTransformer` should process the field or not.

The LogTransformer

You can use this transformer to log data to the console or log files. For example:

```
<entity ...
transformer="LogTransformer"
logTemplate="The name is ${e.name}" logLevel="debug" >
....
</entity>
```

Unlike other transformers, the `LogTransformer` does not apply to any field, so the attributes are applied on the entity itself.

The NumberFormatTransformer

Use this transformer to parse a number from a string, converting it into the specified format, and optionally using a different locale.

`NumberFormatTransformer` will be applied only to fields with an attribute `formatStyle`.

This transformer recognizes the following attributes:

Attribute	Description
<code>formatStyle</code>	The format used for parsing this field. The value of the attribute must be one of (number percent integer currency). This uses the semantics of the Java <code>NumberFormat</code> class.
<code>sourceColName</code>	The column on which the <code>NumberFormat</code> is to be applied. This attribute is absent. The source column and the target column are the same.
<code>locale</code>	The locale to be used for parsing the strings. If this is absent, the ROOT locale is used. It must be specified as language-country. For example, <code>en-US</code> .

For example:

```

<entity name="en" pk="id" transformer="NumberFormatTransformer" ...>
  ...
  <!-- treat this field as UK pounds -->

  <field name="price_uk"
    column="price"
    formatStyle="currency"
    locale="en-UK" />
</entity>

```

The RegexTransformer

The regex transformer helps in extracting or manipulating values from fields (from the source) using Regular Expressions. The actual class name is `org.apache.solr.handler.dataimport.RegexTransformer`. But as it belongs to the default package the package-name can be omitted.

The table below describes the attributes recognized by the regex transformer.

Attribute	Description
regex	The regular expression that is used to match against the column or sourceColName's value(s). If replaceWith is absent, each regex <i>group</i> is taken as a value and a list of values is returned.
sourceColName	The column on which the regex is to be applied. If not present, then the source and target are identical.
splitBy	Used to split a string. It returns a list of values.
groupNames	A comma separated list of field column names, used where the regex contains groups and each group is to be saved to a different field. If some groups are not to be named leave a space between commas.
replaceWith	Used along with regex . It is equivalent to the method <code>new String(<sourceColVal>).replaceAll(<regex>, <replaceWith>)</code> .

Here is an example of configuring the regex transformer:

```
<entity name="foo" transformer="RegexTransformer"
query="select full_name , emailids from foo"/>
...
<field column="full_name" />
<field column="firstName" regex="Mr(\w*)\b.*" sourceColName="full_name" />
<field column="lastName" regex="Mr.*?\b(\w*)" sourceColName="full_name" />

<!-- another way of doing the same -->

<field column="fullName" regex="Mr(\w*)\b(.*)" groupNames="firstName,lastName" />
<field column="mailId" splitBy="," sourceColName="emailids" />
</entity>
```

In this example, `regex` and `sourceColName` are custom attributes used by the transformer. The transformer reads the field `full_name` from the resultset and transforms it to two new target fields, `firstName` and `lastName`. Even though the query returned only one column, `full_name`, in the result set, the Solr document gets two extra fields `firstName` and `lastName` which are "derived" fields. These new fields are only created if the regexp matches.

The `emailids` field in the table can be a comma-separated value. It ends up producing one or more email IDs, and we expect the `mailId` to be a multivalued field in Solr.

Note that this transformer can either be used to split a string into tokens based on a `splitBy` pattern, or to perform a string substitution as per `replaceWith`, or it can assign groups within a pattern to a list of `groupNames`. It decides what it is to do based upon the above attributes `splitBy`, `replaceWith` and `groupNames` which are looked for in order. This first one found is acted upon and other unrelated attributes are ignored.

The ScriptTransformer

The script transformer allows arbitrary transformer functions to be written in any scripting language supported by Java, such as Javascript, JRuby, Jython, Groovy, or BeanShell. Javascript is integrated into Java 6; you'll need to integrate other languages yourself.

Each function you write must accept a `row` variable (which corresponds to a Java `Map<String, Object>`, thus permitting `get`, `put`, `remove` operations). Thus you can modify the value of an existing field or add new fields. The return value of the function is the returned object.

The script is inserted into the DIH configuration file file at the top level and is called once for each row.

Here is a simple example.

```

<dataconfig>

    <!-- simple script to generate a new row, converting a temperature from Fahrenheit
        to Centigrade -->

        <script>
<CDATA
    function f2c(row) {  var tempf, tempc;  tempf = row.get('temp_f');  if (tempf != null) {      tempc = (tempf - 32.0)*5.0/9.0
        row.put('temp_c', temp_c);
    }
    return row;
}
>
</script>
<document>

    <!-- the function is specified as an entity attribute -->

    <entity name="e1" pk="id" transformer="script:f2c" query="select * from X">
    ....
    </entity>
</document>
</dataConfig>

```

The TemplateTransformer

You can use the template transformer to construct or modify a field value, perhaps using the value of other fields. You can insert extra text into the template.

```

<entity name="en" pk="id" transformer="TemplateTransformer" ...>
...
<!-- generate a full address from fields containing the component parts -->
<field column="full_address"
    template="$en.\{street\},$en.\{city\},$en.\{zip\}" />
</entity>

```

Special Commands for the Data Import Handler

You can pass special commands to the DIH by adding any of the variables listed below to any row returned by any component:

Variable	Description
----------	-------------

\$skipDoc	Skip the current document; that is, do not add it to Solr. The value can be the string true false.
\$skipRow	Skip the current row. The document will be added with rows from other entities. The value can be the string true false
\$docBoost	Boost the current document. The boost value can be a number or the toString conversion of a number.
\$deleteDocById	Delete a document from Solr with this ID. The value has to be the uniqueKey value of the document.
\$deleteDocByQuery	Delete documents from Solr using this query. The value must be a Solr Query.

De-Duplication

Preventing duplicate or near duplicate documents from entering an index or tagging documents with a signature/fingerprint for duplicate field collapsing can be efficiently achieved with a low collision or fuzzy hash algorithm. Solr natively supports de-duplication techniques of this type via the `<Signature>` class and allows for the easy addition of new hash/signature implementations. A Signature can be implemented several ways:

Method	Description
MD5Signature	128 bit hash used for exact duplicate detection.
Lookup3Signature	64 bit hash used for exact duplicate detection, much faster than MD5 and smaller to index
TextProfileSignature	Fuzzy hashing implementation from nutch for near duplicate detection. Its tunable but works best on longer text.

Other, more sophisticated algorithms for fuzzy/near hashing can be added later.

- ⚠ Adding in the deduplication process will change the `allowDups` setting so that it applies to an update Term (with `signatureField` in this case) rather than the unique field Term. Of course the `signatureField` could be the unique field, but generally you want the unique field to be unique. When a document is added, a signature will automatically be generated and attached to the document in the specified `signatureField`.

Configuration Options

In `solrconfig.xml`

The `SignatureUpdateProcessorFactory` has to be registered in the `solrconfig.xml` as part of the [`UpdateRequestProcessorChain`](#):

```

<updateRequestProcessorChain name="dedupe">
  <processor class="solr.processor.SignatureUpdateProcessorFactory">
    <bool name="enabled">true</bool>
    <str name="signatureField">id</str>
    <bool name="overwriteDupes">false</bool>
    <str name="fields">name,features,cat</str>
    <str name="signatureClass">solr.processor.Lookup3Signature</str>
  </processor>
</updateRequestProcessorChain>

```

Setting	Default	Description
signatureClass	org.apache.solr.update.processor.Lookup3Signature	A Signature implementation for generating a signature hash.
fields	all fields	The fields to use to generate the signature hash in a comma separated list. By default, all fields on the document will be used.
signatureField	signatureField	The name of the field used to hold the fingerprint/signature. Be sure the field is defined in schema.xml.
enabled	true	Enable/disable deduplication factory processing

In schema.xml

If you are using a separate field for storing the signature you must have it indexed:

```

<field name="signature" type="string" stored="true" indexed="true" multiValued="false"
/>

```

Be sure to change your update handlers to use the defined chain, i.e.

```
<requestHandler name="/update" >
  <lst name="defaults">
    <str name="update.chain">dedupe</str>
  </lst>
</requestHandler>
```

-  The update processor can also be specified per request with a parameter of update.chain=dedupe.

Detecting Languages During Indexing

Solr can identify languages and map text to language-specific fields during indexing using the langid UpdateRequestProcessor. Solr supports two implementations of this feature:

- Tika's language detection feature: <http://tika.apache.org/0.10/detection.html>
- LangDetect language detection: <http://code.google.com/p/language-detection/>

You can see a comparison between the two implementations here:

<http://blog.mikemccandless.com/2011/10/accuracy-and-performance-of-googles.html>. In general, the LangDetect implementation supports more languages with higher performance.

For specific information on each of these language identification implementations, including a list of supported languages for each, see the relevant project websites. For more information about the langid UpdateRequestProcessor, see the Solr wiki: <http://wiki.apache.org/solr/LanguageDetection>. For more information about language analysis in Solr, see [Language Analysis](#).

Configuring Language Detection

You can configure the langid UpdateRequestProcessor in `solrconfig.xml`. Both implementations take the same parameters, which are described in the following section. At a minimum, you must specify the fields for language identification and a field for the resulting language code.

Configuring Tika Language Detection

Here is an example of a minimal Tika langid configuration in `solrconfig.xml`:

```
<processor
  class="org.apache.solr.update.processor.TikaLanguageIdentifierUpdateProcessorFactory">
  <lst name="defaults">
    <str name="langid.fl">title,subject,text,keywords</str>
    <str name="langid.langField">language_s</str>
  </lst>
</processor>
```

Configuring LangDetect Language Detection

Here is an example of a minimal LangDetect langid configuration in `solrconfig.xml`:

```
<processor
  class="org.apache.solr.update.processor.LangDetectLanguageIdentifierUpdateProcessorFactory"
<lst name="defaults">
  <str name="langid.fl">title,subject,text,keywords</str>
  <str name="langid.langField">language_s</str>
</lst>
</processor>
```

langid Parameters

As previously mentioned, both implementations of the `langid` `UpdateRequestProcessor` take the same parameters.

Parameter	Type	Default	Required	Description
<code>langid</code>	Boolean	true	no	Enables and disables language detection.
<code>langid.fl</code>	string	none	yes	A comma- or space-delimited list of fields to be processed by <code>langid</code> .
<code>langid.langField</code>	string	none	yes	Specifies the field for the returned language code.
<code>langid.langsField</code>	multivalued string	none	no	Specifies the field for a list of returned language codes. If you use <code>langid.map.individual</code> , each detected language will be added to this field.
<code>langid.overwrite</code>	Boolean	false	no	Specifies whether the content of the <code>langField</code> and <code>langsField</code> fields will be overwritten if they already contain values.

langid.threshold	float	0.5	no	Specifies a threshold value between 0 and 1 that the language identification score must reach before langid accepts it. With longer text fields, a high threshold such as 0.8 will give good results. For shorter text fields, you may need to lower the threshold for language identification, though you will be risking somewhat lower quality results. We recommend experimenting with your data to tune your results.
langid.whitelist	string	none	no	Specifies a list of allowed language identification codes. Use this in combination with langid.map to ensure that you only index documents into fields that are in your schema.
langid.map	Boolean	false	no	Enables field name mapping. If true, Solr will map field names for all fields listed in langid.fl.
langid.map.fl	string	none	no	A comma-separated list of fields for langid.map that is different than the fields specified in langid.fl.
langid.map.keepOrig	Boolean	false	no	If true, Solr will copy the field during the field name mapping process, leaving the original field in place.
langid.map.individual	Boolean	false	no	If true, Solr will detect and map languages for each field individually.
langid.map.individual.fl	string	none	no	A comma-separated list of fields for use with langid.map.individual that is different than the fields specified in langid.fl.

langid.fallbackFields	string	none	no	If no language is detected that meets the langid.threshold score, or if the detected language is not on the langid.whitelist, this field specifies language codes to be used as fallback values. If no appropriate fallback languages are found, Solr will use the language code specified in langid.fallback.
langid.fallback	string	none	no	Specifies a language code to use if no language is detected or specified in langid.fallbackFields.
langid.map.lcmap	string	none	no	A space-separated list specifying a language code map. For example, you might use this to make Chinese, Japanese, and Korean language fields to a CJK field in your schema, or to map American and British English fields to a single EN field.
langid.map.pattern	Java regular expression	none	no	By default, fields are mapped as <field>_<language>. To change this pattern, you can specify a Java regular expression in this parameter.
langid.map.replace	Java replace	none	no	By default, fields are mapped as <field>_<language>. To change this pattern, you can specify a Java replace in this parameter.
langid.enforceSchema	Boolean	true	no	If false, the langid processor does not validate field names against your schema. This may be useful if you plan to rename or delete fields later in the UpdateChain.

Content Streams

When Solr RequestHandlers are accessed using path based URLs, the `SolrQueryRequest` object containing the parameters of the request may also contain a list of ContentStreams containing bulk data for the request. (The name `SolrQueryRequest` is a bit misleading: it is involved in all requests, regardless of whether it is a query request or an update request.)

Stream Sources

Currently RequestHandlers can get content streams in a variety of ways:

- For multipart file uploads, each file is passed as a stream.
- For POST requests where the content-type is not `application/x-www-form-urlencoded`, the raw POST body is passed as a stream. The full POST body is parsed as parameters and included in the Solr parameters.
- The contents of parameter `stream.body` is passed as a stream.
- If remote streaming is enabled and URL content is called for during request handling, the contents of each `stream.url` and `stream.file` parameters are fetched and passed as a stream.

By default, curl sends a `contentType="application/x-www-form-urlencoded"` header. If you need to test a `SolrContentHeader` content stream, you will need to set the content type with the "-H" flag.

RemoteStreaming

Remote streaming lets you send the contents of a URL as a stream to a given SolrRequestHandler. You could use remote streaming to send a remote or local file to an update plugin. For security reasons, remote streaming is disabled in the `solrconfig.xml` included in the example directory.



If you enable streaming, be aware that this allows *anyone* to send a request to any URL or local file. If dump is enabled, it will allow anyone to view any file on your system.

```
<!--Make sure your system has authentication before enabling remote streaming!-->
<requestParsers enableRemoteStreaming="true" multipartUploadLimitInKB="2048" />
```

Debugging Requests

The example `solrconfig.xml` includes a "dump" RequestHandler:

```
<requestHandler name="/debug/dump" class="solr.DumpRequestHandler" />
```

This handler simply outputs the contents of the SolrQueryRequest using the specified writer type wt. This is a useful tool to help understand what streams are available to the RequestHandlers.

UIMA Integration

You can integrate the Apache Unstructured Information Management Architecture ([UIMA](#)) with Solr. UIMA lets you define custom pipelines of Analysis Engines that incrementally add metadata to your documents as annotations.

For more information about Solr UIMA integration, see <https://wiki.apache.org/solr/SolrUIMA>.

Configuring UIMA

The SolrUIMA UpdateRequestProcessor is a custom update request processor that takes documents being indexed, sends them to a UIMA pipeline, and then returns the documents enriched with the specified metadata. To configure UIMA for Solr, follow these steps:

1. Copy `apache-solr-uima-3.x.0.jar` (under `/apache-solr-3.x.0/dist/`) and its libraries (under `contrib/uima/lib`) to a Solr libraries directory, or set `<lib/>` tags in `solrconfig.xml` appropriately to point to those jar files:

```
<lib dir="..../contrib/uima/lib" />
<lib dir="..../dist/" regex="apache-solr-uima-\d.*\.jar" />
```

2. Modify `schema.xml`, adding your desired metadata fields specifying proper values for type, indexed, stored, and multiValued options. For example:

```
<field name="language" type="string" indexed="true" stored="true"
required="false"/>
<field name="concept" type="string" indexed="true" stored="true"
multiValued="true" required="false"/>
<field name="sentence" type="text" indexed="true" stored="true"
multiValued="true" required="false" />
```

3. Add the following snippet to `solrconfig.xml`:

```
<updateRequestProcessorChain name="uima">
  <processor
    class="org.apache.solr.uima.processor.UIMAUpdateRequestProcessorFactory">
    <lst name="uimaConfig">
      <lst name="runtimeParameters">
        <str name="keyword_apikey">VALID_ALCHEMYAPI_KEY</str>
        <str name="concept_apikey">VALID_ALCHEMYAPI_KEY</str>
        <str name="lang_apikey">VALID_ALCHEMYAPI_KEY</str>
        <str name="cat_apikey">VALID_ALCHEMYAPI_KEY</str>
        <str name="entities_apikey">VALID_ALCHEMYAPI_KEY</str>
        <str name="oc_licenseID">VALID_OPENCALAIS_KEY</str>
      </lst>
    </lsts>
  </processor>
</updateRequestProcessorChain>
```

```
<str
name="analysisEngine">/org/apache/uima/desc/OverridingParamsExtServicesAE.xml</str>
<!-- Set to true if you want to continue indexing even if text processing fails.
Default is false. That is, Solr throws RuntimeException and
never indexed documents entirely in your session. -->
<bool name="ignoreErrors">true</bool>
<!-- This is optional. It is used for logging when text processing fails.
If logField is not specified, uniqueKey will be used as logField.
<str name="logField">id</str>
-->
<lst name="analyzeFields">
<bool name="merge">false</bool>
<arr name="fields">
<str>text</str>
</arr>
</lst>
<lst name="fieldMappings">
<lst name="type">
<str name="name">org.apache.uima.alchemy.ts.concept.ConceptFS</str>
<lst name="mapping">
<str name="feature">text</str>
<str name="field">concept</str>
</lst>
</lst>
<lst name="type">
<str name="name">org.apache.uima.alchemy.ts.language.LanguageFS</str>
<lst name="mapping">
<str name="feature">language</str>
<str name="field">language</str>
</lst>
</lst>
<lst name="type">
<str name="name">org.apache.uima.SentenceAnnotation</str>
<lst name="mapping">
<str name="feature">coveredText</str>
<str name="field">sentence</str>
</lst>
</lst>
</lst>
</processor>
<processor class="solr.LogUpdateProcessorFactory" />
<processor class="solr.RunUpdateProcessorFactory" />
</updateRequestProcessorChain>
```



VALID_ALCHEMYAPI_KEY is your AlchemyAPI Access Key. You need to register an AlchemyAPI Access key to use AlchemyAPI services:
<http://www.alchemyapi.com/api/register.html>.

VALID_OPENCALAIS_KEY is your Calais Service Key. You need to register a Calais Service key to use the Calais services: <http://www.opencalais.com/apikey>.

analysisEngine must contain an AE descriptor inside the specified path in the classpath.

analyzeFields must contain the input fields that need to be analyzed by UIMA. If merge=true then their content will be merged and analyzed only once.

Field mapping describes which features of which types should go in a field.

4. In your solrconfig.xml replace the existing default UpdateRequestHandler or create a new UpdateRequestHandler:

```
<requestHandler name="/update" class="solr.XmlUpdateRequestHandler">
  <lst name="defaults">
    <str name="update.processor">uima</str>
  </lst>
</requestHandler>
```

Once you are done with the configuration your documents will be automatically enriched with the specified fields when you index them.

Searching

This section describes how Solr works with search requests. It covers the following topics:

- [Overview of Searching in Solr](#): An introduction to searching with Solr.
- [Relevance](#): Conceptual information about understanding relevance in search results.
- [Query Syntax and Parsing](#): A brief conceptual overview of query syntax and parsing. It also contains the following sub-sections:
 - [Common Query Parameters](#): No matter the query parser, there are several parameters that are common to all of them.
 - [The Standard Query Parser](#): Detailed information about the standard Lucene query parser.
 - [The DisMax Query Parser](#): Detailed information about Solr's DisMax query parser.
 - [The Extended DisMax Query Parser](#): Detailed information about Solr's Extended DisMax (eDisMax) Query Parser.
 - [Local Parameters in Queries](#): How to add local arguments to queries.
 - [Other Parsers](#): More parsers designed for use in specific situations.
- [Highlighting](#): Detailed information about Solr's highlighting utilities.
- [MoreLikeThis](#): Detailed information about Solr's similar results query component.
- [Faceting](#): Detailed information about categorizing search results based on indexed terms.
- [Suggester](#): Detailed information about Solr's powerful autosuggest component.
- [Function Queries](#): Detailed information about parameters for generating relevancy scores using values from one or more numeric fields.
- [Spatial Search](#): How to use Solr's spatial search capabilities.
- [The Terms Component](#): Detailed information about accessing indexed terms and the documents that include them.
- [The Term Vector Component](#): How to get term information about specific documents.
- [The Stats Component](#): How to return information from numeric fields within a document set.
- [The Query Elevation Component](#): How to force documents to the top of the results for certain queries.
- [Response Writers](#): Detailed information about configuring and using Solr's response writers.
- [Near Real Time Searching](#): How to include documents in search results nearly immediately after they are indexed.
- [RealTime Get](#): How to get the latest version of a document without opening a searcher.

- [Result Grouping](#): Detailed information about grouping results based on common field values.
- [Spell Checking](#): Detailed information about Solr's spelling checker.

Overview of Searching in Solr

Solr offers a rich, flexible set of features for search. To understand the extent of this flexibility, it's helpful to begin with an overview of the steps and components involved in a Solr search.

When a user runs a search in Solr, the search query is processed by a **request handler**. A request handler is a Solr plug-in that defines the logic to be used when Solr processes a request. Solr supports a variety of request handlers. Some are designed for processing search queries, while others manage tasks such as index replication.

Search applications select a particular request handler by default. In addition, applications can be configured to allow users to override the default selection in preference of a different request handler.

To process a search query, a request handler calls a **query parser**, which interprets the terms and parameters of a query. Different query parsers support different syntax. The default query parser is the [DisMax](#) query parser. Solr also includes an earlier "standard" (Lucene) query parser, and an [Extended DisMax](#) (eDisMax) query parser. The [standard](#) query parser's syntax allows for greater precision in searches, but the DisMax query parser is much more tolerant of errors. The DisMax query parser is designed to provide an experience similar to that of popular search engines such as Google, which rarely display syntax errors to users. The Extended DisMax query parser is an improved version of DisMax that handles the full Lucene query syntax while still tolerating syntax errors. It also includes several additional features.

In addition, there are [common query parameters](#) that are accepted by all query parsers.

Input to a query parser can include:

- search strings---that is, *terms* to search for in the index
- *parameters for fine-tuning the query* by increasing the importance of particular strings or fields, by applying Boolean logic among the search terms, or by excluding content from the search results
- *parameters for controlling the presentation of the query response*, such as specifying the order in which results are to be presented or limiting the response to particular fields of the search application's schema.

Search parameters may also specify a **query filter**. As part of a search response, a query filter runs a query against the entire index and caches the results. Because Solr allocates a separate cache for filter queries, the strategic use of filter queries can improve search performance. (Despite their similar names, query filters are not related to analysis filters. Query filters perform queries at search time against data already in the index, while analysis filters, such as Tokenizers, parse content for indexing, following specified rules).

A search query can request that certain terms be highlighted in the search response; that is, the selected terms will be displayed in colored boxes so that they "jump out" on the screen of search results. **Highlighting** can make it easier to find relevant passages in long documents returned in a search. Solr supports multi-term highlighting. Solr includes a rich set of search parameters for controlling how terms are highlighted.

Search responses can also be configured to include **snippets** (document excerpts) featuring highlighted text. Popular search engines such as Google and Yahoo! return snippets in their search results: 3-4 lines of text offering a description of a search result.

To help users zero in on the content they're looking for, Solr supports two special ways of grouping search results to aid further exploration: faceting and clustering.

Faceting is the arrangement of search results into categories (which are based on indexed terms). Within each category, Solr reports on the number of hits for relevant term, which is called a facet constraint. Faceting makes it easy for users to explore search results on sites such as movie sites and product review sites, where there are many categories and many items within a category.

The image below shows an example of faceting from the CNET Web site, which was the first site to use Solr.

The screenshot shows the CNET Digital cameras search results page. At the top, a header reads "Digital cameras". Below it, a section titled "Refine your results" contains four facets: "Manufacturer", "Resolution", "Zoom range", and "More". The "Manufacturer" facet shows options like Canon USA (5), Sony (2), Nikon (2), Olympus (6), and Pentax (2). The "Resolution" facet shows 6 megapixels (3) and 8 megapixels and up (14). The "Zoom range" facet shows 3X to 4X (11) and 8X to 12X (1). The "More" facet lists LCD size, Image stabilizer, Flash memory, Still image format, and Maximum ISO, with a "See all" link. To the left of the facets, a callout box explains that "Manufacturer is a facet, a way of categorizing the results" and "Canon, Sony, and Nikon are constraints, or facet values". At the bottom, a breadcrumb trail shows "you selected: \$400 - \$500, SLR, remove all". The main search results area displays "17 results" with a table showing items like the "Canon EOS Rebel XS (silver, with 18-55mm lens)" for "\$459 to \$699 at 15 stores". A "Regular search results list" button is visible. Navigation links "1 2 next" are at the bottom right. Callout boxes explain the "facet count or constraint count shows how many results match each value" and the "breadcrumb trail shows what constraints have already been applied and allows for their removal".

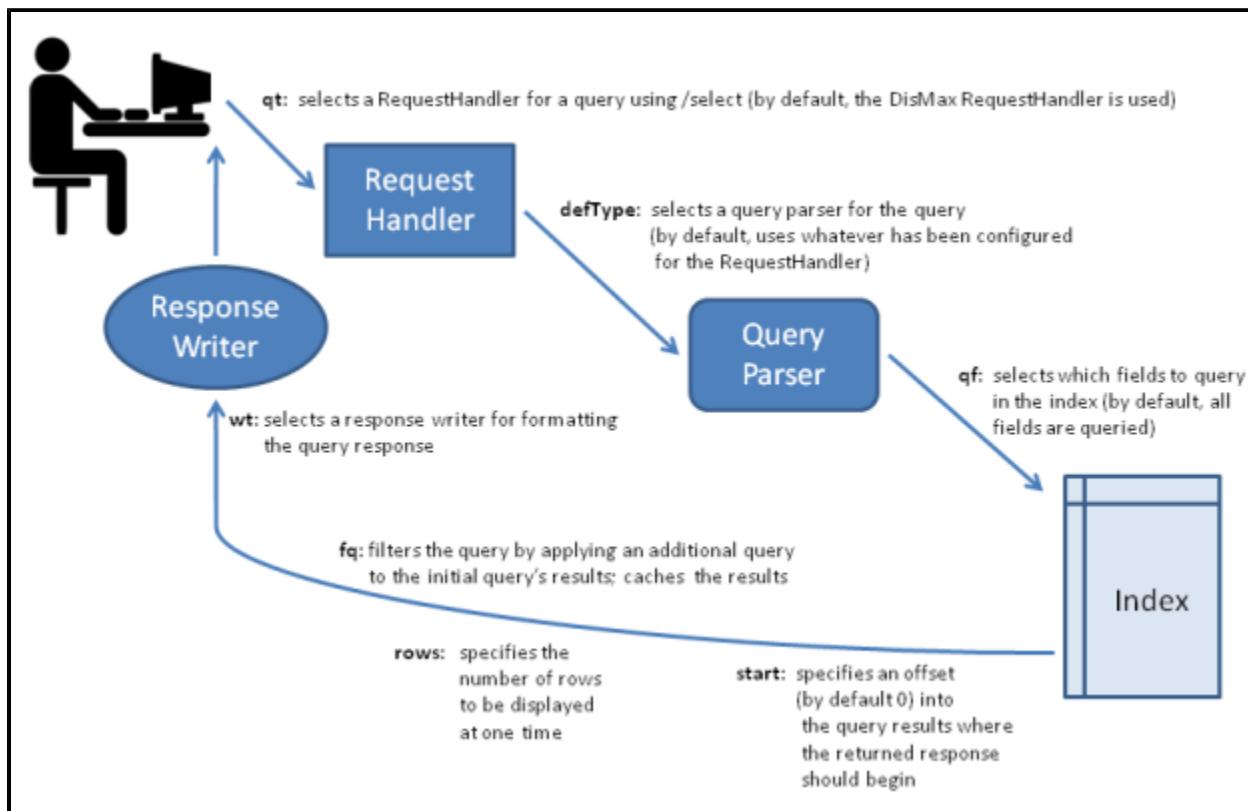
Faceting makes use of fields defined when the search applications were indexed. In the example above, these fields include categories of information that are useful for describing digital cameras: manufacturer, resolution, and zoom range.

Clustering groups search results by similarities discovered when a search is executed, rather than when content is indexed. The results of clustering often lack the neat hierarchical organization found in faceted search results, but clustering can be useful nonetheless. It can reveal unexpected commonalities among search results, and it can help users rule out content that isn't pertinent to what they're really searching for.

Solr also supports a feature called [MoreLikeThis](#), which enables users to submit new queries that focus on particular terms returned in an earlier query. MoreLikeThis queries can make use of faceting or clustering to provide additional aid to users.

A Solr component called a [response writer](#) manages the final presentation of the query response. Solr includes a variety of response writers, including an [XML Response Writer](#) and a [JSON Response Writer](#).

The diagram below summarizes some key elements of the search process.



The Velocity Search UI

Solr includes a sample search UI based on the [VelocityResponseWriter](#) (also known as Solritas) than demonstrates several useful features, such as searching, facetting, highlighting, autocomplete, and geospatial searching.

The VelocityResponseWriter is no longer built into the code. Its jar and dependencies must be added (via <lib> or solr/home lib inclusion), and must be registered in `solrconfig.xml` like this:

```
<queryResponseWriter name="velocity" class="solr.VelocityResponseWriter"/>
```

You can access the Velocity sample Search UI here: <http://localhost:8983/solr/browse>

Solritas

Solritas

localhost:8983/solr/browse

Solr Admin

Apache Solr

Examples: Simple [Spatial](#)

Find: [Submit Query](#) [Reset](#)

Boost by Price

Field Facets

cat

- Electronics (14)
- Memory (3)
- Connector (2)
- Graphics Card (2)
- Hard Drive (2)
- Monitor (2)
- Search (2)
- Software (2)
- Camera (1)
- Copier (1)
- Multifunction Printer (1)
- Music (1)
- Printer (1)
- Scanner (1)

manu_exact

- Apache Software Foundation (2)
- Belkin (2)
- Canon Inc. (2)
- Corsair Microsystems Inc. (2)
- A-DATA Technology Inc. (1)
- ASUS Computer Inc. (1)
- ATI Technologies (1)
- Apple Computer Inc. (1)
- Dell, Inc. (1)
- Maxtor Corp. (1)
- Samsung Electronics Co., Ltd. (1)
- ViewSonic Corp. (1)

Query Facets

17 results found in 98 ms Page 1 of 2

Test with some GB18030 encoded characters [More Like This](#)

Price: \$0.00
Features: No accents here 这是一个功能 This is a feature (translated) 这份文件是很有光泽 This document is very shiny (translated)
In Stock: true

Samsung SpinPoint P120 SP2514N - hard drive - 250 GB - ATA-133 [More Like This](#)

Price: \$92.00
Features: 7200RPM, 8MB cache, IDE Ultra ATA-133 NoiseGuard, SilentSeek technology, Fluid Dynamic Bearing (FDB) motor
In Stock: true

Maxtor DiamondMax 11 - hard drive - 500 GB - SATA-300 [More Like This](#)

Price: \$350.00
Features: SATA 3.0Gb/s, NCQ 8.5ms seek 16MB cache
In Stock: true

Belkin Mobile Power Cord for iPod w/ Dock [More Like This](#)

Price: \$19.95
Features: car power adapter, white
In Stock: false

iPod & iPod Mini USB 2.0 Cable [More Like This](#)

Price: \$11.50
Features: car power adapter for iPod, white

The Velocity Search UI

For more information about the Velocity Search UI, see <https://wiki.apache.org/solr/VelocityResponseWriter>.

Relevance

Relevance is the degree to which a query response satisfies a user who is searching for information.

The relevance of a query response depends on the context in which the query was performed. A single search application may be used in different contexts by users with different needs and expectations. For example, a search engine of climate data might be used by a university researcher studying long-term climate trends, a farmer interested in calculating the likely date of the last frost of spring, a civil engineer interested in rainfall patterns and the frequency of floods, and a college student planning a vacation to a region and wondering what to pack. Because the motivations of these users vary, the relevance of any particular response to a query will vary as well.

How comprehensive should query responses be? Like relevance in general, the answer to this question depends on the context of a search. The cost of *not* finding a particular document in response to a query is high in some contexts, such as a legal e-discovery search in response to a subpoena, and quite low in others, such as a search for a cake recipe on a Web site with dozens or hundreds of cake recipes. When configuring Solr, you should weigh comprehensiveness against other factors such as timeliness and ease-of-use.

The e-discovery and recipe examples demonstrate the importance of two concepts related to relevance:

- **Precision** is the percentage of documents in the returned results that are relevant.
- **Recall** is the percentage of relevant results returned out of all relevant results in the system. Obtaining perfect recall is trivial: simply return every document in the collection for every query.

Returning to the examples above, it's important for an e-discovery search application to have 100% recall returning all the documents that are relevant to a subpoena. It's far less important that a recipe application offer this degree of precision, however. In some cases, returning too many results in casual contexts could overwhelm users. In some contexts, returning fewer results that have a higher likelihood of relevance may be the best approach.

Using the concepts of precision and recall, it's possible to quantify relevance across users and queries for a collection of documents. A perfect system would have 100% precision and 100% recall for every user and every query. In other words, it would retrieve all the relevant documents and nothing else. In practical terms, when talking about precision and recall in real systems, it is common to focus on precision and recall at a certain number of results, the most common (and useful) being ten results.

Through faceting, query filters, and other search components, a Solr application can be configured with the flexibility to help users fine-tune their searches in order to return the most relevant results for users. That is, Solr can be configured to balance precision and recall to meet the needs of a particular user community.

The configuration of a Solr application should take into account:

- the needs of the application's various users (which can include ease of use and speed of response, in addition to strictly informational needs)
- the categories that are meaningful to these users in their various contexts (e.g., dates, product categories, or regions)
- any inherent relevance of documents (e.g., it might make sense to ensure that an official product description or FAQ is always returned near the top of the search results)
- whether or not the age of documents matters significantly (in some contexts, the most recent documents might always be the most important)

Keeping all these factors in mind, it's often helpful in the planning stages of a Solr deployment to sketch out the types of responses you think the search application should return for sample queries. Once the application is up and running, you can employ a series of testing methodologies, such as focus groups, in-house testing, [TREC](#) tests and A/B testing to fine tune the configuration of the application to best meet the needs of its users.

For more information about relevance, see Grant Ingersoll's tech article [Debugging Search Application Relevance Issues](#) which is available on the Lucid Imagination Web site.

Query Syntax and Parsing

Solr supports several query parsers, offering search application designers great flexibility in controlling how queries are parsed.

This section explains how to specify the query parser to be used. It also describes the syntax and features supported by the main query parsers included with Solr and describes some other parsers that may be useful for particular situations. There are some query parameters common to all Solr parsers; these are discussed in the section [Common Query Parameters](#).

The parsers discussed in this Guide are:

- [The Standard Query Parser](#)
- [The DisMax Query Parser](#)
- [The Extended DisMax Query Parser](#)
- [Other Parsers](#)

The query parser plugins are all subclasses of

http://lucene.apache.org/solr/4_0_0/solr-core/org/apache/solr/search/QParserPlugin.html. If you have custom parsing needs, you may want to extend that class to create your own query parser.

For more detailed information about the many query parsers available in Solr, see
<https://wiki.apache.org/solr/SolrQuerySyntax>.

Common Query Parameters

The table below summarizes Solr's common query parameters, which are supported by the [Standard](#), [DisMax](#), and [eDisMax](#) Request Handlers.

LucidWorks strongly recommends that any future SolrRequestHandlers support these parameters, as well.

Parameter	Description
defType	Selects the query parser to be used to process the query.
sort	Sorts the response to a query in either ascending or descending order based on the response's score or another specified characteristic.
start	Specifies an offset (by default, 0) into the responses at which Solr should begin displaying content.
rows	Controls how many rows of responses are displayed at a time (default value: 10)
fq	Applies a filter query to the search results.
fl	With version 3.6, Solr limited the query's responses to a listed set of fields. With version 4.0, this parameter returns only the score.
debug	Request additional debugging information in the response. Specifying the <code>debug=timing</code> parameter returns just the timing information; specifying the <code>debug=results</code> parameter returns "explain" information for each of the documents returned; specifying the <code>debug=query</code> parameter returns all of the debug information.
explainOther	Allows clients to specify a Lucene query to identify a set of documents. If non-blank, the explain info of each document which matches this query, relative to the main query (specified by the <code>q</code> parameter) will be returned along with the rest of the debugging information.
timeAllowed	Defines the time allowed for the query to be processed. If the time elapses before the query response is complete, partial information may be returned.
omitHeader	Excludes the header from the returned results, if set to true. The header contains information about the request, such as the time the request took to complete. The default is false.
wt	Specifies the Response Writer to be used to format the query response.
cache=false	By default, Solr caches the results of all queries and filter queries. Set <code>cache=false</code> to disable caching of the results of a query.

The following sections describe these parameters in detail.

The defType Parameter

The `defType` parameter selects the query parser that Solr should use to process the request. For example:

```
defType=dismax
```

In Solr 1.3 and later, the query parser is set to `dismax` by default.

The sort Parameter

The `sort` parameter arranges search results in either ascending (`asc`) or descending (`desc`) order. The parameter can be used with either numerical or alphabetical content.

Solr can sort query responses according to document scores or the value of any indexed field with a single value (that is, any field whose attributes in `schema.xml` include `multiValued="false"` and `indexed="true"`), provided that:

- the field is non-tokenized (that is, the field has no analyzer and its contents have been parsed into tokens, which would make the sorting inconsistent), or
- the field uses an analyzer (such as the `KeywordTokenizer`) that produces only a single term.

If you want to be able to sort on a field whose contents you want to tokenize to facilitate searching, use the `<copyField>` directive in the `schema.xml` file to clone the field. Then search on the field and sort on its clone.

The table explains how Solr responds to various settings of the `sort` parameter.

Example	Result
	If the <code>sort</code> parameter is omitted, sorting is performed as though the parameter were set to <code>score desc</code> .
<code>score desc</code>	Sorts in descending order from the highest score to the lowest score.
<code>price asc</code>	Sorts in ascending order of the <code>price</code> field
<code>inStock desc, price asc</code>	Sorts by the contents of the <code>inStock</code> field in descending order, then within those results sorts in ascending order by the contents of the <code>price</code> field.

Regarding the `sort` parameter's arguments:

- A sort ordering must include a field name (or `score` as a pseudo field), followed by whitespace (escaped as `+` or `%20` in URL strings), followed by a sort direction (`asc` or `desc`).
- Multiple sort orderings can be separated by a comma, using this syntax: `sort=<field name>+<direction>,<field name>+<direction>],...`

The start Parameter

When specified, the `start` parameter specifies an offset into a query's result set and instructs Solr to begin displaying results from this offset.

The default value is "0". In other words, by default, Solr returns results without an offset, beginning where the results themselves begin.

Setting the `start` parameter to some other number, such as 3, causes Solr to skip over the preceding records and start at the document identified by the offset.

You can use the `start` parameter this way for paging. For example, if the `rows` parameter is set to 10, you could display three successive pages of results by setting `start` to 0, then re-issuing the same query and setting `start` to 10, then issuing the query again and setting `start` to 20.

The rows Parameter

You can use the `rows` parameter to paginate results from a query. The parameter specifies the maximum number of documents from the complete result set that Solr should return to the client at one time.

The default value is 10. That is, by default, Solr returns 10 documents at a time in response to a query.

The fq (Filter Query) Parameter

The `fq` parameter defines a query that can be used to restrict the superset of documents that can be returned, without influencing score. It can be very useful for speeding up complex queries, since the queries specified with `fq` are cached independently of the main query. When a later query uses the same filter, there's a cache hit, and filter results are returned quickly from the cache.

When using the `fq` parameter, keep in mind the following:

- The `fq` parameter can be specified multiple times in a query. Documents will only be included in the result if they are in the intersection of the document sets resulting from each instance of the parameter. In the example below, only documents which have a popularity greater than 10 and have a section of 0 will match.

```
fq=popularity:[10 TO *]&fq=section:0
```

- Filter queries can involve complicated Boolean queries. The above example could also be written as a single `fq` with two mandatory clauses like so:

```
fq=+popularity:[10 TO *]+section:0
```

- The document sets from each filter query are cached independently. Thus, concerning the previous examples: use a single `fq` containing two mandatory clauses if those clauses appear together often, and use two separate `fq` parameters if they are relatively independent. (To learn about tuning cache sizes and making sure a filter cache actually exists, see [The Well-Configured Solr Instance](#).)
- As with all parameters: special characters in an URL need to be properly escaped and encoded as hex values. Online tools are available to help you with URL-encoding. For example: <http://meyerweb.com/eric/tools/dencoder/>.

The `fl` (Field List) Parameter

The `fl` parameter limits the information included in a query response to a specified list of fields. The fields need to have been indexed as stored for this parameter to work correctly.

The field list can be specified as a space-separated or comma-separated list of field names. The string "score" can be used to indicate that the score of each document for the particular query should be returned as a field. The wildcard character "*" selects all the stored fields in a document. You can also add psuedo-fields, functions and transformers to the field list request.

This table shows some basic examples of how to use `fl`:

Field List	Result
<code>id name price</code>	Return only the id, name, and price fields.
<code>id,name,price</code>	Return only the id, name, and price fields.
<code>id name, price</code>	Return only the id, name, and price fields.
<code>id score</code>	Return the id field and the score.
<code>*</code>	Return all the fields in each document. This is the default value of the <code>fl</code> parameter.
<code>* score</code>	Return all the fields in each document, along with each field's score.

Document Transformers

Transformers modify fields returned with the query response. Transformers must first be configured in `solrconfig.xml`. The sample `solrconfig.xml` has a few examples commented out which could be enabled, but others could be added. Then the transformers could be added to the query request and the response will be modified accordingly.

For example, if you have enabled a transformer called "elevated", you could mark all documents that have been elevated with the QueryElevationComponent. One way to do that is to make this entry in `solrconfig.xml`:

```
<transformer name="elevated"  
class="org.apache.solr.response.transform.EditorialMarkerFactory" />
```

Then, you would include `[elevated]` in the part of your request where you define the fields to return:

```
f1=id,title,[elevated]
```

Other common examples are to add "explain" information, add a constant "value" to all documents, or add the "shard" the document has been indexed on. For more information about transformers, see also <http://wiki.apache.org/solr/DocTransformers>.

Field Name Aliases

You can change the name a field is returned with by passing a parameter of `fieldName:displayName`. This will change the name of the field in the response to the `displayName`. For example:

```
f1=id,price:sale_price
```

The debug Parameter

In Solr 4, requesting debugging information with results has been simplified from a suite of related parameters to a single parameter that takes format information as arguments. The parameter is now simply `debug`, with the following arguments:

- `debug=true`: return debug information about the query only.
- `debug=query`: return debug information about the query only.
- `debug=timing`: return debug information about how long the query took to process.
- `debug=results`: return debug information about the results (also known as "explain")

The default behavior is not to include debugging information.

The explainOther Parameter

The `explainOther` parameter specifies a Lucene query in order to identify a set of documents. If this parameter is included and is set to a non-blank value, the query will return debugging information, along with the "explain info" of each document that matches the Lucene query, relative to the main query (which is specified by the `q` parameter). For example:

```
q=supervillians&debugQuery=on&explainOther=id:juggernaut
```

The query above allows you to examine the scoring explain info of the top matching documents, compare it to the explain info for documents matching `id:juggernaut`, and determine why the rankings are not as you expect.

The default value of this parameter is blank, which causes no extra "explain info" to be returned.

The `timeAllowed` Parameter

This parameter specifies the amount of time, in milliseconds, allowed for a search to complete. If this time expires before the search is complete, any partial results will be returned.

The `omitHeader` Parameter

This parameter may be set to either true or false.

If set to true, this parameter excludes the header from the returned results. The header contains information about the request, such as the time it took to complete. The default value for this parameter is false.

The `wt` Parameter

The `wt` parameter selects the Response Writer that Solr should use to format the query's response. For detailed descriptions of Response Writers, see [Response Writers](#).

The `cache=false` Parameter

Solr caches the results of all queries and filter queries by default. To disable result caching, set the `cache=false` parameter.

You can also use the `cost` option to control the order in which non-cached filter queries are evaluated. This allows you to order less expensive non-cached filters before expensive non-cached filters.

For very high cost filters, if `cache=false` and `cost>=100` and the query implements the `PostFilter` interface, a Collector will be requested from that query and used to filter documents after they have matched the main query and all other filter queries. There can be multiple post filters; they are also ordered by cost.

For example:

```
// normal function range query used as a filter, all matching documents generated up
front and cached
fq={!frange l=10 u=100}mul(popularity,price)

// function range query run in parallel with the main query like a traditional lucene
filter
fq={!frange l=10 u=100 cache=false}mul(popularity,price)

// function range query checked after each document that already matches the query and
all other filters.
    Good for really expensive function queries.
fq={!frange l=10 u=100 cache=false cost=100}mul(popularity,price)
```

The Standard Query Parser

Before Solr 1.3, the Standard Request Handler called the standard query parser as the default query parser. In versions since Solr 1.3, the Standard Request Handler calls the DisMax query parser as the default query parser. You can configure Solr to call the standard query parser instead, if you like.

The advantage of the standard query parser is that it enables users to specify very precise queries. The disadvantage is that it is less tolerant of syntax errors than the [DisMax](#) query parser. The DisMax query parser is designed to throw as few errors as possible.

Topics covered in this section:

- Standard Query Parser Parameters
- The Standard Query Parser's Response
- Specifying Terms for the Standard Query Parser
- Specifying Fields in a Query to the Standard Query Parser
- Boolean Operators Supported by the Standard Query Parser
- Grouping Terms to Form Sub-Queries
- Differences between Lucene Query Parser and the Solr Standard Query Parser
- Related Topics

Standard Query Parser Parameters

In addition to the [Common Query Parameters](#), [Faceting Parameters](#), [Highlighting Parameters](#), and [MoreLikeThis Parameters](#), the standard query parser supports the parameters described in the table below.

Parameter	Description
q	Defines a query using standard query syntax. This parameter is mandatory.
q.op	Specifies the default operator for query expressions, overriding the default operator specified in the <code>schema.xml</code> file. Possible values are "AND" or "OR".
df	Specifies a default field, overriding the definition of a default field in the <code>schema.xml</code> file.

Default parameter values are specified in `solrconfig.xml`, or overridden by query-time values in the request.

The Standard Query Parser's Response

By default, the response from the standard query parser contains one `<result>` block, which is unnamed. If the [debugQuery parameter](#) is used, then an additional `<lst>` block will be returned, using the name "debug". This will contain useful debugging info, including the original query string, the parsed query string, and explain info for each document in the `<result>` block. If the [explainOther parameter](#) is also used, then additional explain info will be provided for all the documents matching that query.

Sample Responses

This section presents examples of responses from the standard query parser.

The URL below submits a simple query and requests the XML Response Writer to use indentation to make the XML response more readable.

```
http://yourhost.tld:9999/solr/select?q=id:SP2514N&version=2.1&indent=1
```

Results:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<responseHeader><status>0</status><QTime>1</QTime></responseHeader>
<result numFound="1" start="0">
<doc>
<arr name="cat"><str>electronics</str><str>hard drive</str></arr>
<arr name="features"><str>7200RPM, 8MB cache, IDE Ultra ATA-133</str>
<str>NoiseGuard, SilentSeek technology, Fluid Dynamic Bearing (FDB)
motor</str></arr>
<str name="id">SP2514N</str>
<bool name="inStock">true</bool>
<str name="manu">Samsung Electronics Co. Ltd.</str>
<str name="name">Samsung SpinPoint P120 SP2514N - hard drive - 250 GB - ATA-133</str>
<int name="popularity">6</int>
<float name="price">92.0</float>
<str name="sku">SP2514N</str>
</doc>
</result>
</response>
```

Here's an example of a query with a limited field list.

```
http://yourhost.tld:9999/solr/select?q=id:SP2514N&version=2.1&indent=1&fl=id+name
```

Results:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<responseHeader><status>0</status><QTime>2</QTime></responseHeader>
<result numFound="1" start="0">
<doc>
<str name="id">SP2514N</str>
<str name="name">Samsung SpinPoint P120 SP2514N - hard drive - 250 GB - ATA-133</str>
</doc>
</result>
</response>
```

Specifying Terms for the Standard Query Parser

A query to the standard query parser is broken up into terms and operators. There are two types of terms: single terms and phrases.

- A single term is a single word such as "test" or "hello"
- A phrase is a group of words surrounded by double quotes such as "hello dolly"

Multiple terms can be combined together with Boolean operators to form more complex queries (as described below).

 It is important that the analyzer used for queries parses terms and phrases in a way that is consistent with the way the analyzer used for indexing parses terms and phrases; otherwise, searches may produce unexpected results.

Term Modifiers

Solr supports a variety of term modifiers that add flexibility or precision, as needed, to searches. These modifiers include wildcard characters, characters for making a search "fuzzy" or more general, and so on. The sections below describe these modifiers in detail.

Wildcard Searches

Solr's standard query parser supports single and multiple character wildcard searches within single terms. Wildcard characters can be applied to single terms, but not to search phrases.

Wildcard Search Type	Special Character	Example
Single character (matches a single character)	?	The search string <code>te?t</code> would match both test and text.
Multiple characters (matches zero or more sequential characters)	*	<p>The wildcard search: <code>tes*</code> would match test, testing, and tester.</p> <p>You can also use wildcard characters in the middle of a term. For example:</p> <p><code>te*t</code> would match test and text.</p> <p><code>*est</code> would match pest and test.</p>

 As of Solr 1.4, you can use a * or ? symbol as the first character of a search with the standard query parser.

Fuzzy Searches

Solr's standard query parser supports fuzzy searches based on the Levenshtein Distance or Edit Distance algorithm. Fuzzy searches discover terms that are similar to a specified term without necessarily being an exact match. To perform a fuzzy search, use the tilde ~ symbol at the end of a single-word term. For example, to search for a term similar in spelling to "roam," use the fuzzy search:

```
roam~
```

This search will match terms like foam and roams. It will also match the word "roam" itself.

An optional, additional parameter specifies the degree of similarity required for a match in a fuzzy search. The value must be between 0 and 1. When set closer to 1, the optional parameter causes only terms with a higher similarity to be matched. For example, the search below requires a high degree of similarity to the term "roam" in order for Solr to return a match:

```
roam~0.8
```

If this numerical parameter is omitted, Lucene performs the search as though the parameter were set to 0.5. The sample query above is not very scalable. Upon parsing this query will check the quasi-edit distance for every term in the index. As a result, this query is practical for only very small indexes.

 In many cases, stemming (reducing terms to a common stem) can produce similar effects to fuzzy searches and wildcard searches.

Proximity Searches

A proximity search looks for terms that are within a specific distance from one another.

To perform a proximity search, add the tilde character ~ and a numeric value to the end of a search phrase. For example, to search for a "apache" and "jakarta" within 10 words of each other in a document, use the search:

```
"jakarta apache"~10
```

The distance referred to here is the number of term movements needed to match the specified phrase. In the example above, if "apache" and "jakarta" were 10 spaces apart in a field, but "apache" appeared before "jakarta", more than 10 term movements would be required to move the terms together and position "apache" to the right of "jakarta" with a space in between.

Range Searches

A range search specifies a range of values for a field (a range with an upper bound and a lower bound). The query matches documents whose values for the specified field or fields fall within the range. Range queries can be inclusive or exclusive of the upper and lower bounds. Sorting is done lexicographically, except on numeric fields. For example, the range query below matches all documents whose `mod_date` field has a value between 20020101 and 20030101, inclusive.

```
mod_date:[20020101 TO 20030101]
```

Range queries are not limited to date fields or even numerical fields. You could also use range queries with non-date fields:

```
title:{Aida TO Carmen}
```

This will find all documents whose titles are between Aida and Carmen, but not including Aida and Carmen.

The brackets around a query determine its inclusiveness.

- Square brackets [] denote an inclusive range query that matches values including the upper and lower bound.
- Curly brackets { } denote an exclusive range query that matches values between the upper and lower bounds, but excluding the upper and lower bounds themselves.

Boosting a Term with ^

Lucene/Solr provides the relevance level of matching documents based on the terms found. To boost a term use the caret symbol ^ with a boost factor (a number) at the end of the term you are searching. The higher the boost factor, the more relevant the term will be.

Boosting allows you to control the relevance of a document by boosting its term. For example, if you are searching for

"jakarta apache" and you want the term "jakarta" to be more relevant, you can boost it by adding the ^ symbol along with the boost factor immediately after the term. For example, you could type:

```
jakarta^4 apache
```

This will make documents with the term jakarta appear more relevant. You can also boost Phrase Terms as in the example:

```
"jakarta apache" ^4 "Apache Lucene"
```

By default, the boost factor is 1. Although the boost factor must be positive, it can be less than 1 (for example, it could be 0.2).

Specifying Fields in a Query to the Standard Query Parser

Data indexed in Solr is organized in fields, which are defined in the Solr schema.xml file. Searches can take advantage of fields to add precision to queries. For example, you can search for a term only in a specific field, such as a title field.

The schema.xml file defines one field as a default field. If you do not specify a field in a query, Solr searches only the default field. Alternatively, you can specify a different field or a combination of fields in a query.

To specify a field, type the field name followed by a colon ":" and then the term you are searching for within the field.

For example, suppose an index contains two fields, title and text, and that text is the default field. If you want to find a document called "The Right Way" which contains the text "don't go this way," you could include either of the following terms in your search query:

```
title:"The Right Way" AND text:go
```

```
title:"Do it right" AND go
```

Since text is the default field, the field indicator is not required; hence the second query above omits it.

The field is only valid for the term that it directly precedes, so the query title:Do it right will find only "Do" in the title field. It will find "it" and "right" in the default field (in this case the text field).

Boolean Operators Supported by the Standard Query Parser

Boolean operators allow you to apply Boolean logic to queries, requiring the presence or absence of specific terms or conditions in fields in order to match documents. The table below summarizes the Boolean operators supported by the standard query parser.

Boolean Operator	Alternative Symbol	Description
AND	&&	Requires both terms on either side of the Boolean operator to be present for a match.
NOT	!	Requires that the following term not be present.

OR		Requires that either term (or both terms) be present for a match.
	+	Requires that the following term be present.
	-	Prohibits the following term (that is, matches on fields or documents that do not include that term). The - operator is functional similar to the Boolean operator !. Because it's used by popular search engines such as Google, it may be more familiar to some user communities.

Boolean operators allow terms to be combined through logic operators. Lucene supports AND, "+", OR, NOT and "-" as Boolean operators.

 When specifying Boolean operators with keywords such as AND or NOT, the keywords must appear in all uppercase.

- ① The standard query parser supports all the Boolean operators listed in the table above. The DisMax query parser supports only + and -.

The OR operator is the default conjunction operator. This means that if there is no Boolean operator between two terms, the OR operator is used. The OR operator links two terms and finds a matching document if either of the terms exist in a document. This is equivalent to a union using sets. The symbol || can be used in place of the word OR.

In the schema.xml file, you can specify which symbols can take the place of Boolean operators such as OR. To search for documents that contain either "jakarta apache" or just "jakarta," use the query:

"jakarta apache" jakarta

or

"jakarta apache" OR jakarta

The Boolean Operator +

The + symbol (also known as the "required" operator) requires that the term after the + symbol exist somewhere in a field in at least one document in order for the query to return a match.

For example, to search for documents that must contain "jakarta" and that may or may not contain "lucene," use the following query:

+jakarta lucene

-  This operator is supported by both the standard query parser and the DisMax query parser.

The Boolean Operator AND (&&)

The AND operator matches documents where both terms exist anywhere in the text of a single document. This is equivalent to an intersection using sets. The symbol && can be used in place of the word AND.

To search for documents that contain "jakarta apache" and "Apache Lucene," use either of the following queries:

```
"jakarta apache" AND "Apache Lucene"
```

```
"jakarta apache" && "Apache Lucene"
```

The Boolean Operator NOT (!)

The NOT operator excludes documents that contain the term after NOT. This is equivalent to a difference using sets. The symbol ! can be used in place of the word NOT.

The following queries search for documents that contain the phrase "jakarta apache" but do not contain the phrase "Apache Lucene":

```
"jakarta apache" NOT "Apache Lucene"
```

```
"jakarta apache" ! "Apache Lucene"
```

The Boolean Operator -

The - symbol or "prohibit" operator excludes documents that contain the term after the - symbol.

For example, to search for documents that contain "jakarta apache" but not "Apache Lucene," use the following query:

```
"jakarta apache" -"Apache Lucene"
```

Escaping Special Characters

Solr gives the following characters special meaning when they appear in a query:

```
+ - && || ! ( ) { } [ ] ^ " ~ * ? : /
```

To make Solr interpret any of these characters literally, rather as a special character, precede the character with a backslash character \. For example, to search for (1+1):2 without having Solr interpret the plus sign and parentheses as special characters for formulating a sub-query with two terms, escape the characters by preceding each one with a backslash:

```
\(1\+1\)\:2
```

Grouping Terms to Form Sub-Queries

Lucene/Solr supports using parentheses to group clauses to form sub-queries. This can be very useful if you want to control the Boolean logic for a query.

The query below searches for either "jakarta" or "apache" and "website":

```
(jakarta OR apache) AND website
```

This adds precision to the query, requiring that the term "website" exist, along with either term "jakarta" and "apache."

Grouping Clauses within a Field

To apply two or more Boolean operators to a single field in a search, group the Boolean clauses within parentheses. For example, the query below searches for a title field that contains both the word "return" and the phrase "pink panther":

```
title:(+return +"pink panther")
```

Differences between Lucene Query Parser and the Solr Standard Query Parser

Solr's standard query parser differs from the Lucene Query Parser in the following ways:

- A * may be used for either or both endpoints to specify an open-ended range query
 - field:[* TO 100] finds all field values less than or equal to 100
 - field:[100 TO *] finds all field values greater than or equal to 100
 - field:[* TO *] matches all documents with the field
- Pure negative queries (all clauses prohibited) are allowed (only as a top-level clause)
 - -inStock:false finds all field values where inStock is not false
 - -field:[* TO *] finds all documents without a value for field
- A hook into FunctionQuery syntax. You'll need to use quotes to encapsulate the function if it includes parentheses, as shown in the second example below:
 - val:myfield
 - val:"recip(rord(myfield),1,2,3)"

- Support for any type of query parser. Prior to Solr 4.1, the "magic" field `_query_` needed to be used to nest another query parser. However, with Solr 4.1, other query parsers can be used directly using the local parameters syntax.
 - `{}{
}}`
- Range queries ("[a TO z]"), prefix queries ("a*"), and wildcard queries ("a*b") are constant-scoring (all matching documents get an equal score). The scoring factors TF, IDF, index boost, and "coord" are not used. There is no limitation on the number of terms that match (as there was in past versions of Lucene).

Specifying Dates and Times

If you use the Solr `DateField` type, any queries on those fields (typically range queries) should use the `TrieDate` Field. Here are some examples of valid parameters using syntax appropriate for the `DateField` type:

- `timestamp:[*TO NOW]`
- `createdate:[1976-03-06T23:59:59.999Z TO *]`
- `createdate:[1995-12-31T23:59:59.999Z TO 2007-03-06T00:00:00Z]`
- `pubdate:[NOW-1YEAR/DAY TO NOW/DAY+1DAY]`
- `createdate:[1976-03-06T23:59:59.999Z TO 1976-03-06T23:59:59.999Z+1YEAR]`
- `createdate:[1976-03-06T23:59:59.999Z/YEAR TO 1976-03-06T23:59:59.999Z]`

Related Topics

- [Local Parameters in Queries](#)
- [Other Parsers](#)

The DisMax Query Parser

The DisMax query parser is designed to process simple phrases (without complex syntax) entered by users and to search for individual terms across several fields using different weighting (boosts) based on the significance of each field. Additional options enable users to influence the score based on rules specific to each use case (independent of user input).

In general, the DisMax query parser's interface is more like that of Google than the interface of the 'standard' Solr request handler. This similarity makes DisMax the appropriate query parser for many consumer applications. It accepts a simple syntax, and it rarely produces error messages.

The DisMax query parser supports an extremely simplified subset of the Lucene QueryParser syntax. As in Lucene, quotes can be used to group phrases, and +/- can be used to denote mandatory and optional clauses. All other Lucene query parser special characters (except AND and OR) are escaped to simplify the user experience. The DisMax query parser takes responsibility for building a good query from the user's input using Boolean clauses containing DisMax queries across fields and boosts specified by the user. It also lets the Solr administrator provide additional boosting queries, boosting functions, and filtering queries to artificially affect the outcome of all searches. These options can all be specified as default parameters for the handler in the `solrconfig.xml` file or overridden in the Solr query URL.

Interested in the technical concept behind the DisMax name? DisMax stands for Maximum Disjunction. Here's a definition of a Maximum Disjunction or "DisMax" query:

A query that generates the union of documents produced by its subqueries, and that scores each document with the maximum score for that document as produced by any subquery, plus a tie breaking increment for any additional matching subqueries.

Whether or not you remember this explanation, do remember that the DisMax request handler was primarily designed to be easy to use and to accept almost any input without returning an error.

DisMax Parameters

In addition to the common request parameter, highlighting parameters, and simple facet parameters, the DisMax query parser supports the parameters described below. Like the standard query parser, the DisMax query parser allows default parameter values to be specified in `solrconfig.xml`, or overridden by query-time values in the request.

Parameter	Description
<code>q</code>	Defines the raw input strings for the query.
<code>q.alt</code>	Calls the standard query parser and defines query input strings, when the <code>q</code> parameter is not used.

<code>qf</code>	Query Fields: specifies the fields in the index on which to perform the query. If absent, defaults to <code>df</code> .
<code>mm</code>	Minimum "Should" Match: specifies a minimum number of fields that must match in a query. If no 'mm' parameter is specified in the query, or as a default in <code>solrconfig.xml</code> , the effective value of the <code>q.op</code> parameter (either in the query, as a default in <code>solrconfig.xml</code> , or from the 'defaultOperator' option in <code>schema.xml</code>) is used to influence the behavior. If <code>q.op</code> is effectively AND'ed, then <code>mm=100%</code> ; if <code>q.op</code> is OR'ed, then <code>mm=1</code> . Users who want to force the legacy behavior should set a default value for the 'mm' parameter in their <code>solrconfig.xml</code> file. Users should add this as a configured default for their request handlers. This parameter tolerates miscellaneous white spaces in expressions (e.g., " 3 < -25% 10 < -3\n", "\n-25%\n", "\n3\n").
<code>pf</code>	Phrase Fields: boosts the score of documents in cases where all of the terms in the <code>q</code> parameter appear in close proximity.
<code>ps</code>	Phrase Slop: specifies the number of positions two terms can be apart in order to match the specified phrase.
<code>qs</code>	Query Phrase Slop: specifies the number of positions two terms can be apart in order to match the specified phrase. Used specifically with the <code>qf</code> parameter.
<code>tie</code>	Tie Breaker: specifies a float value (which should be something much less than 1) to use as tiebreaker in DisMax queries.
<code>bq</code>	Boost Query: specifies a factor by which a term or phrase should be "boosted" in importance when considering a match.
<code>bf</code>	Boost Functions: specifies functions to be applied to boosts. (See for details about function queries.)

The sections below explain these parameters in detail.

The `q` Parameter

The `q` parameter defines the main "query" constituting the essence of the search. The parameter supports raw input strings provided by users with no special escaping. The + and - characters are treated as "mandatory" and "prohibited" modifiers for terms. Text wrapped in balanced quote characters (for example, "San Jose") is treated as a phrase. Any query containing an odd number of quote characters is evaluated as if there were no quote characters at all.



The `q` parameter does not support wildcard characters such as *.

The `q.alt` Parameter

If specified, the `q.alt` parameter defines a query (which by default will be parsed using standard query parsing syntax) when the main `q` parameter is not specified or is blank. The `q.alt` parameter comes in handy when you need something like a query to match all documents (don't forget `&rows=0` for that one!) in order to get collection-wise facetting counts.

The `qf` (Query Fields) Parameter

The `qf` parameter introduces a list of fields, each of which is assigned a boost factor to increase or decrease that particular field's importance in the query. For example, the query below:

```
qf="fieldOne^2.3 fieldTwo fieldThree^0.4"
```

assigns `fieldOne` a boost of 2.3, leaves `fieldTwo` with the default boost (because no boost factor is specified), and `fieldThree` a boost of 0.4. These boost factors make matches in `fieldOne` much more significant than matches in `fieldTwo`, which in turn are much more significant than matches in `fieldThree`.

The `mm` (Minimum Should Match) Parameter

When processing queries, Lucene/Solr recognizes three types of clauses: mandatory, prohibited, and "optional" (also known as "should" clauses). By default, all words or phrases specified in the `q` parameter are treated as "optional" clauses unless they are preceded by a "+" or a "-". When dealing with these "optional" clauses, the `mm` parameter makes it possible to say that a certain minimum number of those clauses must match. The DisMax query parser offers great flexibility in how the minimum number can be specified.

The table below explains the various ways that `mm` values can be specified.

Syntax	Example	Description
Positive integer	3	Defines the minimum number of clauses that must match, regardless of how many clauses there are in total.
Negative integer	-2	Sets the minimum number of matching clauses to the total number of optional clauses, minus this value.
Percentage	75%	Sets the minimum number of matching clauses to this percentage of the total number of optional clauses. The number computed from the percentage is rounded down and used as the minimum.
Negative percentage	-25%	Indicates that this percent of the total number of optional clauses can be missing. The number computed from the percentage is rounded down, before being subtracted from the total to determine the minimum number.

An expression beginning with a positive integer followed by a > or < sign and another value	3<90%	Defines a conditional expression indicating that if the number of optional clauses is equal to (or less than) the integer, they are all required, but if it's greater than the integer, the specification applies. In this example: if there are 1 to 3 clauses they are all required, but for 4 or more clauses only 90% are required.
Multiple conditional expressions involving > or < signs	2<-25% 9<-3	Defines multiple conditions, each one being valid only for numbers greater than the one before it. In the example at left, if there are 1 or 2 clauses, then both are required. If there are 3-9 clauses all but 25% are required. If there are more than 9 clauses, all but three are required.

When specifying `mm` values, keep in mind the following:

- When dealing with percentages, negative values can be used to get different behavior in edge cases. 75% and -25% mean the same thing when dealing with 4 clauses, but when dealing with 5 clauses 75% means 3 are required, but -25% means 4 are required.
- If the calculations based on the parameter arguments determine that no optional clauses are needed, the usual rules about Boolean queries still apply at search time. (That is, a Boolean query containing no required clauses must still match at least one optional clause).
- No matter what number the calculation arrives at, Solr will never use a value greater than the number of optional clauses, or a value less than 1. (In other words, no matter how low or how high the calculated result, the minimum number of required matches will never be less than 1 or greater than the number of clauses.)

The default value of `mm` is 100% (meaning that all clauses must match).

The `pf` (Phrase Fields) Parameter

Once the list of matching documents has been identified using the `fq` and `qf` parameters, the `pf` parameter can be used to "boost" the score of documents in cases where all of the terms in the `q` parameter appear in close proximity.

The format is the same as that used by the `qf` parameter: a list of fields and "boosts" to associate with each of them when making phrase queries out of the entire `q` parameter.

The `ps` (Phrase Slop) Parameter

The `ps` parameter specifies the amount of "phrase slop" to apply to queries specified with the `pf` parameter. Phrase slop is the number of positions one token needs to be moved in relation to another token in order to match a phrase specified in a query.

The `qs` (Query Phrase Slop) Parameter

The `qs` parameter specifies the amount of slop permitted on phrase queries explicitly included in the user's query string with the `qf` parameter. As explained above, slop refers to the number of positions one token needs to be moved in relation to another token in order to match a phrase specified in a query.

The `tie` (Tie Breaker) Parameter

The `tie` parameter specifies a float value (which should be something much less than 1) to use as tiebreaker in DisMax queries.

When a term from the user's input is tested against multiple fields, more than one field may match. If so, each field will generate a different score based on how common that word is in that field (for each document relative to all other documents). The `tie` parameter lets you control how much the final score of the query will be influenced by the scores of the lower scoring fields compared to the highest scoring field.

A value of "0.0" makes the query a pure "disjunction max query": that is, only the maximum scoring subquery contributes to the final score. A value of "1.0" makes the query a pure "disjunction sum query" where it doesn't matter what the maximum scoring sub query is, because the final score will be the sum of the subquery scores. Typically a low value, such as 0.1, is useful.

The `bq` (Boost Query) Parameter

The `bq` parameter specifies a raw query string (expressed in Solr query syntax) that will be included with the user's query to influence the score. For example, if you wanted to add a relevancy boost for recent documents:

```
q=cheese bq=date\[NOW/DAY-1YEAR TO NOW/DAY\]
```

You can specify multiple `bq` parameters. If you want your query to be parsed as separate clauses with separate boosts, use multiple `bq` parameters.

The `bf` (Boost Functions) Parameter

The `bf` parameter specifies functions (with optional boosts) that will be included in the user's query to influence the score. Any function supported natively by Solr can be used, along with a boost value. For example:

```
recip(rord(myfield),1,2,3)^1.5
```

Specifying functions with the `bf` parameter is just shorthand for using the `val: "...function..."` syntax in a `bq` parameter.

For example, if you want to show the most recent documents first, use

```
recip(rord(creationDate),1,1000,1000)
```

Examples of Queries Submitted to the DisMax Query Parser

Normal results for the word "video" using the StandardRequestHandler with the default search field:

```
http://localhost:8983/solr/select/?q=video&fl=name+score
```

The "dismax" handler is configured to search across the text, features, name, sku, id, manu, and cat fields all with varying boosts designed to ensure that "better" matches appear first, specifically: documents which match on the name and cat fields get higher scores.

```
http://localhost:8983/solr/select/?defType=dismax&q=video
```

Note that this instance is also configured with a default field list, which can be overridden in the URL.

```
http://localhost:8983/solr/select/?defType=dismax&q=video&fl=*,score
```

You can also override which fields are searched on and how much boost each field gets.

```
http://localhost:8983/solr/select/?defType=dismax&q=video&qt=features^20.0+text^0.3
```

You can boost results that have a field that matches a specific value.

```
http://localhost:8983/solr/select/?defType=dismax&q=video&bq=cat:electronics^5.0
```

Another instance of the handler is registered using the qt "instock" and has slightly different configuration options, notably: a filter for (you guessed it) inStock:true).

```
http://localhost:8983/solr/select/?defType=dismax&q=video&fl=name,score,inStock
```

```
http://localhost:8983/solr/select/?defType=dismax&q=video&qt=instock&fl=name,score,i
```

One of the other really cool features in this handler is robust support for specifying the "BooleanQuery.minimumNumberShouldMatch" you want to be used based on how many terms are in your user's query. These allows flexibility for typos and partial matches. For the dismax handler, one and two word queries require that all of the optional clauses match, but for three to five word queries one missing word is allowed.

```
http://localhost:8983/solr/select/?defType=dismax&q=belkin+ipod
```

```
http://localhost:8983/solr/select/?defType=dismax&q=belkin+ipod+gibberish
```

```
http://localhost:8983/solr/select/?defType=dismax&q=belkin+ipod+apple
```

Just like the StandardRequestHandler, it supports the debugQuery option to viewing the parsed query, and the score explanations for each document.

```
http://localhost:8983/solr/select/?defType=dismax&q=belkin+ipod+gibberish&debugQuery
```

```
http://localhost:8983/solr/select/?defType=dismax&q=video+card&debugQuery=true
```

The Extended DisMax Query Parser

The Extended DisMax (eDisMax) query parser is an improved version of the [DisMax query parser](#). In addition to supporting all the DisMax query parser parameters, Extended Dismax:

- supports the full Lucene query parser syntax.
- supports queries such as AND, OR, NOT, -, and +.
- treats "and" and "or" as "AND" and "OR" in Lucene syntax mode.
- respects the 'magic field' names `_val_` and `_query_`. These are not a real fields in `schema.xml`, but if used it helps do special things (like a function query in the case of `_val_` or a nested query in the case of `_query_`). If `_val_` is used in a term or phrase query, the value is parsed as a function.
- includes improved smart partial escaping in the case of syntax errors; fielded queries, +/-, and phrase queries are still supported in this mode.
- improves proximity boosting by using word shingles; you do not need the query to match all words in the document before proximity boosting is applied.
- includes advanced stopword handling: stopwords are not required in the mandatory part of the query but are still used in the proximity boosting part. If a query consists of all stopwords, such as "to be or not to be", then all words are required.
- includes improved boost function: in Extended DisMax, the `boost` function is a multiplier rather than an addend, improving your boost results; the additive boost functions of DisMax (`bf` and `bq`) are also supported.
- supports pure negative nested queries: queries such as `+foo (-foo)` will match all documents.
- lets you specify which fields the end user is allowed to query, and to disallow direct fielded searches.

 The Extended DisMax query parser is still under active development, in fact many changes were introduced for Solr 4. However, many organizations are already using it in production with great success.

Extended DisMax Parameters

In addition to all the [DisMax parameters](#), Extended DisMax includes these query parameters:

The `boost` Parameter

A multivalued list of strings parsed as queries with scores multiplied by the score from the main query for all matching documents. This parameter is shorthand for wrapping the query produced by eDisMax using the `BoostQParserPlugin`

The `lowercaseOperators` Parameter

A Boolean parameter indicating if lowercase "and" and "or" should be treated the same as operators "AND" and "OR".

The `ps` Parameter

Default amount of slop on phrase queries built with `pf`, `pf2` and/or `pf3` fields (affects boosting).

The `pf2` Parameter

A multivalued list of fields with optional weights, based on pairs of word shingles.

The `ps2` Parameter

Default amount of slop on phrase queries built with `pf`, `pf2` and/or `pf3` fields (affects boosting).

New with Solr 4, it is similar to `ps` but sets default slop factor for `pf2`. If not specified, `ps` is used.

The `pf3` Parameter

A multivalued list of fields with optional weights, based on triplets of word shingles. Similar to `pf`, except that instead of building a phrase per field out of all the words in the input, it builds a set of phrases for each field out of each triplet of word shingles.

The `ps3` Parameter

New with Solr 4. As with `ps` but sets default slop factor for `pf3`. If not specified, `ps` will be used.

The `stopwords` Parameter

A Boolean parameter indicating if the `StopFilterFactory` configured in the query analyzer should be respected when parsing the query: if it is false, then the `StopFilterFactory` in the query analyzer is ignored.

The `uf` Parameter

Specifies which schema fields the end user is allowed to explicitly query. This parameter supports wildcards. The default is to allow all fields, equivalent to `&uf=*`. To allow only title field, use `&uf=title`. To allow title and all fields ending with `_s`, use `&uf=title,*_s`. To allow all fields except title, use `&uf=-title`. To disallow all fielded searches, use `&uf=-*`.

Examples of Queries Submitted to the Extended DisMax Query Parser

Boost the result of the query term "hello" based on the document's popularity:

```
[http://localhost:8983/solr/select/?defType=edismax&q=hello&pf=text&qf=text&boost=populari
```

Search for iPods OR video:

```
[http://localhost:8983/solr/select/?defType=edismax&q=iPod OR video]
```

Search across multiple fields, specifying (via boosts) how important each field is relative each other:

```
[http://localhost:8983/solr/select/?q=video&defType=edismax&qf=features^20.0+text^0.3]
```

You can boost results that have a field that matches a specific value:

```
[http://localhost:8983/solr/select/?q=video&defType=edismax&qf=features^20.0+text^0.3&bq=c]
```

Using the "mm" param, 1 and 2 word queries require that all of the optional clauses match, but for queries with three or more clauses one missing clause is allowed:

```
[http://localhost:8983/solr/select/?q=belkin+ipod&defType=edismax&mm=2]
[http://localhost:8983/solr/select/?q=belkin+ipod+gibberish&defType=edismax&mm=2]
[http://localhost:8983/solr/select/?q=belkin+ipod+apple&defType=edismax&mm=2]
```

Using negative boost

Negative query boosts have been supported at the "Query" object level for a long time (resulting in negative scores for matching documents). Now the QueryParsers have been updated to handle this too.

Using 'slop'

Dismax and Edismax can run queries against all query fields, and also run a query in the form of a phrase against the phrase fields. (This will work only for boosting documents, not actually for matching.) However, that phrase query can have a 'slop,' which is the distance between the terms of the query while still considering it a phrase match. For example:

```
q=foo bar
qf=field1^5 field2^10
pf=field1^50 field2^20
defType=dismax
```

With these parameters, the Dismax Query Parser generates a query that looks something like this:

```
(+(field1:foo^5 OR field2:bar^10) AND (field1:bar^5 OR field2:bar^10))
```

But it also generates another query that will only be used for boosting results:

```
field1:"foo bar"^50 OR field2:"foo bar"^20
```

Thus, any document that has the terms "foo" and "bar" will match; however if some of those documents have both of the terms as a phrase, it will score much higher because it's more relevant.

If you add the parameter `ps` (phrase slop), the second query will instead be:

```
ps=10 field1:"foo bar"~10^50 OR field2:"foo bar"~10^20
```

This means that if the terms "foo" and "bar" appear in the document with less than 10 terms between each other, the phrase will match. For example the doc that says:

```
*Foo* term1 term2 term3 *bar*
```

will match the phrase query.

How does one use phrase slop? Usually it is configured in the request handler (in `solrconfig`).

With query slop (`qs`) the concept is similar, but it applies to explicit phrase queries from the user. For example, if you want to search for a name, you could enter:

```
q="Hans Anderson"
```

A document that contains "Hans Anderson" will match, but a document that contains the middle name "Christian" or where the name is written with the last name first ("Anderson, Hans") won't. For those cases one could configure the query field `qs`, so that even if the user searches for an explicit phrase query, a slop is applied.

Finally, `edismax` contains not only a phrase fields (`pf`) parameters, but also phrase and query fields 2 and 3. You can use those fields for setting different fields or boosts. Each of those can use a different phrase slop.

Using the 'magic fields' `_val_` and `_query_`

If the 'magic field' name `_val_` is used in a term or phrase query, the value is parsed as a function.

The Solr Query Parser's use of `_val_` and `_query_` differs from the Lucene Query Parser in the following ways:

- If the magic field name `_val_` is used in a term or phrase query, the value is parsed as a function.
- It provides a hook into `FunctionQuery` syntax. Quotes are necessary to encapsulate the function when it includes parentheses. For example:

```
_val_:myfield  
_val_:"recip(rord(myfield),1,2,3)"
```

- The Solr Query Parser offers nested query support for any type of query parser (via QParserPlugin). Quotes are often necessary to encapsulate the nested query if it contains reserved characters. For example:

```
{ {_query_:"{\!dismax;qf=myfield}how;now;brown;cow" } }
```

Although not technically a syntax difference, note that if you use the Solr [DateField](#) type, any queries on those fields (typically range queries) should use either the Complete ISO 8601 Date syntax that field supports, or the [DateMath Syntax](#) to get relative dates. For example:

```
timestamp:[ * TO NOW]  
createdate:[1976-03-06T23:59:59.999Z TO *]  
createdate:[1995-12-31T23:59:59.999Z TO 2007-03-06T00:00:00Z]  
pubdate:[NOW-1YEAR/DAY TO NOW/DAY+1DAY]  
createdate:[1976-03-06T23:59:59.999Z TO 1976-03-06T23:59:59.999Z+1YEAR]  
createdate:[1976-03-06T23:59:59.999Z/YEAR TO 1976-03-06T23:59:59.999Z]
```

 TO must be uppercase, or Solr will report a 'Range Group' error.

Local Parameters in Queries

Local parameters are arguments in a Solr request that are specific to a query parameter. Local parameters provide a way to add meta-data to certain argument types such as query strings. (In Solr documentation, local parameters are sometimes referred to as LocalParams.)

Local parameters are specified as prefixes to arguments. Take the following query argument, for example:

```
q=solr rocks
```

We can prefix this query string with local parameters to provide more information to the Standard Query Parser. For example, we can change the default operator type to "AND" and the default field to "title":

```
q={!q.op=AND df=title}solr rocks
```

These local parameters would change the query to require a match on both "solr" and "rocks" while searching the "title" field by default.

Basic Syntax of Local Parameters

To specify a local parameter, insert the following before the argument to be modified:

- Begin with {!
- Insert any number of key=value pairs separated by white space
- End with } and immediately follow with the query argument

You may specify only one local parameters prefix per argument. Values in the key-value pairs may be quoted via single or double quotes, and backslash escaping works within quoted strings.

Query Type Short Form

If a local parameter value appears without a name, it is given the implicit name of "type". This allows short-form representation for the type of query parser to use when parsing a query string. Thus

```
q={!dismax qf=myfield}solr rocks
```

is equivalent to:

```
q={!type=dismax qf=myfield}solr rocks
```

Specifying the Parameter Value with the 'v' Key

A special key of `v` within local parameters is an alternate way to specify the value of that parameter.

```
q={!dismax qf=myfield}solr rocks
```

is equivalent to

```
q={!type=dismax qf=myfield v='solr rocks'}
```

Parameter Dereferencing

Parameter dereferencing or indirection lets you use the value of another argument rather than specifying it directly. This can be used to simplify queries, decouple user input from query parameters, or decouple front-end GUI parameters from defaults set in `solrconfig.xml`.

```
q={!dismax qf=myfield}solr rocks
```

is equivalent to:

```
q={!type=dismax qf=myfield v=$qq}&qq=solr rocks
```

Other Parsers

In addition to the main query parsers discussed earlier, there are several other query parsers that can be used instead of or in conjunction with the main parsers for specific purposes. This section details the other parsers, and gives examples for how they might be used.

Many of these parsers are expressed the same way as [Local Parameters in Queries](#).

Query parsers discussed in this section:

- [Boost Query Parser](#)
- [Field Query Parser](#)
- [Function Query Parser](#)
- [Function Range Query Parser](#)
- [Join Query Parser](#)
- [Lucene Query Parser](#)
- [Nested Query Parser](#)
- [Old Lucene Query Parser](#)
- [Prefix Query Parser](#)
- [Raw Query Parser](#)
- [Spatial Filter Query Parser](#)
- [Surround Query Parser](#)
- [Term Query Parser](#)

Boost Query Parser

`BoostQParser` extends the `QParserPlugin` and creates a boosted query from the input value. The main value is the query to be boosted. Parameter `b` is the function query to use as the boost. The query to be boosted may be of any type.

Examples:

Creates a query "foo" which is boosted (scores are multiplied) by the function query `log(popularity)`:

```
!boost b=log(popularity) foo
```

Creates a query "foo" which is boosted by the date boosting function referenced in `ReciprocalFloatFunction`:

```
!boost b=recip(ms(NOW,mydatefield),3.16e-11,1,1) foo
```

Field Query Parser

The `FieldQParser` extends the `QParserPlugin` and creates a field query from the input value, applying text analysis and constructing a phrase query if appropriate. The parameter `f` is the field to be queried.

Example:

```
{!field f=myfield}Foo Bar
```

This example creates a phrase query with "foo" followed by "bar" (assuming the analyzer for myfield is a text field with an analyzer that splits on whitespace and lowercase terms). This is generally equivalent to the Lucene query parser expression `myfield: "Foo Bar"`.

Function Query Parser

The `FunctionQParser` extends the `QParserPlugin` and creates a function query from the input value. This is only one way to use function queries in Solr; for another, more integrated, approach, see the section on [Function Queries](#).

Example:

```
{!func}log(foo)
```

Function Range Query Parser

The `FunctionRangeQParser` extends the `QParserPlugin` and creates a range query over a function. This is also referred to as `frange`, as seen in the examples below.

Other parameters:

Parameter	Description
l	The lower bound, optional
u	The upper bound, optional
incl	Include the lower bound: true/false, optional, default=true
incu	Include the upper bound: true/false, optional, default=true

Examples:

```
{!frange l=1000 u=50000}myfield
```

```
fq={!frange l=0 u=2.2} sum(user_ranking,editor_ranking)
```

Both of these examples are restricting the results by a range of values found in a declared field or a function query. In the second example, we're doing a sum calculation, and then defining only values between 0 and 2.2 should be returned to the user.

For more information about range queries over functions, see Yonik Seeley's introductory blog post [Ranges over Functions in Solr 1.4](#), hosted at SearchHub.org.

Join Query Parser

`JoinQParser` extends the `QParserPlugin`. It allows normalizing relationships between documents with a join operation. This is different from in concept of a join in a relational database because no information is being truly joined. An appropriate SQL analogy would be an "inner query".

Examples:

Find all products containing the word "ipod", join them against manufacturer docs and return the list of manufacturers:

```
{!join+from=manu_id_s+to=id}ipod
```

Find all manufacturer docs named "belkin", join them against product docs, and filter the list to only products with a price less than \$12:

```
{!join+from=id+to=manu_id_s}compName_s:Belkin&fq=price:[*+TO+12]
```

For more information about join queries, see the Solr Wiki page on [Joins](#). Erick Erickson has also written a blog post about join performance called [Solr and Joins](#), hosted by SearchHub.org.

Lucene Query Parser

The `LuceneQParser` extends the `QParserPlugin` by parsing Solr's variant on the Lucene `QueryParser` syntax. This is effectively the same query parser that is used in Lucene. It uses the operators `q.op`, the default operator ("OR" or "AND") and `df`, the default field name.

Example:

```
{!lucene q.op=AND df=text sort='price asc'}myfield:foo +bar -baz
```

For more information about the syntax for the Lucene Query Parser, see the [Lucene javadocs](#).

Nested Query Parser

The `NestedParser` extends the `QParserPlugin` and creates a nested query, with the ability for that query to redefine its type via local parameters. This is useful in specifying defaults in configuration and letting clients indirectly reference them.

Example:

```
{!query defType=func v=$q1}
```

If the `q1` parameter is `price`, then the query would be a function query on the `price` field. If the `q1` parameter is `{!lucene}inStock:true}}` then a term query is created from the Lucene syntax string that matches documents with `inStock=true`. These parameters would be defined in `solrconfig.xml`, in the `defaults` section:

```
<lst name="defaults">
  <str name="q1">{!lucene}inStock:true</str>
</lst>
```

For more information about the possibilities of nested queries, see Yonik Seeley's blog post [Nested Queries in Solr](#), hosted by SearchHub.org.

Old Lucene Query Parser

`OldLuceneQParser` extends the `QParserPlugin` by parsing Solr's variant of Lucene's `QueryParser` syntax, including the deprecated sort specification after the query.

Example:

```
{!lucenePlusSort} myfield:foo +bar -baz;price asc
```

Prefix Query Parser

`PrefixQParser` extends the `QParserPlugin` by creating a prefix query from the input value. Currently no analysis or value transformation is done to create this prefix query. The parameter is `f`, the field. The string after the prefix declaration is treated as a wildcard query.

Example:

```
{!prefix f=myfield}foo
```

This would be generally equivalent to the Lucene query parser expression `myfield:foo*`.

Raw Query Parser

`RawQParser` extends the `QParserPlugin` by creating a term query from the input value without any text analysis or transformation. This is useful in debugging, or when raw terms are returned from the terms component (this is not the default). The only parameter is `f`, which defines the field to search.

Example:

```
{!raw f=myfield}Foo Bar
```

This example constructs the query: `TermQuery(Term("myfield" , "Foo Bar"))`.

For easy filter construction to drill down in facetting, the [TermQParserPlugin](#) is recommended. For full analysis on all fields, including text fields, you may want to use the [FieldQParserPlugin](#).

Spatial Filter Query Parser

`SpatialFilterQParser` extends the `QParserPlugin` by creating a spatial Filter based on the type of spatial point used. The field must implement [SpatialQueryable](#). All units are in Kilometers.

This query parser takes the following parameters:

Parameter	Description
<code>sfeld</code>	The field on which to filter. Required.
<code>pt</code>	The point to use as a reference. Must match the dimension of the field. Required.
<code>d</code>	The distance in km. Required.

The distance measure used currently depends on the `FieldType`. `LatLonType` defaults to using `haversine`, `PointType` defaults to Euclidean (2-norm).

This example shows the syntax:

```
{!geofilt sfield=<location_field> pt=<lat,lon> d=<distance>}
```

Here are some examples with values configured:

```
fq={!geofilt sfield=store pt=10.312,-20.556 d=3.5}
```

```
fq={!geofilt sfield=store}&pt=10.312,-20&d=3.5
```

```
fq={!geofilt}&sfield=store&pt=10.312,-20&d=3.5
```

If using `geofilt` with `LatLonType`, it is capable of producing scores equal to the computed distance from the point to the field, making it useful as a component of the main query or a boosting query.

There is more information about spatial searches available in the section [Spatial Search](#).

Surround Query Parser

SurroundQParser extends the QParserPlugin. This provides support for the Surround query syntax, which provides proximity search functionality. There are two operators: `w` creates an ordered span query and `n` creates an unordered one. Both operators take a numeric value to indicate distance between two terms. The default is 1, and the maximum is 99. Note that the query string is not analyzed in any way.

Example:

```
{!surround 3w(foo, bar)}
```

This example would find documents where the terms "foo" and "bar" were no more than 3 terms away from each other (i.e., no more than 2 terms between them).

This query parser will also accept boolean operators (AND, OR, and NOT, in either upper- or lowercase), wildcards, quoting for phrase searches, and boosting. The `w` and `n` operators can also be expressed in upper- or lowercase.

More information about Surround queries can be found at
<http://wiki.apache.org/solr/SurroundQueryParser>.

Term Query Parser

TermQParser extends the QParserPlugin by creating a single term query from the input value equivalent to `readableToIndexed()`. This is useful for generating filter queries from the external human readable terms returned by the faceting or terms components. The only parameter is `f`, for the field.

Example:

```
{!term f=weight}1.5
```

For text fields, no analysis is done since raw terms are already returned from the faceting and terms components. To apply analysis to text fields as well, see the [Field Query Parser](#), above.

If no analysis or transformation is desired for any type of field, see the [Raw Query Parser](#), above.

Highlighting

Solr provides a collection of highlighting utilities which can be called by various Request Handlers to include "highlighted" matches in field values. These highlighting utilities may be used with the [DisMax](#), [Extended DisMax](#), or [standard](#) query parsers.

 Only text that has been both indexed and stored may be highlighted.

Some parameters may be overridden on a per-field basis with the following syntax:
`f.<fieldName>.originalParam=<value>`. For example: `f.contents.hl.snippets=2`

The table below describes Solr's parameters for highlighting.

Parameter	Description
hl	<p>When set to "true", enables highlighted snippets to be generated in the query response. If set to "false" or to a blank or missing value, disables highlighting.</p> <p>The default value is blank, which disables highlighting.</p>
hl.q	Specifies an overriding query term for highlighting. If <code>hl.q</code> is specified, the highlighter will use that term rather than the main query term.
hl.fl	<p>Specifies a list of fields to highlight. Accepts a comma- or space-delimited list of fields for which Solr should generate highlighted snippets. If left blank, highlights the <code>defaultSearchField</code> (or the field specified the <code>df</code> parameter if used) for the <code>StandardRequestHandler</code>. For the <code>DisMaxRequestHandler</code>, the <code>qf</code> fields are used as defaults.</p> <p>A '*' can be used to match field globs, such as 'text_*' or even '*' to highlight on all fields where highlighting is possible. When using '*', consider adding <code>hl.requireFieldMatch=true</code>.</p> <p>The default value is blank.</p>

hl.snippets	<p>Specifies maximum number of highlighted snippets to generate per field. Note: it is possible for any number of snippets from zero to this value to be generated. This parameter accepts per-field overrides.</p> <p>The default value is "1".</p>
hl.fragsize	<p>Specifies the size, in characters, of fragments to consider for highlighting. "0" indicates that the whole field value should be used (no fragmenting). This parameter accepts per-field overrides.</p> <p>The default value is "100".</p>
hl.mergeContinuous	<p>Instructs Solr to collapse contiguous fragments into a single fragment. "true" indicates contiguous fragments will be collapsed into single fragment. This parameter accepts per-field overrides.</p> <p>The default value is "false", which is also the backward-compatible setting.</p>
hl.requireFieldMatch	<p>If set to true, highlights terms only if they appear in the specified field. Normally, terms are highlighted in all requested fields regardless of which field matched the query.</p> <p>The default value is "false".</p>
hl.maxAnalyzedChars	<p>Specifies the number of characters into a document that Solr should look for suitable snippets.</p> <p>The default value is "51200".</p>
hl.alternateField	<p>Specifies a field to be used as a backup default summary if Solr cannot generate a snippet (because no terms match). This parameter accepts per-field overrides.</p> <p>By default, Solr does not select a field for a backup summary.</p>
hl.maxAlternateFieldLength	<p>Specifies the maximum number of characters of the field to return. Any value less than or equal to 0 means the field's length is unlimited.</p> <p>The default value is unlimited.</p> <p>Requires the use of the <code>hl.alternateField</code> parameter.</p>

hl.formatter	Selects a formatter for the highlighted output. Currently the only legal value is "simple", which surrounds a highlighted term with a customizable pre- and post-text snippet. This parameter accepts per-field overrides. The default value is "simple".
hl.simple.pre hl.simple.post	Specifies the text that should appear before and after a highlighted term when using the simple formatter. This parameter accepts per-field overrides. The default values are "" and "".
hl.fragmenter	Specifies a text snippet generator for highlighted text. The standard fragmenter is <code>gap</code> (which is so called because it creates fixed-sized fragments with gaps for multi-valued fields). Another option is <code>regex</code> , which tries to create fragments that resemble a specified regular expression. The <code>hl.fragmenter</code> parameter accepts per-field overrides. The default value is <code>gap</code> .
hl.useFastVectorHighlighter	The <code>FastVectorHighlighter</code> is a <code>TermVector</code> -based highlighter that offers higher performance than the standard highlighter in many cases. To use the <code>FastVectorHighlighter</code> , set this parameter to <code>true</code> . You must also turn on <code>termVectors</code> , <code>termPositions</code> , and <code>termOffsets</code> . Lastly, you should use a boundary scanner to prevent the <code>FastVectorHighlighter</code> from truncating your terms. In most cases, using the <code>breakIterator</code> boundary scanner will give you excellent results. See the following topic for more details about boundary scanners.
hl.phraseLimit	To improve the performance of the <code>FastVectorHighlighter</code> , you can set a limit on the number (<code>int</code>) of phrases to be analyzed for highlighting. The default value for this parameter is <code>integer.MAX_VALUE</code> .
hl.boundaryScanner	Specifies one of two boundary scanners to use with the <code>FastVectorHighlighter</code> : <code>simple</code> or <code>breakIterator</code> . See the following topic for more information about the boundary scanners.

hl.usePhraseHighlighter	If set to "true," instructs Solr to use the Lucene SpanScorer class to highlight phrase terms only when they appear within the query phrase in the document. The default is "true."
hl.highlightMultiTerm	If set to "true," instructs Solr to highlight phrase terms that appear in multi-term queries. The default is "true."
hl.regex.slop	<p>Specifies the factor by which the <code>regex</code> fragmenter can stray from the ideal fragment size (given by <code>hl fragsize</code>) to accommodate a regular expression. For instance, a slop of 0.2 with <code>fragsize</code> of 100 should yield fragments between 80 and 120 characters in length. It is usually good to provide a slightly smaller <code>fragsize</code> when using the <code>regex</code> fragmenter.</p> <p>The default value is 0.6.</p>
hl.regex.pattern	Specifies the regular expression for fragmenting. This could be used to extract sentences.
hl.regex.maxAnalyzedChars	Instructs Solr to analyze only this many characters from a field when using the <code>regex</code> fragmenter (after which, the fragmenter produces fixed-sized fragments). Applying a complicated <code>regex</code> to a huge field is computationally expensive.
hl.preserveMulti	If true , multi-valued fields will return all values in the order they were saved in the index. If false , the default, only values that match the highlight request will be returned. This parameter is new for Solr 4.1 .

Using Boundary Scanners with the Fast Vector Highlighter

The Fast Vector Highlighter will occasionally truncate highlighted words. To prevent this, implement a boundary scanner in `solrconfig.xml`, then use the `hl.boundaryScanner` parameter to specify the boundary scanner for highlighting.

Solr supports two boundary scanners: `breakIterator` and `simple`.

The `breakIterator` Boundary Scanner

The `breakIterator` boundary scanner offers excellent performance right out of the box by taking locale and boundary type into account. In most cases you will want to use the `breakIterator` boundary scanner. To implement the `breakIterator` boundary scanner, add this code to the highlighting section of your `solrconfig.xml` file, adjusting the type, language, and country values as appropriate to your application:

```
<boundaryScanner name="breakIterator"
  class="solr.highlight.BreakIteratorBoundaryScanner">
  <lst name="defaults">
    <str name="hl.bs.type">WORD</str>
    <str name="hl.bs.language">en</str>
    <str name="hl.bs.country">US</str>
  </lst>
</boundaryScanner>
```

Possible values for the `hl.bs.type` parameter are WORD, LINE, SENTENCE, and CHARACTER.

The simple Boundary Scanner

The `simple` boundary scanner scans term boundaries for a specified maximum character value and for common delimiters such as punctuation marks. The `simple` boundary scanner may be useful for some custom To implement the `simple` boundary scanner, add this code to the highlighting section of your `solrconfig.xml` file, adjusting the values as appropriate to your application:

```
<boundaryScanner name="simple" class="solr.highlight.SimpleBoundaryScanner"
  default="true">
  <lst name="defaults">
    <str name="hl.bs.maxScan">10</str>
    <str name="hl.bs.chars">.,!?\t\n</str>
  </lst>
</boundaryScanner>
```

MoreLikeThis

The MoreLikeThis search component enables users to query for documents similar to a document in their result list. It does this by using terms from the original document to find similar documents in the index.

There are three ways to use MoreLikeThis. The first, and most common, is to use it as a request handler. In this case, you would send text to the MoreLikeThis request handler as needed (as in when a user clicked on a "similar documents" link). The second is to use it as a search component. This is less desirable since it performs the MoreLikeThis analysis on every document returned. This may slow search results. The final approach is to use it as a request handler but with externally supplied text. This case, also referred to as the MoreLikeThisHandler, will supply information about similar documents in the index based on the text of the input document.

Covered in this section:

- [How MoreLikeThis Works](#)
- [Common Parameters for MoreLikeThis](#)
- [Parameters for the MoreLikeThisComponent.](#)
- [Parameters for the MoreLikeThisHandler](#)
- [Related Topics](#)

How MoreLikeThis Works

MoreLikeThis constructs a Lucene query based on terms in a document. It does this by pulling terms from the defined list of fields (see the `mlt.fl` parameter, below). For best results, the fields should have stored term vectors in `schema.xml`. For example:

```
<field name="cat" ... termVectors="true" />
```

If term vectors are not stored, MoreLikeThis will generate terms from stored fields. A `uniqueKey` must also be stored in order for MoreLikeThis to work properly.

The next phase filters terms from the original document using thresholds defined with the MoreLikeThis parameters. Finally, a query is run with these terms, and any other query parameters that have been defined (see the `mlt.qf` parameter, below) and a new document set is returned.

 In Solr 4.1, MoreLikeThis supports distributed search.

Common Parameters for MoreLikeThis

The table below summarizes the `MoreLikeThis` parameters supported by Lucene/Solr. These parameters can be used with any of the three possible `MoreLikeThis` approaches.

Parameter	Description
<code>mlt.fl</code>	Specifies the fields to use for similarity. If possible, these should have stored <code>termVectors</code> .
<code>mlt.mintf</code>	Specifies the Minimum Term Frequency, the frequency below which terms will be ignored in the source document.
<code>mlt.mindf</code>	Specifies the Minimum Document Frequency, the frequency at which words will be ignored which do not occur in at least this many documents.
<code>mlt.maxdf</code>	Specifies the Maximum Document Frequency, the frequency at which words will be ignored which occur in more than this many documents. New in Solr 4.1
<code>mlt.minwl</code>	Sets the minimum word length below which words will be ignored.
<code>mlt.maxwl</code>	Sets the maximum word length above which words will be ignored.
<code>mlt.maxqtf</code>	Sets the maximum number of query terms that will be included in any generated query.
<code>mlt.maxntp</code>	Sets the maximum number of tokens to parse in each example document field that is not stored with <code>TermVector</code> support.
<code>mlt.boost</code>	Specifies if the query will be boosted by the interesting term relevance. It can be either "true" or "false".
<code>mlt.qf</code>	Query fields and their boosts using the same format as that used by the <code>DisMaxRequestHandler</code> . These fields must also be specified in <code>mlt.fl</code> .

Parameters for the `MoreLikeThisComponent`.

Using `MoreLikeThis` as a search component returns similar documents for each document in the response set. In addition to the common parameters, these additional options are available:

Parameter	Description
<code>mlt</code>	If set to true, activates the <code>MoreLikeThis</code> component and enables Solr to return <code>MoreLikeThis</code> results.
<code>mlt.count</code>	Specifies the number of similar documents to be returned for each result. The default value is 5.

Parameters for the `MoreLikeThisHandler`

The table below summarizes parameters accessible through the `MoreLikeThisHandler`. It supports faceting, paging, and filtering using common query parameters, but does not work well with alternate query parsers.

Parameter	Description
<code>mlt.match.include</code>	Specifies whether or not the response should include the matched document. If set to false, the response will look like a normal select response.
<code>mlt.match.offset</code>	Specifies an offset into the main query search results to locate the document on which the <code>MoreLikeThis</code> query should operate. By default, the query operates on the first result for the <code>q</code> parameter.
<code>mlt.interestingTerms</code>	Controls how the <code>MoreLikeThis</code> component presents the "interesting" terms (the top TF/IDF terms) for the query. Supports three settings. The setting <code>list</code> lists the terms. The setting <code>none</code> lists no terms. The setting <code>details</code> lists the terms along with the boost value used for each term. Unless <code>mlt.boost=true</code> , all terms will have <code>boost=1.0</code> .

Related Topics

- [RequestHandlers and SearchComponents in SolrConfig](#)

Faceting

As described in the section [Overview of Searching in Solr](#), faceting is the arrangement of search results into categories based on indexed terms. Searchers are presented with the indexed terms, along with numerical counts of how many matching documents were found for each term. Faceting makes it easy for users to explore search results, narrowing in on exactly the results they are looking for.

Topics covered on this page:

- [General Parameters](#)
- [Field-Value Faceting Parameters](#)
- [Range Faceting](#)
- [Date Faceting Parameters](#)
- [Local Parameters for Faceting](#)
- [Pivot \(Decision Tree\) Faceting](#)
- [Facets and Time Zone](#)
- [Related Topics](#)

General Parameters

The table below summarizes the general parameters for controlling facetting.

Parameter	Description
<code>facet</code>	If set to true, enables facetting.
<code>facet.query</code>	Specifies a Lucene query to generate a facet count.

These parameters are described in the sections below.

The facet Parameter

If set to "true," this parameter enables facet counts in the query response. If set to "false" to a blank or missing value, this parameter disables faceting. None of the other parameters listed below will have any effect unless this parameter is set to "true." The default value is blank.

The facet.query Parameter

This parameter allows you to specify an arbitrary query in the Lucene default syntax to generate a facet count. By default, Solr's faceting feature automatically determines the unique terms for a field and returns a count for each of those terms. Using `facet.query`, you can override this default behavior and select exactly which terms or expressions you would like to see counted. In a typical implementation of faceting, you will specify a number of `facet.query` parameters. This parameter can be particularly useful for numeric-range-based facets or prefix-based facets.

You can set the `facet.query` parameter multiple times to indicate that multiple queries should be used as separate facet constraints.

To use facet queries in a syntax other than the default syntax, prefix the facet query with the name of the query notation. For example, to use the hypothetical `myfunc` query parser, you could set the `facet.query` parameter like so:

```
facet.query={!myfunc}name~fred
```

Field-Value Faceting Parameters

Several parameters can be used to trigger faceting based on the indexed terms in a field.

When using this parameter, it is important to remember that "term" is a very specific concept in Lucene: it relates to the literal field/value pairs that are indexed after any analysis occurs. For text fields that include stemming, lowercasing, or word splitting, the resulting terms may not be what you expect. If you want Solr to perform both analysis (for searching) and faceting on the full literal strings, use the `copyField` directive in the `schema.xml` file to create two versions of the field: one `Text` and one `String`. Make sure both are `indexed="true"`. (For more information about the `copyField` directive, see [Documents, Fields, and Schema Design](#).)

The table below summarizes Solr's field value faceting parameters.

Parameter	Description
<code>facet.field</code>	Identifies a field to be treated as a facet.
<code>facet.prefix</code>	Limits the terms used for faceting to those that begin with the specified prefix.
<code>facet.sort</code>	Controls how faceted results are sorted.

<code>facet.limit</code>	Controls how many constraints should be returned for each facet.
<code>facet.offset</code>	Specifies an offset into the facet results at which to begin displaying facets.
<code>facet.mincount</code>	Specifies the minimum counts required for a facet field to be included in the response.
<code>facet.missing</code>	Controls whether Solr should compute a count of all matching results which have no value for the field, in addition to the term-based constraints of a facet field.
<code>facet.method</code>	Selects the algorithm or method Solr should use when faceting a field.
<code>facet.enum.cache.minDF</code>	Specifies the minimum document frequency (the number of documents matching a term) for which the <code>filterCache</code> should be used when determining the constraint count for that term.

These parameters are described in the sections below.

The `facet.field` Parameter

The `facet.field` parameter identifies a field that should be treated as a facet. It iterates over each Term in the field and generate a facet count using that Term as the constraint. This parameter can be specified multiple times in a query to select multiple facet fields.

 If you do not set this parameter to at least one field in the schema, none of the other parameters described in this section will have any effect.

The `facet.prefix` Parameter

The `facet.prefix` parameter limits the terms on which to facet to those starting with the given string prefix. This does not limit the query in any way, only the facets that would be returned in response to the query.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.prefix`.

The `facet.sort` Parameter

This parameter determines the ordering of the facet field constraints.

 The true/false values for this parameter were deprecated in Solr 1.4.

<code>facet.sort</code>	Results
Setting	
count	Sort the constraints by count (highest count first).
index	Return the constraints sorted in their index order (lexicographic by indexed term). For terms in the ASCII range, this will be alphabetically sorted.

The default is `count` if `facet.limit` is greater than 0, otherwise, the default is `index`.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.sort`

The `facet.limit` Parameter

This parameter specifies the maximum number of constraint counts (essentially, the number of facets for a field that are returned) that should be returned for the facet fields. A negative value means that Solr will return unlimited number of constraint counts.

The default value is 100.

This parameter can be specified on a per-field basis to apply a distinct limit to each field with the syntax of `f.<fieldname>.facet.limit`.

The `facet.offset` Parameter

The `facet.offset` parameter indicates an offset into the list of constraints to allow paging.

The default value is 0.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.offset`.

The `facet.mincount` Parameter

The `facet.mincount` parameter specifies the minimum counts required for a facet field to be included in the response. If a field's counts are below the minimum, the field's facet is not returned.

The default value is 0.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.mincount`.

The `facet.missing` Parameter

If set to true, this parameter indicates that, in addition to the Term-based constraints of a facet field, a count of all results that match the query but which have no facet value for the field should be computed and returned in the response.

The default value is false.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.missing`.

The `facet.method` Parameter

The `facet.method` parameter selects the type of algorithm or method Solr should use when faceting a field.

Setting	Results
enum	Enumerates all terms in a field, calculating the set intersection of documents that match the term with documents that match the query. This method is recommended for faceting multi-valued fields that have only a few distinct values. The average number of values per document does not matter. For example, faceting on a field with U.S. States such as Alabama, Alaska, ... Wyoming would lead to fifty cached filters which would be used over and over again. The <code>filterCache</code> should be large enough to hold all the cached filters.
fc	Calculates facet counts by iterating over documents that match the query and summing the terms that appear in each document. This is currently implemented using an <code>UnInvertedField</code> cache if the field either is multi-valued or is tokenized (according to <code>FieldType.isTokened()</code>). Each document is looked up in the cache to see what terms/values it contains, and a tally is incremented for each value. This method is excellent for situations where the number of indexed values for the field is high, but the number of values per document is low. For multi-valued fields, a hybrid approach is used that uses term filters from the <code>filterCache</code> for terms that match many documents. The letters <code>fc</code> stand for field cache.
fcs	Per-segment field faceting for single-valued string fields. Enable with <code>facet.method=fcs</code> and control the number of threads used with the <code>threads</code> local parameter. This parameter allows faceting to be faster in the presence of rapid index changes.

The default value is `fc` (except for fields using the `BoolField` field type) since it tends to use less memory and is faster when a field has many unique terms in the index.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.method`.

The facet.enum.cache.minDf Parameter

This parameter indicates the minimum document frequency (the number of documents matching a term) for which the filterCache should be used when determining the constraint count for that term. This is only used with the `facet.method=enum` method of faceting.

A value greater than zero decreases the filterCache's memory usage, but increases the time required for the query to be processed. If you are faceting on a field with a very large number of terms, and you wish to decrease memory usage, try setting this parameter to a value between 25 and 50, and run a few tests. Then, optimize the parameter setting as necessary.

The default value is 0, causing the filterCache to be used for all terms in the field.

This parameter can be specified on a per-field basis with the syntax of `f.<fieldname>.facet.enum.cache.minDF`.

Range Faceting

You can use Range Faceting on any date field or any numeric field that supports range queries. This is particularly useful for stitching together a series of range queries (as facet by query) for things like prices. As of Solr 3.1, Range Faceting is preferred over [Date Faceting](#) (described below).

Parameter	Description
<code>facet.range</code>	Specifies the field to facet by range.
<code>facet.range.start</code>	Specifies the start of the facet range.
<code>facet.range.end</code>	Specifies the end of the facet range.
<code>facet.range.gap</code>	Specifies the span of the range as a value to be added to the lower bound.
<code>facet.range.hardend</code>	A boolean parameter that specifies how Solr handles a range gap that cannot be evenly divided between the range start and end values. If true, the last range constraint will have the <code>facet.range.end</code> value an upper bound. If false, the last range will have the smallest possible upper bound greater than <code>facet.range.end</code> such that the range is the exact width of the specified range gap. The default value for this parameter is false.
<code>facet.range.include</code>	Specifies inclusion and exclusion preferences for the upper and lower bounds of the range. See the <code>facet.range.include</code> topic for more detailed information.
<code>facet.range.other</code>	Specifies counts for Solr to compute in addition to the counts for each facet range constraint.

The facet.range Parameter

The `facet.range` parameter defines the field for which Solr should create range facets. For example:

```
facet.range=price&facet.range=age
```

The `facet.range.start` Parameter

The `facet.range.start` parameter specifies the lower bound of the ranges. You can specify this parameter on a per field basis with the syntax of `f.<fieldname>.facet.range.start`. For example:

```
f.price.facet.range.start=0.0&f.age.facet.range.start=10
```

The `facet.range.end` Parameter

The `facet.range.end` specifies the upper bound of the ranges. You can specify this parameter on a per field basis with the syntax of `f.<fieldname>.facet.range.end`. For example:

```
f.price.facet.range.end=1000.0&f.age.facet.range.start=99
```

The `facet.range.gap` Parameter

The span of each range expressed as a value to be added to the lower bound. For date fields, this should be expressed using the [DateMathParser syntax](#) (such as, `facet.range.gap=%2B1DAY ... '+1DAY'`). You can specify this parameter on a per-field basis with the syntax of `f.<fieldname>.facet.range.gap`. For example:

```
f.price.facet.range.gap=100&f.age.facet.range.gap=10
```

Gaps can also be variable width by passing in a comma separated list of the gap size to be used. The last gap specified will be used to fill out all remaining gaps if the number of gaps given does not go evenly into the range. Variable width gaps are useful, for example, in spatial applications where one might want to facet by distance into three buckets: walking (0-5KM), driving (5-100KM), or other (100KM+). For example:

```
facet.date.gap=1,2,3,10
```

This creates 4+ buckets of size, 1, 2, 3 and then 0 or more buckets of 10 days each, depending on the start and end values.

The `facet.range.hardend` Parameter

The `facet.range.hardend` parameter is a Boolean parameter that specifies how Solr should handle cases where the `facet.range.gap` does not divide evenly between `facet.range.start` and `facet.range.end`. If **true**, the last range constraint will have the `facet.range.end` value as an upper bound. If **false**, the last range will have the smallest possible upper bound greater than `facet.range.end` such that the range is the exact width of the specified range gap. The default value for this parameter is false. This parameter can be specified on a per field basis with the syntax `f.<fieldname>.facet.range.hardend`.

The `facet.range.include` Parameter

By default, the ranges used to compute range facetting between `facet.range.start` and `facet.range.end` are inclusive of their lower bounds and exclusive of the upper bounds. The "before" range defined with the `facet.range.other` parameter is exclusive and the "after" range is inclusive. This default, equivalent to "lower" below, will not result in double counting at the boundaries. You can use the `facet.range.include` parameter to modify this behavior using the following options:

Option	Description
lower	All gap-based ranges include their lower bound.
upper	All gap-based ranges include their upper bound.
edge	The first and last gap ranges include their edge bounds (lower for the first one, upper for the last one) even if the corresponding upper/lower option is not specified.
outer	The "before" and "after" ranges will be inclusive of their bounds, even if the first or last ranges already include those boundaries.
all	Includes all options: lower, upper, edge, outer.

You can specify this parameter on a per field basis with the syntax of `f.<fieldname>.facet.range.include`, and you can specify it multiple times to indicate multiple choices.

-  To ensure you avoid double-counting, do not choose both `lower` and `upper`, do not choose `outer`, and do not choose `all`.

The `facet.range.other` Parameter

The `facet.range.other` parameter specifies that in addition to the counts for each range constraint between `facet.range.start` and `facet.range.end`, counts should also be computed for these options:

Option	Description
before	All records with field values lower than lower bound of the first range.
after	All records with field values greater than the upper bound of the last range.
between	All records with field values between the start and end bounds of all ranges.
none	Do not compute any counts.
all	Compute counts for before, between, and after.

This parameter can be specified on a per field basis with the syntax of `f.<fieldname>.facet.range.other`. In addition to the `all` option, this parameter can be specified multiple times to indicate multiple choices, but `none` will override all other options.

facet.range.hardend

A Boolean parameter instructing Solr what to do in the event that `facet.range.gap` does not divide evenly between `facet.range.start` and `facet.range.end`. If this is true, the last range constraint will have an upper bound of `facet.range.end`; if false, the last range will have the smallest possible upper bound greater than `facet.range.end` such that the range is exactly `facet.range.gap` wide.

The default is false.

This parameter can be specified on a per field basis.

facet.range.other

This param indicates that in addition to the counts for each range constraint between `facet.range.start` and `facet.range.end`, counts should also be computed for...

- before all records with field values lower than lower bound of the first range
- after all records with field values greater than the upper bound of the last range
- between all records with field values between the start and end bounds of all ranges
- none compute none of this information
- all shortcut for before, between, and after

This parameter can be specified on a per field basis.

In addition to the `all` option, this parameter can be specified multiple times to indicate multiple choices -- but `none` will override all other options.

facet.range.include

By default, the ranges used to compute range facetting between `facet.range.start` and `facet.range.end` are inclusive of their lower bounds and exclusive of the upper bounds. The "before" range is exclusive and the "after" range is inclusive. This default, equivalent to `lower` below, will *not* result in double counting at the boundaries. This behavior can be modified by the `facet.range.include` param, which can be any combination of the following options...

- `lower` = all gap based ranges include their lower bound
- `upper` = all gap based ranges include their upper bound
- `edge` = the first and last gap ranges include their edge bounds (i.e., lower for the first one, upper for the last one) even if the corresponding upper/lower option is not specified
- `outer` = the "before" and "after" ranges will be inclusive of their bounds, even if the first or last ranges already include those boundaries.
- `all` = shorthand for lower, upper, edge, outer

This parameter can be specified on a per field basis.

This parameter can be specified multiple times to indicate multiple choices.

If you want to ensure you don't double-count, don't choose both lower & upper, don't choose outer, and don't choose all.

Date Faceting Parameters

As of Solr 3.1, date facetting has been deprecated in favor of [Range Faceting](#), which provides more flexibility with dates and numeric fields. Date Faceting can be used, however. The response structure is slightly different, but the functionality is equivalent (except that it supports numeric fields as well as dates).

Several parameters can be used to trigger facetting based on Date ranges computed using simple [DateMathParser](#) expressions.

When using Date Faceting, the `facet.date`, `facet.date.start`, `facet.date.end`, and `facet.date.gap` parameters are all mandatory.

Name	What it does
<code>facet.date</code>	Allows you to specify names of fields (of type <code>DateField</code> , described in the section, Field Types Included with Solr) which should be treated as date facets. Can be specified multiple times to indicate multiple date facet fields.
<code>facet.date.start</code>	The lower boundary for the first date range for all Date Faceting on this field. This should be a single date expression which may use the DateMathParser syntax. Can be specified on a per field basis.

facet.date.end	The minimum upper boundary for the last date range for all Date Faceting on this field. This should be a single date expression which may use the DateMathParser syntax. Can be specified on a per field basis.
facet.date.gap	The size of each date range expressed as an interval to be added to the lower bound using the DateMathParser syntax. Can be specified on a per field basis. Example: facet.date.gap=%2B1DAY (+1DAY)
facet.date.other	<p>Indicates that in addition to the counts for each date range constraint between facet.date.start and facet.date.end, counts should also be computed for:</p> <ul style="list-style-type: none">• before: all records with field values lower than lower bound of the first range• after: all records with field values greater than the upper bound of the last range• between: all records with field values between the start and end bounds of all ranges• none: compute none of this information• all: shortcut for before, between, and after. <p>Can be specified on a per field basis. In addition to the all option, this parameter can be specified multiple times to indicate multiple choices, but none will override all other options.</p>

facet.date.include	<p>By default, the ranges used to compute date faceting between <code>facet.date.start</code> and <code>facet.date.end</code> are all inclusive of both endpoints, while the "before" and "after" ranges are not inclusive. This behavior can be modified by the <code>facet.date.include</code> parameter, which can be any combination of the following options:</p> <ul style="list-style-type: none">• <code>lower</code> = all gap based ranges include their lower bound• <code>upper</code> = all gap based ranges include their upper bound• <code>edge</code> = the first and last gap ranges include their edge bounds (ie: lower for the first one, upper for the last one) even if the corresponding upper/lower option is not specified• <code>outer</code> = the "before" and "after" ranges will be inclusive of their bounds, even if the first or last ranges already include those boundaries.• <code>all</code> = shorthand for lower, upper, edge, outer <p>This parameter can be specified on a per field basis, and can be specified multiple times to indicate multiple choices.</p>
--------------------	---

Local Parameters for Faceting

The [LocalParams syntax](#) provides a method of adding metadata to other parameter values, much like XML attributes.

Tagging and Excluding Filters

You can tag specific filters and exclude those filters when faceting. This is useful when doing multi-select faceting.

Consider the following example query with faceting:

```
q=mainquery&fq=status:public&fq=doctype:pdf&facet=on&facet.field=doctype
```

Because everything is already constrained by the filter `doctype:pdf`, the `facet.field=doctype` facet command is currently redundant and will return 0 counts for everything except `doctype:pdf`.

To implement a multi-select facet for `doctype`, a GUI may want to still display the other `doctype` values and their associated counts, as if the `doctype:pdf` constraint had not yet been applied. For example:

```
==== Document Type ====
[ ] Word (42)
[x] PDF (96)
[ ] Excel(11)
[ ] HTML (63)
```

To return counts for doctype values that are currently not selected, tag filters that directly constrain doctype, and exclude those filters when faceting on doctype.

```
q=mainquery&fq=status:public&fq={!tag=dt}doctype:pdf&facet=on&facet.field={!ex=dt}do
```

Filter exclusion is supported for all types of facets. Both the `tag` and `ex` local parameters may specify multiple values by separating them with commas.

Changing the Output Key

To change the output key for a facetting command, specify a new name with the `key` local parameter. For example:

```
facet.field={!ex=dt key=mylabel}doctype
```

The parameter setting above causes the results to be returned under the key "mylabel" rather than "doctype" in the response. This can be helpful when facetting on the same field multiple times with different exclusions.

Pivot (Decision Tree) Faceting

Pivoting is a summarization tool that lets you automatically sort, count, total or average data stored in a table. It displays the results in a second table showing the summarized data. Pivot facetting lets you create a summary table of the results from a query across numerous documents. With Solr 4, pivot facetting supports nested facet queries, not just facet fields.

Another way to look at it is that the query produces a Decision Tree, in that Solr tells you "for facet A, the constraints/counts are X/N, Y/M, etc. If you were to constrain A by X, then the constraint counts for B would be S/P, T/Q, etc.". In other words, it tells you in advance what the "next" set of facet results would be for a field if you apply a constraint from the current facet results.

facet.pivot

The `facet.pivot` parameter defines the fields to use for the pivot. Multiple `facet.pivot` values will create multiple "facet_pivot" sections in the response. Separate each list of fields with a comma.

As of Solr 4.1, local parameters can be used in pivot facet queries.

facet.pivot.mincount

The `facet.pivot.mincount` parameter defines the minimum number of documents that need to match in order for the facet to be included in results. The default is 1.

For example, we can use Solr's example data set to make a query like this:

```
http://localhost:8983/solr/select?q=*&facet.pivot=cat,popularity,inStock&facet.pivot=pop&facet=true&facet.field=cat&facet.limit=5&rows=0&wt=json&indent=true&facet.pivot.mincount=
```

This query will returns the data below, with the pivot faceting results found in the section "facet_pivot":

```
"facet_counts":{  
    "facet_queries":{},  
    "facet_fields":{  
        "cat": [  
            "electronics", 14,  
            "currency", 4,  
            "memory", 3,  
            "connector", 2,  
            "graphics card", 2]  
    },  
    "facet_dates":{},  
    "facet_ranges":{},  
    "facet_pivot":{  
        "cat,popularity,inStock": [ {  
            "field": "cat",  
            "value": "electronics",  
            "count": 14,  
            "pivot": [ {  
                "field": "popularity",  
                "value": 6,  
                "count": 5,  
                "pivot": [ {  
                    "field": "inStock",  
                    "value": true,  
                    "count": 5 } ]  
            ]  
        } ]  
    }  
}
```

Pivot faceting supports distributed searching.

Facets and Time Zone

You can construct your query to use the `tz` parameter, which overrides the time zone used when rounding dates in DateMath expressions for the entire request. This applies to all date range queries and date facetting.

For instance, a request that does not include date parameters, like this:

```
http://localhost:8983/solr/select?indent=on&version=2.2&q=*&%3A*&fq=&start=0&rows=10&fl=*&%2&qt=&wt=&explanOther=&hl.fl=&
```

can be modified by adding facet parameters for dates:

```
&facet=true&facet.date=manufacturedate_dt&facet.date.start=2005-01-31T15:00:00Z  
&facet.date.end=2006-05-31T15:00:00Z&facet.date.gap=%2B1MONTH/DAY&facet.date.tz=Asia/Tokyo
```

The result would be a date facet request with monthly gap in Tokyo time (GMT+9:00).

One problem may be that some time zones observe Daylight Savings Time, and some don't. For instance, consider this query:

```
http://localhost:8983/solr/select?indent=on&version=2.2&q=*&%3A*&fq=&start=0&rows=10&fl=*&%2&explanOther=&hl.fl=&start = 2007-07-30T00:00:00.000Z/DAY&end =  
2007-07-31T00:00:00.000Z/DAY  
&gap = \+1DAY&tz=Asia/Singapore
```

The way to make sure that your return reflects Daylight Savings Time is to simply add '/DAY' to the end of each parameter, which will round the times to the beginning of the day. Since the rounding takes `tz` into account, this query will come out to 8:00 (the difference between GMT and Singapore time).

```
&start=2007-07-30T00:00:00.000Z/DAY&end=2007-07-31T00:00:00.000Z/DAY  
&gap=+1DAY&tz=Asia/Singapore
```

Related Topics

- [SimpleFacetParameters](#) from the Solr Wiki.

Result Grouping

Result Grouping groups documents with a common field value into groups and returns the top documents for each group. For example, if you searched for "DVD" on an electronic retailer's e-commerce site, you might be returned three categories such as "TV and Video," "Movies," and "Computers," with three results per category. In this case, the query term "DVD" appeared in all three categories, so Solr groups them together in order to increase relevancy for the user.

Result Grouping is separate from [Faceting](#). Though it is conceptually similar, facetting returns all relevant results and allows the user to refine the results based on the facet category. For example, if you searched for "shoes" on a footwear retailer's e-commerce site, you would be returned all results for that query term, along with selectable facets such as "size," "color," "brand," and so on.

However, with Solr 4 you can also group facets. The grouped facetting works with the first `group.field` parameter, and other `group.field` parameters are ignored. Grouped facetting only supports `facet.field` for string based fields that are not tokenized and are not multivalued.

Grouped facetting currently doesn't support date and pivot facetting, but it does support range facetting.

Grouped facetting differs from non grouped facets ($\text{sum of all facets} == (\text{total of products with that property})$) as shown in the following example:

Object 1

- name: Phaser 4620a
- ppm: 62
- product_range: 6

Object 2

- name: Phaser 4620i
- ppm: 65
- product_range: 6

Object 3

- name: ML6512
- ppm: 62
- product_range: 7

If you ask Solr to group these documents by "product_range", then the total amount of groups is 2, but the facets for ppm are 2 for 62 and 1 for 65.

Request Parameters

Result Grouping takes the following request parameters. Any number of these request parameters can be included in a single request:

Parameter	Type	Description
group	Boolean	If true, query results will be grouped.
group.field	string	The name of the field by which to group results. The field must be single-valued, and either be indexed or a field type that has a value source and works in a function query, such as ExternalFileField. It must also be a string-based field, such as StrField or TextField
group.func	query	Group based on the unique values of a function query. Supported only in Sol4r 4.0.
group.query	query	Return a single group of documents that match the given query.
rows	integer	The number of groups to return. The default value is 10.
start	integer	Specifies an initial offset for the list of groups.
group.limit	integer	Specifies the number of results to return for each group. The default value is 1.
group.offset	integer	Specifies an initial offset for the document list of each group.
sort	sortspec	Specifies how Solr sorts the groups relative to each other. For example, <code>sort=popularity desc</code> will cause the groups to be sorted according to the highest popularity document in each group. The default value is <code>score desc</code> .
group.sort	sortspec	Specifies how Solr sorts documents within a single group. The default value is <code>score desc</code> .
group.format	grouped/simple	If this parameter is set to simple, the grouped documents are presented in a single flat list, and the start and rows parameters affect the numbers of documents instead of groups.
group.main	Boolean	If true, the result of the first field grouping command is used as the main result list in the response, using <code>group.format=simple</code> .

group.ngroups	Boolean	If true, Solr includes the number of groups that have matched the query in the results. The default value is false.
group.truncate	Boolean	If true, facet counts are based on the most relevant document of each group matching the query. The default value is false.
group.facet	Boolean	Determines whether to compute grouped facets for the field facets specified in facet.field parameters. Grouped facets are computed based on the first specified group. As with normal field faceting, fields shouldn't be tokenized (otherwise counts are computed for each token). Grouped faceting supports single and multivalued fields. Default is false. New with Solr 4.
group.cache.percent	integer between 0 and 100	Setting this parameter to a number greater than 0 enables caching for result grouping. Result Grouping executes two searches; this option caches the second search. The default value is 0. Testing has shown that group caching only improves search time with Boolean, wildcard, and fuzzy queries. For simple queries like term or "match all" queries, group caching degrades performance.

Any number of group commands (`group.field`, `group.func`, `group.query`) may be specified in a single request.

Grouping is also supported for distributed searches. Currently `group.truncate` and `group.func` are the only parameters that aren't supported for distributed searches.

Examples

All of the following examples work with the data provided in the Solr Example directory.

Grouping Results by Field

In this example, we will group results based on the `manu_exact` field, which specifies the manufacturer of the items in the sample dataset.

```
http://localhost:8983/solr/select?wt=json&indent=true&fl=id,name&q=solr+memory&group
```

```
{  
...  
"grouped":{  
    "manu_exact":{  
        "matches":6,  
        "groups": [ {  
            "groupValue": "Apache Software Foundation",  
            "doclist": { "numFound":1, "start":0, "docs": [  
                {  
                    "id": "SOLR1000",  
                    "name": "Solr, the Enterprise Search Server" } ]  
            } },  
            {  
                "groupValue": "Corsair Microsystems Inc.",  
                "doclist": { "numFound":2, "start":0, "docs": [  
                    {  
                        "id": "VS1GB400C3",  
                        "name": "CORSAIR ValueSelect 1GB 184-Pin DDR SDRAM Unbuffered DDR 400 (PC  
3200) System Memory - Retail" } ]  
                } },  
                {  
                    "groupValue": "A-DATA Technology Inc.",  
                    "doclist": { "numFound":1, "start":0, "docs": [  
                        {  
                            "id": "VDBDB1A16",  
                            "name": "A-DATA V-Series 1GB 184-Pin DDR SDRAM Unbuffered DDR 400 (PC  
3200) System Memory - OEM" } ]  
                    } },  
                    {  
                        "groupValue": "Canon Inc.",  
                        "doclist": { "numFound":1, "start":0, "docs": [  
                            {  
                                "id": "0579B002",  
                                "name": "Canon PIXMA MP500 All-In-One Photo Printer" } ]  
                        } },  
                        {  
                            "groupValue": "ASUS Computer Inc.",  
                            "doclist": { "numFound":1, "start":0, "docs": [  
                                {  
                                    "id": "EN7800GTX/2DHTV/256M",  
                                    "name": "ASUS Extreme N7800GTX/2DHTV (256 MB)" } ]  
                            } }  
            ]  
        } }  
}
```

The response indicates that there are six total matches for our query. For each unique value of `group.field`, Solr returns a `docList` with the top scoring document. The `docList` also includes the total number of matches in that group as the `numFound` value. The groups are sorted by the score of the top document within each group.

We can run the same query with the request parameter `group.main=true`. This will format the results as a single flat document list. This flat format does not include as much information as the normal result grouping query results, but it may be easier for existing Solr clients to parse.

```
http://localhost:8983/solr/select?wt=json&indent=true&fl=id,name,manufacturer&q=solr
```

```
{
  "responseHeader": {
    "status": 0,
    "QTime": 1,
    "params": {
      "fl": "id, name, manufacturer",
      "indent": "true",
      "q": "solr memory",
      "group.field": "manu_exact",
      "group.main": "true",
      "group": "true",
      "wt": "json" } },
  "grouped": {},
  "response": { "numFound": 6, "start": 0, "docs": [
    {
      "id": "SOLR1000",
      "name": "Solr, the Enterprise Search Server" },
    {
      "id": "VS1GB400C3",
      "name": "CORSAIR ValueSelect 1GB 184-Pin DDR SDRAM Unbuffered DDR 400 (PC 3200) System Memory - Retail" },
    {
      "id": "VDBDB1A16",
      "name": "A-DATA V-Series 1GB 184-Pin DDR SDRAM Unbuffered DDR 400 (PC 3200) System Memory - OEM" },
    {
      "id": "0579B002",
      "name": "Canon PIXMA MP500 All-In-One Photo Printer" },
    {
      "id": "EN7800GTX/2DHTV/256M",
      "name": "ASUS Extreme N7800GTX/2DHTV (256 MB)" } ] }
}
```

Grouping by Query

In this example, we will use the `group.query` parameter to find the top three results for "memory" in two different price ranges: 0.00 to 99.99, and over 100.

```
http://localhost:8983/solr/select?wt=json&indent=true&fl=name,price&q=memory&group=t
```

```
{  
  "responseHeader": {  
    "status": 0,  
    "QTime": 42,  
    "params": {  
      "fl": "name,price",  
      "indent": "true",  
      "q": "memory",  
      "group.limit": "3",  
      "group.query": ["price:[0 TO 99.99]",  
                     "price:[100 TO *]"],  
      "group": "true",  
      "wt": "json" } },  
  "grouped": {  
    "price:[0 TO 99.99)": {  
      "matches": 5,  
      "doclist": { "numFound": 1, "start": 0, "docs": [  
        {  
          "name": "CORSAIR ValueSelect 1GB 184-Pin DDR SDRAM Unbuffered DDR 400 (PC  
3200) System Memory - Retail",  
          "price": 74.99}  
      ]},  
      "price:[100 TO *)": {  
        "matches": 5,  
        "doclist": { "numFound": 3, "start": 0, "docs": [  
          {  
            "name": "CORSAIR XMS 2GB (2 x 1GB) 184-Pin DDR SDRAM Unbuffered DDR 400 (PC  
3200) Dual Channel  
Kit System Memory - Retail",  
            "price": 185.0},  
          {  
            "name": "Canon PIXMA MP500 All-In-One Photo Printer",  
            "price": 179.99},  
          {  
            "name": "ASUS Extreme N7800GTX/2DHTV (256 MB)",  
            "price": 479.95}  
        ]}  
      }  
    }  
  }  
}
```

In this case, Solr found five matches for "memory," but only returns four results grouped by price. This is because one result for "memory" did not have a price assigned to it.

Distributed Result Grouping

Solr also supports result grouping on distributed indexes. If you are using result grouping on the "/select" request handler, you must provide the `shards` parameter described here. If you are using result grouping on a request handler other than "/select", you must also provide the `shards.qt` parameter:

Parameter	Description
<code>shards</code>	Specifies the shards in your distributed indexing configuration. For more information about distributed indexing, see Distributed Search with Index Sharding
<code>shards.qt</code>	Specifies the request handler Solr uses for requests to shards. This parameter is not required for the /select request handler.

For example:

```
http://localhost:8983/solr/select?wt=json&indent=true&fl=id,name,manufacturer&q=solr
```

Spell Checking

The SpellCheck component is designed to provide inline query suggestions based on other, similar, terms. The basis for these suggestions can be terms in a field in Solr, externally created text files, or fields in other Lucene indexes.

Topics covered in this section:

- [Configuring the SpellCheckComponent](#)
- [Spell Check Parameters](#)
- [Distributed SpellCheck](#)

Configuring the SpellCheckComponent

Define Spell Check in solrconfig.xml

The first step is to specify the source of terms in `solrconfig.xml`. There are three approaches to spell checking in Solr, discussed below.

IndexBasedSpellChecker

The `IndexBasedSpellChecker` uses a Solr index as the basis for a parallel index used for spell checking. It requires defining a field as the basis for the index terms; a common practice is to copy terms from some fields (such as `title`, `body`, etc.) to another field created for spell checking. Here is a simple example of configuring `solrconfig.xml` with the `IndexBasedSpellChecker`:

```
<searchComponent name="spellcheck" class="solr.SpellCheckComponent">
  <lst name="spellchecker">
    <str name="classname">solr.IndexBasedSpellChecker</str>
    <str name="spellcheckIndexDir">./spellchecker</str>
    <str name="field">content</str>
    <str name="buildOnCommit">true</str>
  </lst>
</searchComponent>
```

The first element defines the `searchComponent` to use the `solr.SpellCheckComponent`. The `classname` is the specific implementation of the `SpellCheckComponent`, in this case `solr.IndexBasedSpellChecker`. Defining the `classname` is optional; if not defined, it will default to `IndexBasedSpellChecker`.

The `spellcheckIndexDir` defines the location of the directory that holds the spellcheck index, while the `field` defines the source field (defined in `schema.xml`) for spell check terms. When choosing a field for the spellcheck index, it's best to avoid a heavily processed field to get more accurate results. If the field has many word variations from processing synonyms and/or stemming, the dictionary will be created with those variations in addition to more valid spelling data.

Finally, `buildOnCommit` defines whether to build the spell check index at every commit (that is, every time new documents are added to the index). It is optional, and can be omitted if you would rather set it to `false`.

DirectSolrSpellChecker

The `DirectSolrSpellChecker` uses terms from the Solr index without building a parallel index like the `IndexBasedSpellChecker`. It is considered experimental and still in development, but is being used widely. This spell checker has the benefit of not having to be built regularly, meaning that the terms are always up-to-date with terms in the index. Here is how this might be configured in `solrconfig.xml`

```
<searchComponent name="spellcheck" class="solr.SpellCheckComponent">
  <lst name="spellchecker">
    <str name="name">default</str>
    <str name="field">name</str>
    <str name="classname">solr.DirectSolrSpellChecker</str>
    <str name="distanceMeasure">internal</str>
    <float name="accuracy">0.5</float>
    <int name="maxEdits">2</int>
    <int name="minPrefix">1</int>
    <int name="maxInspections">5</int>
    <int name="minQueryLength">4</int>
    <float name="maxQueryFrequency">0.01</float>
    <float name="thresholdTokenFrequency">.01</float>
  </lst>
</searchComponent>
```

When choosing a field to query for this spell checker, you want one which has relatively little analysis performed on it (particularly analysis such as stemming). Note that you need to specify a field to use for the suggestions, so like the `IndexBasedSpellChecker`, you may want to copy data from fields like `title`, `body`, etc., to a field dedicated to providing spelling suggestions.

Many of the parameters relate to how this spell checker should query the index for term suggestions. The `distanceMeasure` defines the metric to use during the spell check query. The value "internal" uses the default Levenshtein metric, which is the same metric used with the other spell checker implementations.

Because this spell checker is querying the main index, you may want to limit how often it queries the index to be sure to avoid any performance conflicts with user queries. The `accuracy` setting defines the threshold for a valid suggestion, while `maxEdits` defines the number of changes to the term to allow. Since most spelling mistakes are only 1 letter off, setting this to 1 will reduce the number of possible suggestions (the default, however, is 2); the value can only be 1 or 2. `minPrefix` defines the minimum number of characters the terms should share. Setting this to 1 means that the spelling suggestions will all start with the same letter, for example.

The `maxInspections` parameter defines the maximum number of possible matches to review before returning results; the default is 5. `minQueryLength` defines how many characters must be in the query before suggestions are provided; the default is 4. `maxQueryFrequency` sets the maximum threshold for the number of documents a term must appear in before being considered as a suggestion. This can be a percentage (such as .01, or 1%) or an absolute value (such as 4). A lower threshold is better for small indexes. Finally, `thresholdTokenFrequency` sets the minimum number of documents a term must appear in, and can also be expressed as a percentage or an absolute value.

FileBasedSpellChecker

The `FileBasedSpellChecker` uses an external file as a spelling dictionary. This can be useful if using Solr as a spelling server, or if spelling suggestions don't need to be based on actual terms in the index. In `solrconfig.xml`, you would define the `searchComponent` as so:

```
<searchComponent name="spellcheck" class="solr.SpellCheckComponent">
  <lst name="spellchecker">
    <str name="classname">solr.FileBasedSpellChecker</str>
    <str name="name">file</str>
    <str name="sourceLocation">spellings.txt</str>
    <str name="characterEncoding">UTF-8</str>
    <str name="spellcheckIndexDir">./spellcheckerFile</str>
  </lst>
</searchComponent>
```

The differences here are the use of the `sourceLocation` to define the location of the file of terms and the use of `characterEncoding` to define the encoding of the terms file.

- ① In the previous example, `name` is used to name this specific definition of the spellchecker. Multiple definitions can co-exist in a single `solrconfig.xml`, and the `name` helps to differentiate them when they are defined in the `schema.xml`. If only defining one spellchecker, no name is required.

WordBreakSolrSpellChecker

A parallel implementation, `WordBreakSolrSpellChecker` offers suggestions by combining adjacent query terms and/or breaking terms into multiple words. It is a `SpellCheckComponent` enhancement, leveraging Lucene's `WordBreakSpellChecker`. It can detect spelling errors resulting from misplaced whitespace without the use of shingle-based dictionaries and provides collation support for word-break errors, including cases where the user has a mix of single-word spelling errors and word-break errors in the same query. It also provides shard support.

Here is how it might be configured in `solrconfig.xml`:

```
<searchComponent name="spellcheck" class="solr.SpellCheckComponent">
  <lst name="spellchecker">
    <str name="name">wordbreak</str>
    <str name="classname">solr.WordBreakSolrSpellChecker</str>
    <str name="field">lowerfilt</str>
    <str name="combineWords">true</str>
    <str name="breakWords">true</str>
    <int name="maxChanges">10</int>
  </lst>
</searchComponent>
```

Some of the parameters will be familiar from the discussion of the other spell checkers, such as `name`, `classname`, and `field`. New for this spell checker is `combineWords`, which defines whether words should be combined in a dictionary search (default is true); `breakWords`, which defines if words should be broken during a dictionary search (default is true); and `maxChanges`, an integer which defines how many times the spell checker should check collation possibilities against the index (default is 10).

The spellchecker can be configured with a traditional checker (ie: `DirectSolrSpellChecker`). The results are combined and collations can contain a mix of corrections from both spellcheckers.

Add It to a Request Handler

Queries will be sent to a [RequestHandler](#). If every request should generate a suggestion, then you would add the following to the `requestHandler` that you are using:

```
<str name="spellcheck">true</str>
```

One of the possible parameters is the `spellcheck.dictionary` to use, and multiples can be defined. With multiple dictionaries, all specified dictionaries are consulted and results are interleaved. Collations are created with combinations from the different spellcheckers, with care taken that multiple overlapping corrections do not occur in the same collation.

Here is an example with multiple dictionaries:

```
<requestHandler name="spellCheckWithWordbreak"
class="org.apache.solr.handler.component.SearchHandler">
<lst name="defaults">
<str name="spellcheck.dictionary">default</str>
<str name="spellcheck.dictionary">wordbreak</str>
<str name="spellcheck.count">20</str>
</lst>
<arr name="last-components">
<str>spellcheck</str>
</arr>
</requestHandler>
```

Spell Check Parameters

The SpellCheck component accepts the parameters described in the table below. All of these parameters can be overridden by specifying `spellcheck.collateParam.xx` where `xx` is the parameter you are overriding.

Parameter	Description
<code>spellcheck</code>	Turns on or off SpellCheck suggestions for the request. If true , then spelling suggestions will be generated.
<code>spellcheck.q</code> or <code>q</code>	Selects the query to be spellchecked.
<code>spellcheck.build</code>	Instructs Solr to build a dictionary for use in spellchecking.
<code>spellcheck.collate</code>	Causes Solr to build a new query based on the best suggestion for each term in the submitted query.
<code>spellcheck.maxCollations</code>	This parameter specifies the maximum number of collations to return.
<code>spellcheck.maxCollationTries</code>	This parameter specifies the number of collation possibilities for Solr to try before giving up.

<code>spellcheck.maxCollationEvaluations</code>	This parameter specifies the maximum number of word correction combinations to rank and evaluate prior to deciding which collation candidates to test against the index.
<code>spellcheck.collateExtendedResult</code>	If true, returns an expanded response detailing the collations found. If <code>spellcheck.collate</code> is false, this parameter will be ignored.
<code>spellcheck.count</code>	Specifies the maximum number of spelling suggestions to be returned.
<code>spellcheck.dictionary</code>	Specifies the dictionary that should be used for spellchecking.
<code>spellcheck.extendedResults</code>	Causes Solr to return additional information about spellcheck results, such as the frequency of each original term in the index (<code>origFreq</code>) as well as the frequency of each suggestion in the index (<code>frequency</code>). Note that this result format differs from the non-extended one as the returned suggestion for a word is actually an array of lists, where each list holds the suggested term and its frequency.
<code>spellcheck.onlyMorePopular</code>	Limits spellcheck responses to queries that are more popular than the original query.
<code>spellcheck.maxResultsForSuggest</code>	The maximum number of hits the request can return in order to both generate spelling suggestions and set the "correctlySpelled" element to "false".
<code>spellcheck.alternativeTermCount</code>	The count of suggestions to return for each query term existing in the index and/or dictionary.
<code>spellcheck.reload</code>	Reloads the spellchecker.
<code>spellcheck.accuracy</code>	Specifies an accuracy value to help decide whether a result is worthwhile.
<code>spellcheck.<DICT_NAME>.key</code>	Specifies a key/value pair for the implementation handling a given dictionary.

The `spellcheck` Parameter

This parameter turns on SpellCheck suggestions for the request. If `true`, then spelling suggestions will be generated.

The `spellcheck.q` or `q` Parameter

This parameter specifies the query to spellcheck. If `spellcheck.q` is defined, then it is used; otherwise the original input query is used. The `spellcheck.q` parameter is intended to be the original query, minus any extra markup like field names, boosts, and so on. If the `q` parameter is specified, then the `SpellingQueryConverter` class is used to parse it into tokens; otherwise the `WhitespaceTokenizer` is used. The choice of which one to use is up to the application. Essentially, if you have a spelling "ready" version in your application, then it is probably better to use `spellcheck.q`. Otherwise, if you just want Solr to do the job, use the `q` parameter.

 The `SpellingQueryConverter` class does not deal properly with non-ASCII characters. In this case, you have either to use `spellcheck.q`, or implement your own `QueryConverter`.

The `spellcheck.build` Parameter

If set to `true`, this parameter creates the dictionary that the `SolrSpellChecker` will use for spell-checking. In a typical search application, you will need to build the dictionary before using the `SolrSpellChecker`. However, it's not always necessary to build a dictionary first. For example, you can configure the spellchecker to use a dictionary that already exists.

The dictionary will take some time to build, so this parameter should not be sent with every request.

The `spellcheck.reload` Parameter

If set to `true`, this parameter reloads the spellchecker. The results depend on the implementation of `SolrSpellChecker.reload()`. In a typical implementation, reloading the spellchecker means reloading the dictionary.

The `spellcheck.count` Parameter

This parameter specifies the maximum number of suggestions that the spellchecker should return for a term. If this parameter isn't set, the value defaults to 1. If the parameter is set but not assigned a number, the value defaults to 5. If the parameter is set to a positive integer, that number becomes the maximum number of suggestions returned by the spellchecker.

The `spellcheck.onlyMorePopular` Parameter

If `true`, Solr will return suggestions that result in more hits for the query than the existing query. Note that this will return more popular suggestions even when the given query term is present in the index and considered "correct".

The `spellcheck.maxResultsForSuggest` Parameter

For example, if this is set to 5 and the user's query returns 5 or fewer results, the spellchecker will report "correctlySpelled=false" and also offer suggestions (and collations if requested). Setting this greater than zero is useful for creating "did-you-mean?" suggestions for queries that return a low number of hits.

The `spellcheck.alternativeTermCount` Parameter

Specify the number of suggestions to return for each query term existing in the index and/or dictionary. Presumably, users will want fewer suggestions for words with `docFrequency>0`. Also setting this value turns "on" context-sensitive spell suggestions.

The `spellcheck.extendedResults` Parameter

This parameter causes Solr to include additional information about the suggestion, such as the frequency in the index.

The `spellcheck.collate` Parameter

If `true`, this parameter directs Solr to take the best suggestion for each token (if one exists) and construct a new query from the suggestions. For example, if the input query was "jawa class lording" and the best suggestion for "jawa" was "java" and "lording" was "loading", then the resulting collation would be "java class loading".

The `spellcheck.collate` parameter only returns collations that are guaranteed to result in hits if re-queried, even when applying original `fq` parameters. This is especially helpful when there is more than one correction per query.



This only returns a query to be used. It does not actually run the suggested query.

The `spellcheck.maxCollations` Parameter

The maximum number of collations to return. The default is `1`. This parameter is ignored if `spellcheck.collate` is false.

The `spellcheck.maxCollationTries` Parameter

This parameter specifies the number of collation possibilities for Solr to try before giving up. Lower values ensure better performance. Higher values may be necessary to find a collation that can return results. The default value is `0`, which maintains backwards-compatible (Solr 1.4) behavior (do not check collations). This parameter is ignored if `spellcheck.collate` is false.

The `spellcheck.maxCollationEvaluations` Parameter

This parameter specifies the maximum number of word correction combinations to rank and evaluate prior to deciding which collation candidates to test against the index. This is a performance safety-net in case a user enters a query with many misspelled words. The default is **10,000** combinations, which should work well in most situations.

The **spellcheck.collateExtendedResult** Parameter

If **true**, this parameter returns an expanded response format detailing the collations Solr found. The default value is **false** and this is ignored if `spellcheck.collate` is false.

The **spellcheck.dictionary** Parameter

This parameter causes Solr to use the dictionary named in the parameter's argument. The default setting is "default". This parameter can be used to invoke a specific spellchecker on a per request basis.

The **spellcheck.accuracy** Parameter

Specifies an accuracy value to be used by the spell checking implementation to decide whether a result is worthwhile or not. The value is a float between 0 and 1. Defaults to `Float.MIN_VALUE`.

The **spellcheck.<DICT_NAME>.key** Parameter

Specifies a key/value pair for the implementation handling a given dictionary. The value that is passed through is just `key=value` (`spellcheck.<DICT_NAME>.` is stripped off).

For example, given a dictionary called `foo`, `spellcheck.foo.myKey=myValue` would result in `myKey=myValue` being passed through to the implementation handling the dictionary `foo`.

Example

This example shows the results of a simple query that defines a query using the `spellcheck.q` parameter. The query also includes a `spellcheck.build=true` parameter, which is needs to be called only once in order to build the index. `spellcheck.build` should not be specified with for each request.

```
http://localhost:8983/solr/spellCheckCompRH?q=*&spellcheck.q=hell%20ultrashar&spel
```

Results:

```

<lst name="spellcheck">
    <lst name="suggestions">
        <lst name="hell">
            <int name="numFound">1</int>
            <int name="startOffset">0</int>
            <int name="endOffset">4</int>
            <arr name="suggestion">
                <str>dell</str>
            </arr>
        </lst>
        <lst name="ultrashar">
            <int name="numFound">1</int>
            <int name="startOffset">5</int>
            <int name="endOffset">14</int>
            <arr name="suggestion">
                <str>ultrasharp</str>
            </arr>
        </lst>
    </lst>
</lst>

```

Distributed SpellCheck

The `SpellCheckComponent` also supports spellchecking on distributed indexes. If you are using the `SpellCheckComponent` on a request handler other than `"/select"`, you must provide the following two parameters:

Parameter	Description
<code>shards</code>	Specifies the shards in your distributed indexing configuration. For more information about distributed indexing, see Distributed Search with Index Sharding
<code>shards.qt</code>	Specifies the request handler Solr uses for requests to shards. This parameter is not required for the <code>/select</code> request handler.

For example:

`http://localhost:8983/solr/select?q=*&spellcheck=true&spellcheck.build=true&spellcheck.count=5`

In case of a distributed request to the `SpellCheckComponent`, the shards are requested for at least five suggestions even if the `spellcheck.count` parameter value is less than five. Once the suggestions are collected, they are ranked by the configured distance measure (Levenshtein Distance by default) and then by aggregate frequency.

Suggester

Solr includes an autosuggest component called Suggester, which is built on the [SpellCheck search component](#). The autocomplete suggestions that Suggester provides come from a dictionary that is either based on the main index or on a dictionary file that you provide. It is common to provide only the top-N suggestions, either ranked alphabetically or according to their usefulness for an average user (such as popularity or the number of returned results).

Because this feature is based on the [SpellCheck search component](#), configuring Suggester is similar to configuring spell checking. Unlike the SpellCheck Component, however, Suggester has no direct indexing option at this time.

In `solrconfig.xml`, we need to add a search component and a request handler.

Covered in this section:

- [Adding the Suggest Search Component](#)
- [Adding the Suggest Request Handler](#)
- [Defining a Field for Suggester](#)
- [Related Topics](#)

Adding the Suggest Search Component

The first step is to add a search component to `solrconfig.xml` to extend the SpellChecker. Here is some sample code that could be used.

```
<searchComponent class="solr.SpellCheckComponent" name="suggest">
  <lst name="spellchecker">
    <str name="name">suggest</str>
    <str name="classname">org.apache.solr.spelling.suggest.Suggerster</str>
    <str name="lookupImpl">org.apache.solr.spelling.suggest.tst.TSTLookup</str>
    <str name="field">name</str>  <!-- the indexed field to derive suggestions from
-->
    <float name="threshold">0.005</float>
    <str name="buildOnCommit">true</str>
  <!--
    <str name="sourceLocation">american-english</str>
-->
  </lst>
</searchComponent>
```

One of the most important parameters is the `lookupImpl`, which is described in more detail below. In this example, the `sourceLocation` is commented out, which means that a dictionary file will not be used. Instead, the field defined with the `field` parameter will be used as the dictionary. We've included the unused `sourceLocation` in the example to demonstrate it's usage.

Suggerster Search Component Parameters

The Suggerster search component takes the following configuration parameters:

Parameter	Description
searchComponent name	Arbitrary name for the search component.
name	A symbolic name for this spellchecker. You can refer to this name in the URL parameters and in the SearchHandler configuration.
classname	The full class name of the component: <code>org.apache.solr.spelling.Suggerster</code>

lookupImpl	<p>Lookup implementation. Choose one of these four:</p> <p><code>org.apache.solr.suggest.fst.FSTLookup</code>: automaton-based lookup. This implementation is slower to build, but provides the lowest memory cost. We recommend using this implementation unless you need more sophisticated matching results, in which case you should use the Jaspell implementation.</p> <p><code>org.apache.solr.suggest.wfst.WFSTLookup</code>: weighted automaton representation; an alternative to FSTLookup for more fine-grained ranking. WFSTLookup does not use buckets, but instead a shortest path algorithm. Note that it expects weights to be whole numbers. If weight is missing it's assumed to be 1.0. Weights affect the sorting of matching suggestions when <code>spellcheck.onlyMorePopular=true</code> is selected: weights are treated as "popularity" score, with higher weights preferred over suggestions with lower weights.</p> <p><code>org.apache.solr.suggest.jaspell.JaspellLookup</code>: a more complex lookup based on a ternary trie from the JaSpell project. Use this implementation if you need more sophisticated matching results.</p> <p><code>org.apache.solr.suggest.tst.TSTLookup</code>: a simple compact ternary trie based lookup.</p> <p>All four implementations will likely run at similar speed when requests are made through HTTP. Direct benchmarks of these classes indicate that FSTLookup provides better performance compared to the other three methods, and at a much lower memory cost. We recommend using the FSTLookup implementation unless you need more sophisticated matching, in which case you should use the JaspellLookup implementation or FSTLookupFactory.</p>
------------	--

buildOnCommit or buildOnOptimize	False by default. If true then the lookup data structure will be rebuilt after commit. If false , then the lookup data will be built only when requested by URL parameter <code>spellcheck.build=true</code> . Use <code>buildOnCommit</code> to rebuild the dictionary with every commit, or <code>buildOnOptimize</code> to build the dictionary only when the index is optimized.
	<p> Currently implemented lookups keep their data in memory, so unlike spellchecker data, this data is discarded on core reload and not available until you invoke the build command, either explicitly or implicitly during a commit.</p>
queryConverter	Allows defining an alternate converter that can parse phrases in dictionary files. It passes the whole string to the query analyzer rather than analyzing it for spelling. Define it in <code>solrconfig.xml</code> as <code><queryConverter name="queryConverter" class="org.apache.solr.spelling.SuggestQueryConverter" /></code> .
sourceLocation	The path to the dictionary file. If this value is empty then the main index will be used as a source of terms and weights.
field	If <code>sourceLocation</code> is empty then terms from this field in the index will be used when building the trie.
threshold	<p>A value between zero and one representing the minimum fraction of the total documents where a term should appear in order to be added to the lookup dictionary.</p> <p>When you use the index as the dictionary, you may encounter many invalid or uncommon terms. The <code>threshold</code> parameter addresses this issue. By setting the <code>threshold</code> parameter to a value just above zero, you can greatly reduce the number of unusable terms in your dictionary while maintaining most of the common terms. The example above sets the <code>threshold</code> value to 0.5%. The <code>threshold</code> parameter does not affect file-based dictionaries.</p>

Using a Dictionary File

If using a dictionary file, it should be a plain text file in UTF-8 encoding. Blank lines and lines that start with a '#' are ignored. The remaining lines must consist of either a string without literal TAB (\u0007) characters, or a string and a TAB separated floating-point weight. You can use both single terms and phrases in a dictionary file.

```
# This is a sample dictionary file.

acquire
accidentally\t2.0
accommodate\t3.0
```

Adding the Suggest Request Handler

After adding the search component, a request handler must be added to `solrconfig.xml`. This request handler will set a number of parameters for serving suggestion requests and incorporate the "suggest" search component defined in the previous step. Because the Suggester is based on the SpellCheckComponent, the request handler shares many of the same parameters.

```
<requestHandler class="org.apache.solr.handler.component.SearchHandler"
name="/suggest">
  <lst name="defaults">
    <str name="spellcheck">true</str>
    <str name="spellcheck.dictionary">suggest</str>
    <str name="spellcheck.onlyMorePopular">true</str>
    <str name="spellcheck.count">5</str>
    <str name="spellcheck.collate">true</str>
  </lst>
  <arr name="components">
    <str>suggest</str>
  </arr>
</requestHandler>
```

Suggester Request Handler Parameters

The Suggester request handler takes the following configuration parameters:

Parameter	Description
<code>spellcheck=true</code>	This parameter should always be true, because we always want to run the Suggester for queries submitted to this handler.
<code>spellcheck.dictionary</code>	The name of the dictionary component configured in the search component.
<code>spellcheck.onlyMorePopular</code>	If true, then suggestions will be sorted by weight ("popularity"), which is the recommended setting. The <code>count</code> parameter will effectively limit this to a top-N list of best suggestions. If false, suggestions are sorted alphabetically.

spellcheck.count	Specifies the number of suggestions for Solr to return.
spellcheck.collate	If true, Solr provides a query collated with the first matching suggestion.

Defining a Field for Suggester

Any field can be used as the basis of the dictionary (if not using an explicit dictionary file). You may want to create a custom field for this purpose, and use the copy fields feature to copy text from various fields to the dedicated "suggester" field.

```
<field indexed="true" multiValued="true" name="suggestions" stored="false"
type="textSpell"/>
```

You may want to define a custom fieldType in schema.xml to prevent over-analysis of the content of a field for use in suggestions. For example, if you have some analysis that stems terms, you wouldn't want the stemmed terms in the suggestion list, since the stemmed forms of words would be presented to users. Here is an example that could be used:

```
<fieldType class="solr.TextField" name="textSpell" positionIncrementGap="100">
<analyzer>
    <tokenizer class="solr.StandardTokenizerFactory"/>
    <filter class="solr.StandardFilterFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
</analyzer>
</fieldType>
```

Related Topics

- [RequestHandlers and SearchComponents in SolrConfig](#)
- [Solr Field Types](#)
- [Copying Fields](#)

Function Queries

Function queries enable you to generate a relevancy score using the actual value of one or more numeric fields. Function queries are supported by the [DisMax](#), [Extended DisMax](#), and [standard](#) query parsers.

Function queries use *functions*. The functions can be a constant (numeric or string literal), a field, another function or a parameter substitution argument. You can use these functions to modify the ranking of results for users. These could be used to change the ranking of results based on a user's location, or some other calculation.

Function query topics covered in this section:

- [Using Function Query](#)
- [Available Functions](#)
- [Example Function Queries](#)
- [Sort By Function](#)
- [Related Topics](#)

Using Function Query

Functions must be expressed as function calls (for example, `sum(a,b)` instead of simply `a+b`).

There are three principal ways of including function queries in a Solr query:

- Introduce a function query with the `_val_` keyword. For example:

```
_val_:mynumericfield _val_:"recip(rord(myfield),1,2,3)"
```

- Use a parameter that has an explicit type of `FunctionQuery`, such as the `DisMax` query parser's [bf \(boost function\) parameter](#). Note that the `bf` parameter actually takes a list of function queries separated by white space and each with an optional boost. Make sure you eliminate any internal white space in single function queries when using `bf`. For example:

```
q=dismax&bf="ord(popularity)^0.5 recip(rord(price),1,1000,1000)^0.3"
```

- Add the results of FunctionQueries as fields to a document. For instance, for:

```
&fl=sum(x, y),id,a,b,c,score
```

the output would be:

```
...
<str name="id">foo</str>
<float name="sum(x,y)">40</float>
<float name="score">0.343</float>
...
```

Only functions with fast random access are recommended.

Available Functions

The table below summarizes the functions available for function queries.

Function	Description	Syntax Examples
abs	Returns the absolute value of the specified value or function.	abs(x) abs(-5)
and	Returns a value of true if and only if both of its operands are true.	and("blue" , "sky") is true only if the sky is blue.
constant	Specifies a floating point constant.	1.5 _val_:1.5
def	def is short for default. Returns the value of field "field", or if the field does not exist, returns the default value specified. and yields the first value where exists() == true .)	def(rating,5): This def() function returns the rating or if no rating specified in the doc, returns 5 def(myfield, 1.0): equivalent to if(exists(myfield),myfield,1.0)
div	Divides one value or function by another. div(x,y) divides x by y.	div(1,y) div(sum(x,100),max(y,1))

dist	<p>Return the distance between two vectors (points) in an n-dimensional space. Takes in the power, plus two or more ValueSource instances and calculates the distances between the two vectors. Each ValueSource must be a number. There must be an even number of ValueSource instances passed in and the method assumes that the first half represent the first vector and the second half represent the second vector.</p>	<p><code>dist(2, x, y, 0, 0)</code>: calculates the Euclidean distance between (0,0) and (x,y) for each document</p> <p><code>dist(1, x, y, 0, 0)</code>: calculates the Manhattan (taxicab) distance between (0,0) and (x,y) for each document</p> <p><code>dist(2, x,y,z,0,0,0)</code>: Euclidean distance between (0,0,0) and (x,y,z) for each document.</p> <p><code>dist(1,x,y,z,e,f,g)</code>: Euclidean distance between (x,y,z) and (e,f,g) where each letter is a field name</p>
docfreq(field,val)	<p>Returns the number of documents that contain the term in the field. This is a constant (the same value for all documents in the index).</p> <p>You can quote the term if it's more complex, or do parameter substitution for the term value.</p>	<pre>docfreq(text,solr) docfreq(text,'solr') & defType=func&q=docfreq(text,\$myterm)&myterm=sc</pre>
exists	Returns TRUE if any member of the field exists.	<p><code>exists(author)</code> returns TRUE for any document has a value in the "author" field.</p> <p><code>exists(query(price:5.00))</code> returns TRUE if "price" contains "5.00".</p>

field	Returns the numeric field value of an indexed (not multi-valued) field with a maximum of one value per document. The syntax is simply the field name by itself. 0 is returned for documents without a value in the field.	<code>_val_:myFloatField</code>
hsin	The Haversine distance calculates the distance between two points on a sphere when traveling along the sphere. The values must be in radians. <code>hsin</code> also takes a Boolean argument to specify whether the function should convert its output to radians.	<code>hsin(2, true, x, y, 0, 0)</code>
fl	Defines the fields of a document that are returned in response to a query. You need to specify <code>score</code> as one of the fields to return if you want any calculated values to be returned in the score.	<code>fl=*,score</code>

idf	<p>Inverse document frequency; a measure of whether the term is common or rare across all documents.</p> <p>Obtained by dividing the total number of documents by the number of documents containing the term, and then taking the logarithm of that quotient. See also tf.</p>	<p>idf(list1,326.55): measures the inverse of the frequency of the occurrence of the number 326.55 in list1.</p>
if	<p>Enables conditional function queries.</p> <p>IF(test; value1; value2) where:</p>	<pre>if(color=="red"; 100; if(color=="green"; 50; 25)) :</pre> <p>This function checks the document field "color", and if is "red" returns 100, if it is "green" returns 50, else returns 25.</p>

- test is or refers to a logical value or expression that returns a logical value (TRUE or FALSE).
- value1 is the value that is returned by the function if test yields TRUE.
- value2 is the value that is returned by the function if test yields FALSE.
- If value2 is omitted it is assumed to be FALSE; if value1 is also omitted it is assumed to be TRUE.
An expression can be any function which outputs boolean values, or even functions returning numeric values, in which case value 0 will be interpreted as false, or strings, in which case empty string is interpreted as false.

linear	Implements $m*x+c$ where m and c are constants and x is an arbitrary function. This is equivalent to <code>sum(product(m,x),c)</code> , but slightly more efficient as it is implemented as a single function.	<code>linear(x,m,c)</code> <code>linear(x,2,4) returns 2*x+4</code>
log	Returns the log base 10 of the specified function.	<code>log(x)log(sum(x,100))</code>
map	Maps any values of the function x that fall within min and max inclusive to the specified target. The arguments $\text{min}, \text{max}, \text{target}$ are constants. The function outputs the field's value if it does not fall between min and max .	<code>map(x,min,max,target)</code> <code>map(x,0,0,1) - changes any values of 0 to 1. This can be useful in handling default 0 values.</code> <code>map(x,min,max,target,altarg)</code> <code>map(x,0,0,1,0) - changes any values of 0 to 1 and if the value is not zero it can be set to the value of the 5 argument instead of defaulting to the field's value.</code>
max	Returns the max of another function and a constant, which are specified as arguments: <code>max(x,c)</code> . The max function is useful for "bottoming out" another function at some constant.	<code>max(myfield,0)</code>

maxdoc	Returns the number of documents in the index, including those that are marked as deleted but have not yet been purged. This is a constant (the same value for all documents in the index).	maxdoc(list1)
ms	Returns milliseconds of difference between its arguments. Dates are relative to the Unix or POSIX time epoch, midnight, January 1, 1970 UTC. Arguments may be numerically indexed date fields such as TrieDate (the default in 1.4), or date math based on a constant date or NOW. ms() : Equivalent to ms(NOW), number of milliseconds since the epoch. ms(a) : Returns the number of milliseconds since the epoch that the argument represents.	ms(NOW/DAY) ms(2000-01-01T00:00:00Z) ms(mydatefield) ms(a,b) : Returns the number of milliseconds that b occurs before a (that is, a - b) Examples: ms(NOW,mydatefield) ms(mydatefield,2000-01-01T00:00:00Z) ms(datefield1,datefield2)
norm(<i>field</i>)	Returns the "norm" stored in the index, the product of the index time boost and the length normalization factor, according to the Similarity for the field.	norm(text)
not	A logically negative value.	if (NOT value1) [...]: TRUE only when value1 is false.

numdocs	Returns the number of documents in the index, not including those that are marked as deleted but have not yet been purged. This is a constant (the same value for all documents in the index).	numdocs(list1)
or	A logical disjunction.	(value1 OR value2): TRUE if either value1 or value is true.

ord	<p>Returns the ordinal of the indexed field value within the indexed list of terms for that field in Lucene index order (lexicographically ordered by unicode value), starting at 1. In other words, for a given field, all values are ordered lexicographically; this function then returns the offset of a particular value in that ordering. The field must have a maximum of one value per document (not multi-valued). 0 is returned for documents without a value in the field.</p> <div style="border: 1px solid orange; padding: 10px; background-color: #fffacd;">  ord() depends on the position in an index and can change when other documents are inserted or deleted. </div> <p>See also <code>rord</code> below.</p>	<pre>ord(myIndexedField) _val_: "ord(myIndexedField)"</pre> <p>Example: If there were only three values ("apple", "banana", "pear") for a particular field, then: <code>ord("apple")=1</code> <code>ord("banana")=2</code> <code>ord("pear")=3</code></p>
pow	<p>Raises the specified base to the specified power. <code>pow(x,y)</code> raises x to the power of y.</p>	<pre>pow(x,y) pow(x,log(y)) pow(x,0.5): the same as sqrt</pre>

product	Returns the product of multiple values or functions, which are specified in a comma-separated list.	<code>product(x,y,...)</code> <code>product(x,2)</code> <code>product(x,y)</code>
query	Returns the score for the given subquery, or the default value for documents not matching the query. Any type of subquery is supported through either parameter de-referencing <code>\$otherparam</code> or direct specification of the query string in the Local Parameters through the <code>v</code> key.	<code>query(subquery, default)</code> <code>q=product(popularity, query({!dismax v='solr rocks'})}</code> : returns the product of the popularity and the score of the DisMax query. <code>q=product(popularity, query(\$qq))&qq={!dismax}solr rocks</code> : equivalent to the previous query, using parameter de-referencing. <code>q=product(popularity, query(\$qq,0.1))&qq={!dismax}solr rocks</code> : specifies a default score of 0.1 for documents that don't match the DisMax query.

recip	<p>Performs a reciprocal function with <code>recip(myfield,m,a,b)</code> implementing $a/(m*x+b)$. m, a, b are constants, and x is any arbitrarily complex function.</p> <p>When a and b are equal, and $x \geq 0$, this function has a maximum value of 1 that drops as x increases. Increasing the value of a and b together results in a movement of the entire function to a flatter part of the curve.</p> <p>These properties can make this an ideal function for boosting more recent documents when x is <code>rord(datefield)</code>.</p>	<pre>recip(myfield,m,a,b) recip(rord(creationDate),1,1000,1000)</pre>
rord	Returns the reverse ordering of that returned by <code>ord</code> .	<pre>rord(myDateField) _val_: "rord(myDateField)"</pre> <p><code>rord(myDateField)</code>: a metric for how old a document is. The youngest document will return 1. The oldest document will return the total number of documents.</p>

scale	<p>Scales values of the function <code>x</code> such that they fall between the specified <code>minTarget</code> and <code>maxTarget</code> inclusive. The current implementation traverses all of the function values to obtain the min and max, so it can pick the correct scale.</p> <p>The current implementation cannot distinguish when documents have been deleted or documents that have no value. It uses 0.0 values for these cases. This means that if values are normally all greater than 0.0, one can still end up with 0.0 as the min value to map from. In these cases, an appropriate <code>map()</code> function could be used as a workaround to change 0.0 to a value in the real range, as shown here:</p> <pre>scale(map(x,0,0,5),1,2)</pre>	<code>scale(x,minTarget,maxTarget)</code> <code>scale(x,1,2)</code> : scales the values of <code>x</code> such that all values will be between 1 and 2 inclusive.
-------	--	--

sqedist	The Square Euclidean distance calculates the 2-norm (Euclidean distance) but does not take the square root, thus saving a fairly expensive operation. It is often the case that applications that care about Euclidean distance do not need the actual distance, but instead can use the square of the distance. There must be an even number of ValueSource instances passed in and the method assumes that the first half represent the first vector and the second half represent the second vector.	sqedist(x_td, y_td, 0, 0)
sqrt	Returns the square root of the specified value or function.	sqrt(x)sqrt(100)sqrt(sum(x,100))

strdist	<p>Calculate the distance between two strings. Uses the Lucene spell checker StringDistance interface and supports all of the implementations available in that package, plus allows applications to plug in their own via Solr's resource loading capabilities. strdist takes (string1, string2, distance measure). Possible values for distance measure are:</p> <p>jw: Jaro-Winkler</p> <p>edit: Levenshtein or Edit distance</p> <p>ngram: The NGramDistance, if specified, can optionally pass in the ngram size too. Default is 2.</p> <p>FQN: Fully Qualified class Name for an implementation of the StringDistance interface. Must have a no-arg constructor.</p>	strdist("SOLR",id,edit)
sub	Returns x-y from sub(x,y).	sub(myfield,myfield2) sub(100,sqrt(myfield))

sum	Returns the sum of multiple values or functions, which are specified in a comma-separated list.	<code>sum(x,y,...)</code> <code>sum(x,1)</code> <code>sum(x,y)</code> <code>sum(sqrt(x),log(y),z,0.5)</code>
sumtotaltermfreq	Returns the sum of <code>totaltermfreq</code> values for all terms in the field in the entire index (i.e., the number of indexed tokens for that field). (Aliases <code>sumtotaltermfreq</code> to <code>sttf</code> .)	If <code>doc1:(fieldX:A B C)</code> and <code>doc2:(fieldX:A A A A)</code> : <code>docFreq(fieldX:A) = 2</code> (A appears in 2 docs) <code>freq(doc1, fieldX:A) = 4</code> (A appears 4 times in doc 2) <code>totalTermFreq(fieldX:A) = 5</code> (A appears 5 times across all docs) <code>sumTotalTermFreq(fieldX) = 7</code> in fieldX, there are 1 As, 1 B, 1 C
termfreq	Returns the number of times the term appears in the field for that document.	<code>termfreq(text, 'memory')</code>
tf	Term frequency; returns the term frequency factor for the given term, using the Similarity for the field. The <code>tf-idf</code> value increases proportionally to the number of times a word appears in the document, but is offset by the frequency of the word in the document, which helps to control for the fact that some words are generally more common than others. See also <code>idf</code> .	<code>tf(text, 'solr')</code>

top	<p>Causes the function query argument to derive its values from the top-level IndexReader containing all parts of an index.</p> <p>For example, the ordinal of a value in a single segment will be different from the ordinal of that same value in the complete index.</p> <p>The <code>ord()</code> and <code>rord()</code> functions implicitly use <code>top()</code>, and hence <code>ord(foo)</code> is equivalent to <code>top(ord(foo))</code>.</p>	
Totaltermfreq	Returns the number of times the term appears in the field in the entire index. (Aliases <code>totaltermfreq</code> to <code>ttf</code> .)	<code>ttf(text, 'memory')</code>
xor()	Logical exclusive disjunction, or one or the other but not both.	<code>xor(field1,field2)</code> returns TRUE if either <code>field1</code> or <code>field2</code> is true; FALSE if both are true.

Example Function Queries

To give you a better understanding of how function queries can be used in Solr, suppose an index stores the dimensions in meters `x,y,z` of some hypothetical boxes with arbitrary names stored in field `boxname`. Suppose we want to search for box matching name `findbox` but ranked according to volumes of boxes. The query parameters would be:

```
q=boxname:findbox_val_:"product(product(x,y),z)
```

This query will rank the results based on volumes. In order to get the computed volume, you will need to request the `score`, which will contain the resultant volume:

```
&fl=*, score
```

Suppose that you also have a field storing the weight of the box as `weight`. To sort by the density of the box and return the value of the density in score, you would submit the following query:

```
http://localhost:8983/solr/select/?q=boxname:findbox_val_div(weight,product(product(x,y),z,x y z weight score)
```

Sort By Function

You can sort your query results by the output of a function. For example, to sort results by distance, you could enter:

```
http://localhost:8983/solr/select?q=*:*&sort=dist(2, point1, point2) desc
```

Sort by function also supports pseudo-fields: fields can be generated dynamically and return results as though it was normal field in the index. For example,

`&fl=id,sum(x, y),score`

would return:

```
<str name="id">foo</str>
<float name="sum(x,y)">40</float>
<float name="score">0.343</float>
```

Related Topics

- [FunctionQuery](#)

Spatial Search

Solr supports location data for use in spatial or geospatial searches. Using spatial search, you can:

- Represent spatial data in the index
- Filter by location based on a bounding box or circle
- Sort by distance
- Score and boost by distance

With Solr 4, there are two ways to work with spatial searches. The approach introduced in Solr 3.x remains available for use, but a new approach based on a new module in Lucene 4 has also been introduced.

Solr 3 Spatial

There are three function queries that support spatial search: `dist`, to determine the distance between two points; `hdist`, to calculate the distance between two points on a sphere; and `sqdist`, to calculate the square Euclidean distance between two points. For more information about these function queries, see the section on [Function Queries](#).

For more information on Solr spatial search, see <http://wiki.apache.org/solr/SpatialSearch>.

Spatial Search Features

Solr includes three useful tools for working with spatial queries: `geofilt`, a geospatial filter; `bbox`, a geospatial bounding-box filter; and `geodist`, a geospatial distance function.

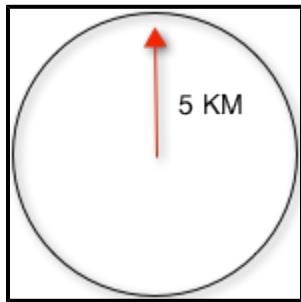
Spatial Search Parameters

The following parameters are used for spatial search:

Parameter	Description
d	distance, in kilometers
pt	a lat/lon coordinate point
sfield	a spatial field, by default a location (lat/lon) field type.

geofilt

The `geofilt` filter allows you to retrieve results based on the distance from a given point. For example, to find all results for a product search within five kilometers of the lat/lon point, you could enter `&q=*:&fq={!geofilt sfield=store}&pt=45.15,-93.85&d=5`. This filter returns all results within a circle of the given radius around the initial point:



bbox

`bbox` allows you to filter results based on a specified area around a given point. `bbox` takes the same parameters as `geofilt`, but rather than calculating all points in a circle within the given radius from the initial point, it only calculates the lower left and upper right corners of a square that would enclose a circle with the given radius. To return all results within five kilometers of a give point, you could enter `...&q=:&fq={!bbox sfield=store}&pt=45.15,-93.85&d=5`. The resulting bounding box would encompass all points within a five kilometer circle around the initial point, but it would also include some extra points in the corners of the bounding box that fall outside the five kilometer radius. Bounding box filters therefore can return results that fall outside your desired parameters, but they are much less "expensive" to implement.



- ⚠ When a bounding box includes a pole, the `location` field type produces a "bounding bowl" (a spherical cap) that includes all values that are north or south of the latitude of the bounding box corner (the lower left and the upper right) that is closer to the equator. In other words, Solr still calculates what the coordinates of the upper right corner and the lower left corner of the box would be just as in all other filtering cases, but it then take the corner that is closest to the equator (since it goes over the pole it may not be the lower left, despite the name) and filters by latitude only. This returns more matches than a pure bounding box match, but the query is both faster and easier to construct.

geodist

geodist is a distance function that takes three optional parameters: (sfield,latitude,longitude). You can use the geodist function to sort results by distance or score return results.

For example, to sort your results by ascending distance, enter

```
...&q=*:&fq={!geofilt}&sfield=store&pt=45.15,-93.85&d=50&sort=geodist asc.
```

To return the distance as the document score, enter

```
...&q={!func}geodist()&sfield=store&pt=45.15,-93.85&sort=score+asc.
```

Post filtering

Post filtering is an option available for spatial queries qualifying bbox and geofilt with LatLonType, which specifies latitude and longitude. LatLonType is passed as numbers in the query, as shown in Example 1 below.

Filtering is usually done in parallel with or before the main query. Post filters are applied after the main query. This is important when the filter itself is very time-consuming, so it's better to always apply it to matching documents instead of all documents.

More Examples

Here are a few more useful examples of what you can do with spatial search in Solr.

Use as a Sub-Query to Expand Search Results

Here we will query for results in Jacksonville, Florida, or within 50 kilometers of 45.15,-93.85 (near Buffalo, Minnesota):

```
&q=*:&fq=(state:"FL" AND city:"Jacksonville") OR  
_query_:"{!geofilt}"&sfield=store&pt=45.15,-93.85&d=50&sort=geodist()&asc
```

Facet by Distance

To facet by distance, use the Frange query parser:

```
&q=*:&sfield=store&pt=45.15,-93.85&facet.query={!frange l=0  
u=5}geodist()&facet.query={!frange l=5.001 u=3000}geodist()
```

Boost Nearest Results

Using the [DisMax](#) or [Extended DisMax](#), you can combine spatial search with the boost function to boost the nearest results:

```
&q.alt=*:&fq={!geofilt}&sfield=store&pt=45.15,-93.85&d=50&bf=recip(geodist(),2,200,  
desc
```

Solr 4 Spatial

The new approach to spatial offers several new features and improvements over the former approach:

- New shapes: polygons, line strings, and other new shapes
- Multi-valued indexed fields
- Ability to index non-point shapes as well as point shapes
- Rectangles with user-specified corners. The Solr 3 approach only supports bounding box of a circle
- Multi-value distance sort and score boosting
- Well-Known-Text support when JTS is used (for polygons, etc.)

The new approach incorporates the basic features of the Solr 3 approach, such as lat-lon bounding boxes and circles.

Solr 4 Spatial Field Type

The first step to using Solr 4 spatial is to register a field type in `schema.xml`. There are several options for this field type.

Setting	Description
name	The name of the field type.
class	For most use cases, using <code>solr.SpatialRecursivePrefixTreeFieldType</code> will be sufficient. Since the new spatial module in Lucene is meant to be a framework for different spatial "strategies", another class may be used. See the Spatial4J project for more information.
spatialContextFactory	If polygons or other shapes beyond a point, rectangle or circle are used, the JTS Topology Suite is a required dependency. If you intend to use those shapes, defined the class here.
units	This is required, and currently can only be "degrees".
distErrPct	Defines the default precision of non-point shapes, as a fraction between 0.0 (fully precise) to 0.5. The closer this number is to zero, the more accurate the shape will be. However, more precise indexed shapes use more disk space and take longer to index.
maxDistErr	Defines the highest level of detail required for indexed data. If left blank, the default is one meter, just a bit less than 0.000009 degrees.
geo	If true , the default, latitude and longitude coordinates will be based on WGS84 instead of Euclidean/Cartesian based. If false, the coordinates will be Euclidean/Cartesian-based.

worldBounds	Defines the valid numerical ranges for x and y, in the format of "minX minY maxX maxY". If geo=true, this is assumed "-180 -90 180 90". If geo=false, you should define your boundaries for non-geospatial uses.
distCalculator	Defines the distance calculation algorithm. If geo=true, "haversine" is the default. If geo=false, "cartesian" will be the default. Other possible values are "lawOfCosines", "vincentySphere" and "cartesian^2".
prefixTree	Defines the spatial grid implementation. Since a PrefixTree (such as RecursivePrefixTree) maps the world as a grid, each grid cell is decomposed to another set of grid cells at the next level. Using a "geohash" implementation, there are 32 children at each level. If geo=true, "geohash" is the only option. If geo=false, "quad" could be used for prefixTree, which has 4 children at each level.
maxLevels	Sets the maximum grid depth for indexed data. It may be simpler to use maxDistErr to calculate real distances.

```
<fieldType name="location_rpt" class="solr.SpatialRecursivePrefixTreeFieldType"
  spatialContextFactory="com.spatial4j.core.context.jts.JtsSpatialContextFactory"
  distErrPct="0.025"
  maxDistErr="0.000009"
  units="degrees" />
```

Once the field type has been defined, use it to define a field.

Because this functionality is quite new and some of it is still experimental, please review the Solr Wiki at <http://wiki.apache.org/solr/SolrAdaptersForLuceneSpatial4> for more information about searching and sorting location data.

- ✓ As of Solr 4.1, the new approach to spatial search will also work with the `{!geofilt}` and `{!bbox}` query parsers.

The Terms Component

The Terms Component provides access to the indexed terms in a field and the number of documents that match each term. This can be useful for building an auto-suggest feature or any other feature that operates at the term level instead of the search or document level. Retrieving terms in index order is very fast since the implementation directly uses Lucene's TermEnum to iterate over the term dictionary.

In a sense, this component provides fast field-faceting over the whole index, not restricted by the base query or any filters. The document frequencies returned are the number of documents that match the term, including any documents that have been marked for deletion but not yet removed from the index.

To use the Terms Component, users can pass in a variety of options to control what terms are returned. These parameters are:

Parameter	Description	Syntax
terms	If set to true, terms on the Terms Component. By default, the Terms Component is turned off.	terms={true false}
terms.fl	Specifies the field from which to retrieve terms.	terms.fl=field
terms.lower	Specifies the term at which to start. If not specified, the empty string is used, causing Solr to start at the beginning of the field.	terms.lower=term
terms.lower.incl	If set to true, includes the lower-bound term in the result set. By default, this parameter is set to true.	terms.lower.incl={true false}
terms.mincount	Specifies the minimum document frequency to return in order for a term to be included in a query response. Results are inclusive of the mincount (that is, \geq mincount). This parameter is optional.	terms.mincount=integer

terms.maxcount	Specifies the maximum document frequency a term must have in order to be included in a query response. The default setting is -1, which sets no upper bound. Results are inclusive of the maxcount (that is, <= maxcount). This parameter is optional.	terms.maxcount=integer
terms.prefix	Restricts matches to terms that begin with the specified string.	terms.prefix={string}
terms.limit	Specifies the maximum number of terms to return. The default is 10. If the limit is set to a number less than 0, then no maximum limit is enforced.	terms.limit=integer
terms.upper	Specifies the term to stop at. Any application using the Terms component must set either terms.limit or terms.upper.	terms.upper=upper_term
terms.upper.incl	If set to true, includes the upper bound term in the result set. The default is false.	terms.upper.incl={true false}
terms.raw	If set to true, returns the raw characters of the indexed term, regardless of whether it is human-readable. For instance, the indexed form of numeric numbers is not human-readable. The default is false.	terms.raw={true false}

The output is a list of the terms and their document frequency values.

Examples

The following examples use the sample Solr configuration located in the <Solr>/example directory.

The query below requests the first ten terms in the name field.

```
http://localhost:8983/solr/terms?terms.fl=name
```

Results:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<lst name="responseHeader">
<int name="status">0</int>
<int name="QTime">1</int>
</lst>
<lst name="terms">
<lst name="name">
<int name="0">5</int>
<int name="1">15</int>
<int name="11">5</int>
<int name="120">5</int>
<int name="133">5</int>
<int name="184">15</int>
<int name="19">5</int>
<int name="1900">5</int>
<int name="2">15</int>
<int name="20">5</int>
</lst>
</lst>
</response>
```

The query below requests the first ten terms in the name field, beginning with the first term that begins with the letter a.

<http://localhost:8983/solr/terms?terms.fl=name&terms.lower=a>

Results:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<lst name="responseHeader">
<int name="status">0</int>
<int name="QTime">2</int>
</lst>
<lst name="terms">
<lst name="name">
<int name="a">8</int>
<int name="adata">5</int>
<int name="all">5</int>
<int name="allinon">5</int>
<int name="amber">1</int>
<int name="appl">5</int>
<int name="asus">5</int>
<int name="ata">5</int>
<int name="ati">5</int>
<int name="b">5</int>
</lst>
</lst>
</response>
```

Using the Terms Component for an Auto-Suggest Feature

If the [Suggester](#) doesn't suit your needs, you can use the Terms component in Solr to build a similar feature for your own search application. Simply submit a query specifying whatever characters the user has typed so far as a prefix. For example, if the user has typed "at", the search engine's interface would submit the following query:

```
http://localhost:8983/solr/terms?terms.fl=name&terms.prefix=at
```

Result:

```
<?xml version="1.0" encoding="UTF-8"?>
<response>
<lst name="responseHeader">
<int name="status">0</int>
<int name="QTime">120</int>
</lst>
<lst name="terms">
<lst name="name">
<int name="ata">5</int> <int name="ati">5</int>
</lst>
</lst>
</response>
```

You can use the parameter `omitHeader=true` to omit the response header from the query response, like so:

```
http://localhost:8983/solr/terms?terms.fl=name&terms.prefix=at&indent=true&wt=json&o
```

Result:

```
{
  "terms": [
    "name", [
      "ata", 1,
      "ati", 1]]}
```

Distributed Search Support

The TermsComponent also supports distributed indexes. For the /terms request handler, you must provide the following two parameters:

Parameter	Description
shards	Specifies the shards in your distributed indexing configuration. For more information about distributed indexing, see Distributed Search with Index Sharding .
shards.qt	Specifies the request handler Solr uses for requests to shards.

The Term Vector Component

The Term Vector Component (TVC) is a search component designed to return information about documents. For each document, the TVC can return the term vector, the term frequency, inverse document frequency, position, and offset information. The TVC is stored when setting the `termVector` attribute on a field:

```
<field name="features"
      type="text"
      indexed="true"
      stored="true"
      multiValued="true"
      termVectors="true"
      termPositions="true"
      termOffsets="true" />
```

As with most components, there are a number of options that are outlined in the samples below. All examples are based on the Solr example.

Enabling the the TermVectorComponent

Changes for solrconfig.xml

To enable the Term VectorComponent, you need to configure a `searchComponent` element in your `solrconfig.xml` file, like so:

```
<searchComponent name="tvComponent"
                 class="org.apache.solr.handler.component.TermVectorComponent"/>
```

A request handler configuration using this component could look like this:

```
<requestHandler name="tvrh" class="org.apache.solr.handler.component.SearchHandler">
    <lst name="defaults">
        <bool name="tv">true</bool>
    </lst>
    <arr name="last-components">
        <str>tvComponent</str>
    </arr>
</requestHandler>
```

Invoking the Term Vector Component

The example below shows an invocation of this component:

`http://localhost:8983/solr/select/?q=*&version=2.2&start=0&rows=10&indent=on&qt=`

In the example, the component is associated with a request handler named `tvrh`, but you can associate it with any request handler. To turn on the component for a request, add the `tv=true` parameter (or add it to your RequestHandler defaults configuration).

Example output: See <http://wiki.apache.org/solr/TermVectorComponentExampleEnabled>

Optional Parameters

The example below shows optional parameters for this component:

`http://localhost:8983/solr/select/?q=*&version=2.2&start=0&rows=10&indent=on&qt=`

Boolean Parameters	Description
<code>tv.all</code>	A shortcut that invokes all the parameters listed below.
<code>tv.df</code>	Returns the Document Frequency (DF) of the term in the collection. This can be computationally expensive.
<code>tv.offsets</code>	Returns offset information for each term in the document.
<code>tv.positions</code>	Returns position information.
<code>tv.tf</code>	Returns document term frequency info per term in the document.
<code>tv.tf_idf</code>	Calculates TF*IDF for each term. Requires the parameters <code>tv.tf</code> and <code>tv.df</code> to be "true". This can be computationally expensive. (The results are not shown in example output)

To learn more about TermVector component output, see the Wiki page:

<http://wiki.apache.org/solr/TermVectorComponentExampleOptions>

For schema requirements, see the Wiki page: <http://wiki.apache.org/solr/FieldOptionsByUseCase>

The Term Vector component also accepts these optional parameters:

Parameters	Description
<code>tv.docIds</code>	Returns term vectors for the specified list of Lucene document IDs (not the Solr Unique Key).
<code>tv.fl</code>	Returns term vectors for the specified list of fields. If not specified, the <code>f1</code> parameter is used.

SolrJ and the Term Vector Component

Neither the SolrQuery class nor the QueryResponse class offer specific method calls to set Term Vector Component parameters or get the "termVectors" output. However, there is a patch for it: [SOLR-949](#).

The Stats Component

The Stats component returns simple statistics for numeric, string, and date fields within the document set.

Stats Component Parameters

The Stats Component accepts the following parameters:

Parameter	Description
stats	If true , then invokes the Stats component.
stats.field	Specifies a field for which statistics should be generated. This parameter may be invoked multiple times in a query in order to request statistics on multiple fields. (See the example below.)
stats.facet	Returns sub-results for values within the specified facet.

Statistics Returned

The table below describes the statistics returned by the Stats component.

Name	Description
min	The minimum value in the field.
max	The maximum value in the field.
sum	The sum of all values in the field.
count	The number of non-null values in the field.
missing	The number of null values in the field.
sumOfSquares	Sum of all values squared (useful for stddev).
mean	The average $(v_1 + v_2 + \dots + v_N) / N$
stddev	Standard deviation, measuring how widely spread the values in the data set are.

Example

The query below:

```
http://localhost:8983/solr/select?q=*&stats=true&stats.field=price&stats.field=pop
```

Would produce the following results:

```
<lst name="stats">
  <lst name="stats_fields">
    <lst name="price">
      <double name="min">0.0</double>
      <double name="max">2199.0</double>
      <double name="sum">5251.269999999995</double>
      <long name="count">15</long>
      <long name="missing">11</long>
      <double name="sumOfSquares">6038619.160300001</double>
      <double name="mean">350.0846666666664</double>
      <double name="stddev">547.737557906113</double>
    </lst>
    <lst name="popularity">
      <double name="min">0.0</double>
      <double name="max">10.0</double>
      <double name="sum">90.0</double>
      <long name="count">26</long>
      <long name="missing">0</long>
      <double name="sumOfSquares">628.0</double>
      <double name="mean">3.4615384615384617</double>
      <double name="stddev">3.5578731762756157</double>
    </lst>
  </lst>
</lst>
```

Here are the same results with facetting requested for the field `inStock`, using the parameter `&stats.facet=inStock`.

```
<lst name="{*}stats{*}">
<lst name="{*}stats{*}_fields">
<lst name="price">
<double name="min">0.0</double>
<double name="max">2199.0</double>
<double name="sum">5251.269999999995</double>
<long name="count">15</long>
<long name="missing">11</long>
<double name="sumOfSquares">6038619.160300001</double>
<double name="mean">350.0846666666664</double>
<double name="stddev">547.737557906113</double>
<lst name="facets">
<lst name="inStock">
<lst name="false">
<double name="min">11.5</double>
<double name="max">649.99</double>
<double name="sum">1161.39</double>
<long name="count">4</long>
<long name="missing">0</long>
<double name="sumOfSquares">653369.2551</double>
<double name="mean">290.3475</double>
<double name="stddev">324.63444676281654</double>
</lst>
<lst name="true">
<double name="min">0.0</double>
<double name="max">2199.0</double>
<double name="sum">4089.879999999999</double>
<long name="count ">11</long>
<long name="missing">0</long>
<double name="sumOfSquares">5385249.905200001</double>
<double name="mean">371.8072727272727</double>
<double name="stddev">621.6592938755265</double>
</lst>
</lst>
</lst>
</lst>
</lst>
```

The Stats Component and Faceting

The facet field can be selectively applied. That is if you want stats on field "A" and "B", you can facet a on "X" and B on "Y" using the parameters:

```
&stats.field=A&f.A.stats.facet=X&stats.field=B&f.B.stats.facet=Y
```

 All facet results are returned, so be careful what fields you ask for.

Multi-valued fields and facets may be slow.

Multi-value fields rely on `UnInvertedField.java` for implementation. This is like a `FieldCache`, so be aware of your memory footprint.

The Query Elevation Component

The [Query Elevation Component](#) lets you configure the top results for a given query regardless of the normal Lucene scoring. This is sometimes called "sponsored search," "editorial boosting," or "best bets." This component matches the user query text to a configured map of top results. The text can be any string or non-string IDs, as long as it's indexed. Although this component will work with any QueryParser, it makes the most sense to use with [DisMax](#) or [eDisMax](#).

The [Query Elevation Component](#) is supported by distributed searching.

Configuring the Query Elevation Component

You can configure the Query Elevation Component in the `solrconfig.xml` file. The default configuration looks like this:

```
<searchComponent name="elevator" class="solr.QueryElevationComponent" >
  <!-- pick a fieldType to analyze queries -->
  <str name="queryFieldType">string</str>
  <str name="config-file">elevate.xml</str>
</searchComponent>

<requestHandler name="/elevate" class="solr.SearchHandler" startup="lazy">
  <lst name="defaults">
    <str name="echoParams">explicit</str>
  </lst>
  <arr name="last-components">
    <str>elevator</str>
  </arr>
</requestHandler>
```

Optionally, in the Query Elevation Component configuration you can also specify the following to distinguish editorial results from "normal" results:

```
<str name="editorialMarkerFieldName">foo</str>
```

The Query Elevation Search Component takes the following arguments:

Argument	Description
queryFieldType	Specifies which fieldType should be used to analyze the incoming text. For example, it may be appropriate to use a fieldType with a LowerCaseFilter.

config-file	Path to the file that defines query elevation. This file must exist in <code>\$<instanceDir>/conf/<config-file></code> or <code>\$<dataDir/<config-file></code> . If the file exists in the /conf/ directory it will be loaded once at startup. If it exists in the data directory, it will be reloaded for each IndexReader.
forceElevation	By default, this component respects the requested <code>sort</code> parameter: if the request asks to sort by date, it will order the results by date. If <code>forceElevation=true</code> (the default), results will first return the boosted docs, then order by date.

elevate.xml

Elevated query results are configured in an external XML file specified in the config-file argument. An `elevate.xml` file might look like this:

```
<elevate>
  <query text="AAA">
    <doc id="A" />
    <doc id="B" />
  </query>

  <query text="ipod">
    <doc id="A" />

    <!-- you can optionally exclude documents from a query result -->

    <doc id="B" exclude="true" />
  </query>
</elevate>
```

In this example, the query "AAA" would first return documents A and B, then whatever normally appears for the same query. For the query "ipod", it would first return A, and would make sure that B is not in the result set.

Using the Query Elevation Component

The enableElevation Parameter

For debugging it may be useful to see results with and without the elevated docs. To hide results, use `enableElevation=false`:

```
http://localhost:8983/solr/elevate?q=YYYY&debugQuery=true&enableElevation=true
```

```
http://localhost:8983/solr/elevate?q=YYYY&debugQuery=true&enableElevation=false
```

The forceElevation Parameter

You can force elevation during runtime by adding `forceElevation=true` to the query URL:

```
http://localhost:8983/solr/elevate?q=YYYY&debugQuery=true&enableElevation=true&force
```

The exclusive Parameter

You can force Solr to return only the results specified in the elevation file by adding `exclusive=true` to the URL:

```
http://localhost:8983/solr/elevate?q=YYYY&debugQuery=true&exclusive=true
```

The fq Parameter

Query elevation respects the standard filter query (`fq`) parameter. That is, if the query contains the `fq` parameter, all results will be within that filter even if `elevate.xml` adds other documents to the result set.

Response Writers

A Response Writer generates the formatted response of a search. Solr supports a variety of Response Writers to ensure that query responses can be parsed by the appropriate language or application.

The `wt` parameter selects the Response Writer to be used. The table below lists the most common settings for the `wt` parameter.

wt Parameter Setting	Response Writer Selected
csv	CSVResponseWriter
json	JSONResponseWriter
php	PHPResponseWriter
phps	PHPSerializedResponseWriter
python	PythonResponseWriter
ruby	RubyResponseWriter
xml	XMLResponseWriter
xslt	XSLTResponseWriter

The Standard XML Response Writer

The XML Response Writer is the most general purpose and reusable Response Writer currently included with Solr. It is the format used in most discussions and documentation about the response of Solr queries.

Note that the XSLT Response Writer can be used to convert the XML produced by this writer to other vocabularies or text-based formats.

The behavior of the XML Response Writer can be driven by the following query parameters.

The version Parameter

The `version` parameter determines the XML protocol used in the response. Clients are strongly encouraged to *always* specify the protocol version, so as to ensure that the format of the response they receive does not change unexpectedly when the Solr server is upgraded.

XML Version	Notes	Comments

2.0	An <arr> tag was used for multiValued fields only if there was more than one value.	Not supported in Solr 4.
2.1	An <arr> tag is used for multiValued fields even if there is only one value.	Not supported in Solr 4.
2.2	The format of the responseHeader changed to use the same <lst> structure as the rest of the response.	Supported in Solr 4.

The default value is the latest supported.

The **stylesheet** Parameter

The `stylesheet` parameter can be used to direct Solr to include a `<?xml-stylesheet type="text/xsl" href="..."?>` declaration in the XML response it returns.

The default behavior is not to return any stylesheet declaration at all.



Use of the `stylesheet` parameter is discouraged, as there is currently no way to specify external stylesheets, and no stylesheets are provided in the Solr distributions. This is a legacy parameter, which may be developed further in a future release.

The **indent** Parameter

If the `indent` parameter is used, and has a non-blank value, then Solr will make some attempts at indenting its XML response to make it more readable by humans.

The default behavior is not to indent.

The **XSLT Response Writer**

The XSLT Response Writer applies an XML stylesheet to output. It can be used for tasks such as formatting results for an RSS feed.

tr Parameter

The XSLT Response Writer accepts one parameter: the `tr` parameter, which identifies the XML transformation to use. The transformation must be found in the Solr `conf/xslt` directory.

The Content-Type of the response is set according to the `<xsl:output>` statement in the XSLT transform, for example: `<xsl:output media-type="text/html"/>`

Configuration

The example below, from the default `solrconfig.xml` file, shows how the XSLT Response Writer is configured.

```
<!--
Changes to XSLT transforms are taken into account
every xsltCacheLifetimeSeconds at most.

-->
<queryResponseWriter
  name="xslt"
  class="org.apache.solr.request.XSLTResponseWriter"
>
  <int name="xsltCacheLifetimeSeconds">5</int>
</queryResponseWriter>
```

A value of 5 for `xsltCacheLifetimeSeconds` is good for development, to see XSLT changes quickly. For production you probably want a much higher value.

JSON Response Writer

A very commonly used Response Writer is the `JsonResponseWriter`, which formats output in JavaScript Object Notation (JSON), a lightweight data interchange format specified in RFC 4627. Setting the `wt` parameter to `json` invokes this Response Writer.

With Solr 4, the `JsonResponseWriter` has been changed:

- The default mime type for the writer is now `application/json`.
- The example `solrconfig.xml` has been updated to explicitly use this parameter to set the type to `text/plain`:

```
<queryResponseWriter name="json" class="solr.JSONResponseWriter">
  <!-- For the purposes of the tutorial, JSON response are written as
      plain text so that it's easy to read in *any* browser.
      If you are building applications that consume JSON, just remove
      this override to get the default "application/json" mime type.
  -->
  <str name="content-type">text/plain</str>
</queryResponseWriter>
```

Python Response Writer

Solr has an optional Python response format that extends its JSON output in the following ways to allow the response to be safely evaluated by the python interpreter:

- true and false changed to True and False

- Python unicode strings are used where needed
- ASCII output (with unicode escapes) is used for less error-prone interoperability
- newlines are escaped
- null changed to None

PHP Response Writer and PHP Serialized Response Writer

Solr has a PHP response format that outputs an array (as PHP code) which can be evaluated. Setting the `wt` parameter to `php` invokes the PHP Response Writer.

Example usage:

```
$code = file_get_contents('http://localhost:8983/solr/select?q=iPod&wt=*php*');
eval("$result = " . $code . ";");
print_r($result);
```

Solr also includes a PHP Serialized Response Writer that formats output in a serialized array. Setting the `wt` parameter to `phps` invokes the PHP Serialized Response Writer.

Example usage:

```
$serializedResult =
file_get_contents('http://localhost:8983/solr/select?q=iPod&wt=*php{*}s');
$result = unserialize($serializedResult);
print_r($result);
```

Before you use either the PHP or Serialized PHP Response Writer, you may first need to un-comment these two lines in `solrconfig.xml`:

```
<queryResponseWriter name="php" class="org.apache.solr.request.PHPResponseWriter"/>
<queryResponseWriter name="phps"
class="org.apache.solr.request.PHPSerializedResponseWriter"/>
```

Ruby Response Writer

Solr has an optional Ruby response format that extends its JSON output in the following ways to allow the response to be safely evaluated by Ruby's interpreter:

- Ruby's single quoted strings are used to prevent possible string exploits.
- \ and ' are the only two characters escaped.
- Unicode escapes are not used. Data is written as raw UTF-8.
- nil used for null.
- => is used as the key/value separator in maps.

Here is a simple example of how one may query Solr using the Ruby response format:

```
require 'net/http'
h = Net::HTTP.new('localhost', 8983)
hresp, data = h.get('/solr/select?q=iPod&wt=ruby', nil)
rsp = eval(data)
puts 'number of matches = ' + rsp['response'][ 'numFound' ].to_s
#print out the name field for each returned document
rsp['response'][ 'docs' ].each { |doc| puts 'name field = ' + doc[ 'name' ] }
```

CSV Response Writer

The CSV response writer returns a list of documents in comma-separated values (CSV) format. Other information that would normally be included in a response, such as facet information, is excluded.

The CSV response writer supports multi-valued fields, and the output of this CSV format is compatible with Solr's [CSV update format](#).

CSV Parameters

These parameters specify the CSV format that will be returned. You can accept the default values or specify your own.

Parameter	Default Value
csv.encapsulator	"
csv.escape	None
csv.separator	,
csv.header	Defaults to true. If false, Solr does not print the column headers
csv.newline	\n
csv.null	Defaults to a zero length string. Use this parameter when a document has no value for a particular field.

Multi-Valued Field CSV Parameters

These parameters specify how multi-valued fields are encoded. Per-field overrides for these values can be done using `f.<fieldname>.csv.separator=|`.

Parameter	Default Value
-----------	---------------

csv.mv.encapsulator	None
csv.mv.escape	\
csv.mv.separator	Defaults to the <code>csv.separator</code> value

Example

`http://localhost:8983/solr/select?q=ipod&fl=id,cat,name,popularity,price,score&wt=cs`
returns:

```
id,cat,name,popularity,price,score
IW-02,"electronics,connector",iPod & iPod Mini USB 2.0 Cable,1,11.5,0.98867977
F8V7067-APL-KIT,"electronics,connector",Belkin Mobile Power Cord for iPod w/
Dock,1,19.95,0.6523595
MA147LL/A,"electronics,music",Apple 60 GB iPod with Video Playback
Black,10,399.0,0.2446348
```

Binary Response Writer

Solr also includes a Response Writer that outputs binary format for use with a Java client. See [Client APIs](#) for more details.

Near Real Time Searching

Near Real Time (NRT) search means that documents are available for search almost immediately after being indexed: additions and updates to documents are seen in 'near' real time. Solr 4 no longer blocks updates while a commit is in progress. Nor does it wait for background merges to complete before opening a new search of indexes and returning.

With NRT, you can modify a `commit` command to be a **soft commit**, which avoids parts of a standard commit that can be costly. You will still want to do standard commits to ensure that documents are in stable storage, but **soft commits** let you see a very near real time view of the index in the meantime. However, pay special attention to cache and autowarm settings as they can have a significant impact on NRT performance.

Commits and Optimizing

A commit operation makes index changes visible to new search requests. A **hard commit** uses the transaction log to get the id of the latest document changes, and also calls `fsync` on the index files to ensure they have been flushed to stable storage and no data loss will result from a power failure.

A **soft commit** is much faster since it only makes index changes visible and does not `fsync` index files or write a new index descriptor. If the JVM crashes or there is a loss of power, changes that occurred after the last **hard commit** will be lost. Search collections that have NRT requirements (that want index changes to be quickly visible to searches) will want to soft commit often but hard commit less frequently. A softCommit may be "less expensive" in terms of time, but not free, since it can slow throughput.

An **optimize** is like a **hard commit** except that it forces all of the index segments to be merged into a single segment first. Depending on the use, this operation should be performed infrequently (e.g., nightly), if at all, since it involves reading and re-writing the entire index. Segments are normally merged over time anyway (as determined by the merge policy), and optimize just forces these merges to occur immediately.

Soft commit takes uses two parameters: `maxDocs` and `maxTime`.

Parameter	Description
<code>maxDocs</code>	Integer. Defines the number of documents to queue before pushing them to the index. It works in conjunction with the <code>update_handler_autosoftcommit_max_time</code> parameter in that if either limit is reached, the documents will be pushed to the index.

maxTime	The number of milliseconds to wait before pushing documents to the index. It works in conjunction with the <code>update_handler_autosoftcommit_max_docs</code> parameter in that if either limit is reached, the documents will be pushed to the index.
---------	---

Use `maxDocs` and `maxTime` judiciously to fine-tune your commit strategies.

AutoCommits

An autocommit also uses the parameters `maxDocs` and `maxTime`. However it's useful in many strategies to use both a hard `autocommit` and `autosoftcommit` to achieve more flexible commits.

A common configuration is to do a hard `autocommit` every 1-10 minutes and a `autosoftcommit` every second. With this configuration, new documents will show up within about a second of being added, and if the power goes out, soft commits are lost unless a hard commit has been done.

For example:

```
<autoSoftCommit>
  <maxTime>1000</maxTime>
</autoSoftCommit>
```

It's better to use `maxTime` rather than `maxDocs` to modify an `autoSoftCommit`, especially when indexing a large number of documents through the commit operation. It's also better to turn off `autoSoftCommit` for bulk indexing.

Optional Attributes for commit and optimize

Parameter	Description	
<code>waitSearcher</code>	true, false	Block until a new searcher is opened and registered as the main query searcher, making the changes visible. Default is true.
<code>softCommit</code>	true, false	Perform a soft commit. This will refresh the view of the index faster, but without guarantees that the document is stably stored. Default is false.
<code>expungeDeletes</code>	true, false	Valid for <code>commit</code> only. This parameter purges deleted data from segments. The default is false.
<code>maxSegments = N</code>	integer	Valid for <code>optimize</code> only. Optimize down to at most this number of segments. The default is 1.

Example of `commit` and `optimize` with optional attributes:

```
<commit waitSearcher="false"/>
<commit waitSearcher="false" expungeDeletes="true" />
<optimize waitSearcher="false"/>
```

Passing commit and commitWithin parameters as part of the URL

Update handlers can also get commit-related parameters as part of the update URL. This example adds a small test document and causes an explicit commit to happen immediately afterwards:

```
http://localhost:8983/solr/update?stream.body=<add><doc>
  <field name="id">testdoc</field></doc></add>&commit=true
```

Alternately, you may want to use this:

```
http://localhost:8983/solr/update?stream.body=<optimize/>
```

This example causes the index to be optimized down to at most 10 segments, but won't wait around until it's done (waitFlush=false):

```
curl 'http://localhost:8983/solr/update?optimize=true&maxSegments=10&waitFlush=false'
```

This example adds a small test document with a commitWithin instruction that tells Solr to make sure the document is committed no later than 10 seconds later (this method is generally preferred over explicit commits):

```
curl http://localhost:8983/solr/update?commitWithin=10000
-H "Content-Type: text/xml" --data-binary
'<add><doc><field name="id">testdoc</field></doc></add>'
```

RealTime Get

For index updates to be visible (searchable), some kind of commit must reopen a searcher to a new point-in-time view of the index. The **realtime get** feature allows retrieval (by unique-key) of the latest version of any documents without the associated cost of reopening a searcher. This is primarily useful when using Solr as a NoSQL data store and not just a search index.

Realtime Get currently relies on the update log feature, which is enabled by default. It relies on an update log, which is configured in `solrconfig.xml`, in a section like:

```
<updateLog>
  <str name="dir">${solr.ulog.dir:}</str>
</updateLog>
```

The latest `example solrconfig.xml` should also have a request handler named `/get` already defined.

Start (or restart) the Solr server, and then index a document:

```
curl 'http://localhost:8983/solr/update/json?commitWithin=10000000'
-H 'Content-type:application/json' -d '[{"id":"mydoc", "title":"realtime-get test!"}]'
```

If you do a normal search, this document should not be found:

```
http://localhost:8983/solr/select?q=id:mydoc
...
"response":
{ "numFound":0, "start":0, "docs":[] }
```

However if you use the realtime get handler exposed at `/get`, you should be able to retrieve that document:

```
http://localhost:8983/solr/get?id=mydoc
...
{ "doc":{ "id":"mydoc", "title":"realtime-get test!" } }
```

You can also specify multiple documents at once via the **ids** parameter and a comma separated list of ids, or by using multiple **id** parameters. If you specify multiple ids, or use the **ids** parameter, the response will mimic a normal query response to make it easier for existing clients to parse. Since you've only indexed one document, the following equivalent examples just repeat the same id.

```
http://localhost:8983/solr/get?ids=mydoc,mydoc
http://localhost:8983/solr/get?id=mydoc&id=mydoc
...
{ "response":
  { "numFound":2, "start":0, "docs":
    [ { "id":"mydoc", "title":["realtime-get test!"]},
      { "id":"mydoc", "title":["realtime-get test!"]} ]
  }
}
```

The Well-Configured Solr Instance

This section tells you how to fine-tune your Solr instance for optimum performance. This section covers the following topics:

[Configuring solrconfig.xml](#): Describes how to work with the main configuration file for Solr.

[Core Admin and Configuring solr.xml](#): Describes how to configure your Solr core, or multiple Solr cores within a single instance.

[Lucene IndexWriters](#): Describes how to configure the index writers in the underlying Lucene engine.

[HTTP Request Dispatcher](#): Describes how to configure Solr's response to HTTP requests

[JVM Settings](#): Gives some guidance on best practices for working with Java Virtual Machines.



The focus of this section is on configuring a single Solr instance. To scale a Solr implementation in a cluster environment, see [SolrCloud](#). There are also options to scale through sharding or replication, described in the section [Scaling and Distribution](#).

For more information about factors affecting Solr performance, see
<http://wiki.apache.org/solr/SolrPerformanceFactors>.

Configuring solrconfig.xml

The `solrconfig.xml` file is the configuration file with the most parameters affecting Solr itself. While configuring Solr, you'll work with `solrconfig.xml` often. The file comprises a series of XML statements that set configuration values. In `solrconfig.xml`, you configure important features such as:

- request handlers
- listeners (processes that "listen" for particular query-related events; listeners can be used to trigger the execution of special code, such as invoking some common queries to warm-up caches)
- the Request Dispatcher for managing HTTP communications
- the Admin Web interface
- parameters related to replication and duplication (these parameters are covered in detail in [Legacy Scaling and Distribution](#))

The `solrconfig.xml` file is found in the `solr/conf/` directory. The example file is well-commented, and includes information on best practices for most installations.

We've covered the options in the following sections:

More Information

- The Solr Wiki has a comprehensive page on `solrconfig.xml`, at <http://wiki.apache.org/solr/SolrConfigXml>.
- [6 Sins of solrconfig.xml modifications](#) from `solr.pl`.

DataDir and DirectoryFactory in SolrConfig

Specifying a Location for Index Data with the dataDir Parameter

By default, Solr stores its index data in a directory called /data under the Solr home. If you would like to specify a different directory for storing index data, use the <dataDir> parameter in the solrconfig.xml file. You can specify another directory either with a full pathname or a pathname relative to the current working directory of the servlet container. For example:

```
<dataDir>/var/data/solr/</dataDir>
```

If you are using replication to replicate the Solr index (as described in [Legacy Scaling and Distribution](#)), then the <dataDir> directory should correspond to the index directory used in the replication configuration.

Specifying the DirectoryFactory For Your Index

The default solr.StandardDirectoryFactory is filesystem based, and tries to pick the best implementation for the current JVM and platform. You can force a particular implementation by specifying solr.MMapDirectoryFactory, solr.NIOFSDirectoryFactory, or solr.SimpleFSDirectoryFactory.

```
<directoryFactory name="DirectoryFactory"
                  class="${solr.directoryFactory:solr.StandardDirectoryFactory}"/>
```

The solr.RAMDirectoryFactory is memory based, not persistent, and does not work with replication. Use this DirectoryFactory to store your index in RAM.

```
<directoryFactory class="org.apache.solr.core.RAMDirectoryFactory" />
```

Lib Directives in SolrConfig

Solr allows loading plugins by defining `<lib/>` directives in `solrconfig.xml`.

The plugins are loaded in the order they appear in `solrconfig.xml`. If there are dependencies, list the lowest level dependency jar first.

Regular expressions can be used to provide control loading jars with dependencies on other jars in the same directory. All directories are resolved as relative to the Solr `instanceDir`.

```
<lib dir="../../../../contrib/extraction/lib" regex=".*\.\.jar" />
<lib dir="../../../../dist/" regex="apache-solr-cell-\d.*\.\.jar" />

<lib dir="../../../../contrib/clustering/lib/" regex=".*\.\.jar" />
<lib dir="../../../../dist/" regex="apache-solr-clustering-\d.*\.\.jar" />

<lib dir="../../../../contrib/langid/lib/" regex=".*\.\.jar" />
<lib dir="../../../../dist/" regex="apache-solr-langid-\d.*\.\.jar" />

<lib dir="../../../../contrib/velocity/lib" regex=".*\.\.jar" />
<lib dir="../../../../dist/" regex="apache-solr-velocity-\d.*\.\.jar" />
```

IndexConfig in SolrConfig

The `<indexConfig>` section of `solrconfig.xml` defines low-level behavior of the Lucene index writers. By default, the settings are commented out in the sample `solrconfig.xml` included with Solr, which means the defaults are used. In most cases, the defaults are fine.

```
<indexConfig>
  ...
</indexConfig>
```

 Prior to Solr 4, many of these settings were contained in sections called `mainIndex` and `indexDefaults`. In Solr 4, those sections are deprecated and removed. Any settings that used to be in those sections, now belong in `<indexConfig>`.

Parameters covered in this section:

- [Sizing Index Segments](#)
- [Merging Index Segments](#)
- [Index Locks](#)
- [Other Indexing Settings](#)

Sizing Index Segments

ramBufferSizeMB

Once accumulated document updates exceed this much memory space (defined in megabytes), then the pending updates are flushed. This can also create new segments or trigger a merge. Using this setting is generally preferable to `maxBufferedDocs`. If both `maxBufferedDocs` and `ramBufferSizeMB` are set in `solrconfig.xml`, then a flush will occur when either limit is reached. The default is 100Mb (raised from 32Mb for Solr 4.1).

```
<ramBufferSizeMB>100</ramBufferSizeMB>
```

maxBufferedDocs

Sets the number of document updates to buffer in memory before flushed to disk and added to the current index segment. If the segment fills up, a new one may be created, or a merge may be started. The default Solr configuration leaves this value undefined.

```
<maxBufferedDocs>1000</maxBufferedDocs>
```

maxIndexingThreads

The maximum number of simultaneous threads used to index documents. Once this threshold is reached, additional threads will wait for the others to finish. The default is 8. This parameter is new for Solr 4.1.

```
<maxIndexingThreads>8</maxIndexingThreads>
```

UseCompoundFile

Setting `<useCompoundFile>` to **true** combines the various files of a segment into a single file, although the default is **false**. On systems where the number of open files allowed per process is limited, setting this to **false** may avoid hitting that limit (the open files limit might also be tunable for your OS with the Linux/Unix `ulimit` command, or something similar for other operating systems). In some cases, other internal factors may set a segment to "compound=false", even if this setting is explicitly set to true, so the compounding of the files in a segment may not always happen.

Updating a compound index may incur a minor performance hit for various reasons, depending on the runtime environment. For example, filesystem buffers are typically associated with open file descriptors, which may limit the total cache space available to each index.

This setting may also affect how much data needs to be transferred during index replication operations.

The default is **false**.

```
<useCompoundFile>false</useCompoundFile>
```

Merging Index Segments

mergeFactor

The `mergeFactor` controls how many segments a Lucene index is allowed to have before it is coalesced into one segment. When an update is made to an index, it is added to the most recently opened segment. When that segment fills up (see `maxBufferedDocs` and `ramBufferSizeMB` in the next section), a new segment is created and subsequent updates are placed there.

If creating a new segment would cause the number of lowest-level segments to exceed the `mergeFactor` value, then all those segments are merged together to form a single large segment. Thus, if the merge factor is ten, each merge results in the creation of a single segment that is roughly ten times larger than each of its ten constituents. When there are `mergeFactor` settings for these larger segments, then they in turn are merged into an even larger single segment. This process can continue indefinitely.

Choosing the best merge factor is generally a trade-off of indexing speed vs. searching speed. Having fewer segments in the index generally accelerates searches, because there are fewer places to look. It also can also result in fewer physical files on disk. But to keep the number of segments low, merges will occur more often, which can add load to the system and slow down updates to the index.

Conversely, keeping more segments can accelerate indexing, because merges happen less often, making an update is less likely to trigger a merge. But searches become more computationally expensive and will likely be slower, because search terms must be looked up in more index segments. Faster index updates also means shorter commit turnaround times, which means more timely search results.

The default value in the example `solrconfig.xml` is 10, which is a reasonable starting point.

```
<mergeFactor>10</mergeFactor>
```

mergePolicy

Defines how merging segments is done. The default in Solr is the `TieredMergePolicy`. This default policy merges segments of approximately equal size, subject to an allowed number of segments per tier. Other policies available are the `LogMergePolicy`, `LogByteSizeMergePolicy` and `LogDocMergePolicy`. For more information on these policies, please see the Lucene javadocs at http://lucene.apache.org/core/4_0_0/core/index.html?org/apache/lucene/index/MergePolicy.html.

```
<mergePolicy class="org.apache.lucene.index.TieredMergePolicy">
  <int name="maxMergeAtOnce">10</int>
  <int name="segmentsPerTier">10</int>
</mergePolicy>
```

 When using `TieredMergePolicy`, the setting `maxMergeDocs` is not needed. Since this is the default in Solr, the setting is effectively removed. However, if using another policy, this setting may be useful.

mergeScheduler

The merge scheduler controls how merges are performed. The default `ConcurrentMergeScheduler` performs merges in the background using separate threads. The alternative, `SerialMergeScheduler`, does not perform merges with separate threads.

```
<mergeScheduler class="org.apache.lucene.index.ConcurrentMergeScheduler"/>
```

Index Locks

lockType

The LockFactory options specify its implementation.

`lockType=single` uses `SingleInstanceLockFactory`, and is for a read-only index or when there is no possibility of another process trying to modify the index.

`lockType=native` uses `NativeFSLockFactory` to specify native OS file locking. Do not use when multiple Solr web applications in the same JVM are attempting to share a single index.

`lockType=simple` uses `SimpleFSLockFactory` to specify a plain file for locking.

`native` is the default for Solr3.6 and later versions; otherwise `simple` is the default.

For more information on the nuances of each LockFactory, see
<http://wiki.apache.org/lucene-java/AvailableLockFactories>.

```
<lockType>native</lockType>
```

unlockOnStartup

If `true`, any write or commit locks that have been held will be unlocked on system startup. This defeats the locking mechanism that allows multiple processes to safely access a Lucene index. The default is `false`, and changing this should only be done with care. This parameter is not used if the `lockType` is "none" or "single".

```
<unlockOnStartup>false</unlockOnStartup>
```

writeLockTimeout

The maximum time to wait for a write lock on an IndexWriter. The default is 1000, expressed in milliseconds.

```
<writeLockTimeout>1000</writeLockTimeout>
```

Other Indexing Settings

There are a few other parameters that may be important to configure for your implementation. These settings affect how or when updates are made to an index.

Setting	Description
---------	-------------

termIndexInterval	Controls how often terms are loaded into memory. The default is 128.
reopenReaders	Controls if IndexReaders will be re-opened, instead of closed and then opened, which is often less efficient. The default is true.
deletionPolicy	Controls how commits are retained in case of rollback. The default is SolrDeletionPolicy, which has sub-parameters for the maximum number of commits to keep (<code>maxCommitsToKeep</code>), the maximum number of optimized commits to keep (<code>maxOptimizedCommitsToKeep</code>), and the maximum age of any commit to keep (<code>maxCommitAge</code>), which supports DateMathParser syntax.
infoStream	The InfoStream allows sending detailed debug information from the indexing process to the specified file.

```
<termIndexInterval>128</termIndexInterval>
<reopenReaders>true</reopenReaders>
<deletionPolicy class="solr.SolrDeletionPolicy">
  <str name="maxCommitsToKeep">1</str>
  <str name="maxOptimizedCommitsToKeep">0</str>
  <str name="maxCommitAge">1DAY</str>
</deletionPolicy>
<infoStream file="INFOSTREAM.txt">false</infoStream>
```

-  The `maxFieldLength` parameter was removed in Solr 4. If restricting the length of fields is important to you, you can get similar behavior with the `LimitTokenCountFactory`, which can be defined for the fields you'd like to limit. For example, `<filter class="solr.LimitTokenCountFilterFactory" maxTokenCount="10000" />` would limit the field to 10,000 characters.

UpdateHandlers in SolrConfig

The settings in this section are configured in the `<updateHandler>` element in `solrconfig.xml` and may affect the performance of index updates. These settings affect how updates are done internally. `<updateHandler>` configurations do not affect the higher level configuration of [RequestHandlers](#) that process client update requests.

```
<updateHandler class="solr.DirectUpdateHandler2">
  ...
</updateHandler>
```

Topics covered in this section:

- [Commits](#)
- [Event](#)
- [Listeners](#)
- [Transaction Log](#)

Commits

Data sent to Solr is not searchable until it has been *committed* to the index. The reason for this is that in some cases commits can be slow and they should be done in isolation from other possible commit requests to avoid overwriting data. So, it's preferable to provide control over when data is committed. Several options are available to control the timing of commits.

commit and softCommit

With Solr 4, `commit` is generally used only as a boolean flag sent with a client update request. The command `commit=true` would perform a commit as soon as the data as finished loading to Solr.

You can also set the flag `softCommit=true` to do a 'soft' commit, meaning that Solr will commit your changes quickly but not guarantee that documents are in stable storage. This is an implementation of Near Real Time storage, a feature that boosts document visibility, since you don't have to wait for background merges and storage (to ZooKeeper, if using [SolrCloud](#)) to finish before moving on to something else. A full commit means that, if a server crashes, Solr will know exactly where your data was stored; a soft commit means that the data is stored, but the location information isn't yet stored. The tradeoff is that a soft commit gives you faster visibility because it's not waiting for background merges to finish.

For more information about Near Real Time operations, see [Near Real Time Searching](#).

autoCommit

These settings control how often pending updates will be automatically pushed to the index. An alternative to `autoCommit` is to use `commitWithin`, which can be defined when making the update request to Solr (i.e., when pushing documents), or in an update RequestHandler.

Setting	Description
maxDocs	The number of updates that have occurred since the last commit.
maxTime	The number of milliseconds since the oldest uncommitted update.
openSearcher	Whether to open a new searcher when performing a commit. If this is false , the default, the commit will flush recent index changes to stable storage, but does not cause a new searcher to be opened to make those changes visible

If either of these `maxDocs` or `maxTime` limits are reached, Solr automatically performs a commit operation. If the `autoCommit` tag is missing, then only explicit commits will update the index. The decision whether to use auto-commit or not depends on the needs of your application.

Determining the best auto-commit settings is a tradeoff between performance and accuracy. Settings that cause frequent updates will improve the accuracy of searches because new content will be searchable more quickly, but performance may suffer because of the frequent updates. Less frequent updates may improve performance but it will take longer for updates to show up in queries.

```
<autoCommit>
  <maxDocs>10000</maxDocs>
  <maxTime>1000</maxTime>
  <openSearcher>false</openSearcher>
</autoCommit>
```

You can also specify 'soft' autoCommits in the same way that you can specify 'soft' commits, except that instead of using `autoCommit` you set the `autoSoftCommit` tag.

```
<autoSoftCommit>
  <maxTime>1000</maxTime>
</autoSoftCommit>
```

maxPendingDeletes

This value sets a limit on the number of deletions that Solr will buffer during document deletion. This can affect how much memory is used during indexing.

```
<maxPendingDeletes>100000</maxPendingDeletes>
```

Event Listeners

The UpdateHandler section is also where update-related event listeners can be configured. These can be triggered to occur after a commit or optimize event, or after only an optimize event.

The listener is called with the `RunExecutableListener`, which runs an external executable with the defined set of instructions. The available commands are:

Setting	Description
event	If postCommit , the <code>RunExecutableListener</code> will be run after every commit or optimize. If postOptimize , the <code>RunExecutableListener</code> will be run every optimize only.
exe	The name of the executable to run. It should include the path to the file, relative to Solr home.
dir	The directory to use as the working directory. The default is ".".
wait	Forces the calling thread to wait until the executable returns a response. The default is true .
args	Any arguments to pass to the program. The default is none.
env	Any environment variables to set. The default is none.

Below is the example from `solrconfig.xml`, which shows an example from script-based replication described at <http://wiki.apache.org/solr/CollectionDistribution>:

```
<listener event="postCommit" class="solr.RunExecutableListener">
  <str name="exe">solr/bin/snapshooter</str>
  <str name="dir">.</str>
  <bool name="wait">true</bool>
  <arr name="args"> <str>arg1</str> <str>arg2</str> </arr>
  <arr name="env"> <str>MYVAR=val1</str> </arr>
</listener>
```

Transaction Log

As described in the section [RealTime Get](#), a transaction log is required for that feature. It is configured in the `updateHandler` section of `solrconfig.xml`.

Realtime Get currently relies on the update log feature, which is enabled by default. It relies on an update log, which is configured in `solrconfig.xml`, in a section like:

```
<updateLog>
  <str name="dir">${solr.ulog.dir:}</str>
</updateLog>
```

Query Settings in SolrConfig

The settings in this section affect the way that Solr will process and respond to queries. These settings are all configured in child elements of the `<query>` element in `solrconfig.xml`.

```
<query>
  ...
</query>
```

Topics covered in this section:

- [Caches](#)
- [Query Sizing and Warming](#)
- [Query-Related Listeners](#)

Caches

Solr caches are associated with a specific instance of an Index Searcher, a specific view of an index that doesn't change during the lifetime of that searcher. As long as that Index Searcher is being used, any items in its cache will be valid and available for reuse. Caching in Solr differs from caching in many other applications in that cached Solr objects do not expire after a time interval; instead, they remain valid for the lifetime of the Index Searcher.

When a new searcher is opened, the current searcher continues servicing requests while the new one auto-warms its cache. The new searcher uses the current searcher's cache to pre-populate its own. When the new searcher is ready, it is registered as the current searcher and begins handling all new search requests. The old searcher will be closed once it has finished servicing all its requests.

In Solr, there are three cache implementations: `solr.search.LRUCache`, `solr.search.FastLRUCache`, and `solr.search.LFUCache`.

The acronym LRU stands for Least Recently Used. When an LRU cache fills up, the entry with the oldest last-accessed timestamp is evicted to make room for the new entry. The net effect is that entries that are accessed frequently tend to stay in the cache, while those that are not accessed frequently tend to drop out and will be re-fetched from the index if needed again.

The `FastLRUCache`, which was introduced in Solr 1.4, is designed to be lock-free, so it is well suited for caches which are hit several times in a request.

Both `LRUCache` and `FastLRUCache` use an auto-warm count that supports both integers and percentages which get evaluated relative to the current size of the cache when warming happens.

The `LFUCache` refers to the Least Frequently Used cache. This works in a way similar to the LRU cache, except that when the cache fills up, the entry that has been used the least is evicted.

The Statistics page in the Solr Admin UI will display information about the performance of all the active caches. This information can help you fine-tune the sizes of the various caches appropriately for your particular application. When a Searcher terminates, a summary of its cache usage is also written to the log.

Each cache has settings to define it's initial size (`initialSize`), maximum size (`size`) and number of items to use for during warming (`autowarmCount`). The LRU and FastLRU cache implementations can take a percentage instead of an absolute value for `autowarmCount`.

Details of each cache are described below.

filterCache

This cache is used by `SolrIndexSearcher` for filters (DocSets) for unordered sets of all documents that match a query. The numeric attributes control the number of entries in the cache.

Solr uses the `filterCache` to cache results of queries that use the `fq` search parameter.

Subsequent queries using the same parameter setting result in cache hits and rapid returns of results. See [Searching](#) for a detailed discussion of the `fq` parameter.

Solr also makes this cache for faceting when the configuration parameter `facet.method` is set to `fc`. For a discussion of faceting, see [Searching](#).

```
<filterCache class="solr.LRUCache"
    size="512"
    initialSize="512"
    autowarmCount="128"/>
```

queryResultCache

This cache holds the results of previous searches: ordered lists of document IDs (DocList) based on a query, a sort, and the range of documents requested.

```
<queryResultCache class="solr.LRUCache"
    size="512"
    initialSize="512"
    autowarmCount="128"/>
```

documentCache

This cache holds Lucene Document objects (the stored fields for each document). Since Lucene internal document IDs are transient, this cache is not auto-warmed. The size for the `documentCache` should always be greater than `max_results` times the `max_concurrent_queries`, to ensure that Solr does not need to refetch a document during a request. The more fields you store in your documents, the higher the memory usage of this cache will be.

```
<documentCache class="solr.LRUCache"
    size="512"
    initialSize="512"
    autowarmCount="0" />
```

User Defined Caches

You can also define named caches for your own application code to use. You can locate and use your cache object by name by calling the `SolrIndexSearcher` methods `getCache()`, `cacheLookup()` and `cacheInsert()`. If you want auto-warming of your cache, include a `regenerator` attribute with the fully qualified name of a class that implements `solr.search.CacheRegenerator`.

```
<cache name="myUserCache" class="solr.LRUCache"
    size="4096"
    initialSize="1024"
    autowarmCount="1024"
    regenerator="org.mycompany.mypackage.MyRegenerator" />
```

Query Sizing and Warming

maxBooleanClauses

This sets the maximum number of clauses allowed in a boolean query. This can affect range or prefix queries that expand to a query with a large number of boolean terms. If this limit is exceeded, an exception is thrown.

```
<maxBooleanClauses>1024</maxBooleanClauses>
```

- ➥ This option modifies a global property that effects all Solr cores. If multiple `solrconfig.xml` files disagree on this property, the value at any point in time will be based on the last Solr core that was initialized.

enableLazyFieldLoading

If this parameter is set to true, then fields that are not directly requested will be loaded lazily as needed. This can boost performance if the most common queries only need a small subset of fields, especially if infrequently accessed fields are large in size.

```
<enableLazyFieldLoading>true</enableLazyFieldLoading>
```

useFilterForSortedQuery

This parameter configures Solr to use a filter to satisfy a search. If the requested sort does not include "score", the `filterCache` will be checked for a filter matching the query. For most situations, this is only useful if the same search is requested often with different sort options and none of them ever use "score".

```
<useFilterForSortedQuery>true</useFilterForSortedQuery>
```

queryResultWindowSize

Used with the `queryResultCache`, this will cache a superset of the requested number of document IDs. For example, if the a search in response to a particular query requests documents 10 through 19, and `queryWindowSize` is 50, documents 0 through 49 will be cached.

```
<queryResultWindowSize>20</queryResultWindowSize>
```

queryResultMaxDocsCached

This parameter sets the maximum number of documents to cache for any entry in the `queryResultCache`.

```
<queryResultMaxDocsCached>200</queryResultMaxDocsCached>
```

useColdSearcher

This setting controls whether search requests for which there is not a currently registered searcher should wait for a new searcher to warm up (false) or proceed immediately (true). When set to "false", requests will block until the searcher has warmed its caches.

```
<useColdSearcher>false</useColdSearcher>
```

maxWarmingSearchers

This parameter sets the maximum number of searchers that may be warming up in the background at any given time. Exceeding this limit will raise an error. For read-only slaves, a value of two is reasonable. Masters should probably be set a little higher.

```
<maxWarmingSearchers>2</maxWarmingSearchers>
```

Query-Related Listeners

As described in the section on [Caches](#), new Index Searchers are cached. It's possible to use the triggers for listeners to perform query-related tasks. The most common use of this is to define queries to further "warm" the Index Searchers while they are starting. One benefit of this approach is that field caches are pre-populated for faster sorting.

Good query selection is key with this type of listener. It's best to choose your most common and/or heaviest queries and include not just the keywords used, but any other parameters such as sorting or filtering requests.

There are two types of events that can trigger a listener. A `firstSearcher` event occurs when a new searcher is being prepared but there is no current registered searcher to handle requests or to gain auto-warming data from (i.e., on Solr startup). A `newSearcher` event is fired whenever a new searcher is being prepared and there is a current searcher handling requests.

The listener is always instantiated with the class `solr.QuerySenderListener`, and followed a `NamedList` array. These examples are included with `solrconfig.xml`:

```
<listener event="newSearcher" class="solr.QuerySenderListener">
  <arr name="queries">
    <!--
      <lst><str name="q">solr</str><str name="sort">price asc</str></lst>
      <lst><str name="q">rocks</str><str name="sort">weight asc</str></lst>
    -->
  </arr>
</listener>

<listener event="firstSearcher" class="solr.QuerySenderListener">
  <arr name="queries">
    <lst><str name="q">static firstSearcher warming in solrconfig.xml</str></lst>
  </arr>
</listener>
```

-  The above code sample is the default in `solrconfig.xml`, and a key best practice is to modify these defaults before taking your application to production. While the sample queries are commented out in the section for the "newSearcher", the example is not commented out for the "firstSearcher" event. There is no point in auto-warming your Index Searcher with the query string "static firstSearcher warming in solrconfig.xml" if that is not relevant to your search application.

RequestDispatcher in SolrConfig

The `requestDispatcher` element of `solrconfig.xml` controls the way the Solr servlet's RequestDispatcher implementation responds to HTTP requests. Included are parameters for defining if it should handle `/select` urls (for Solr 1.1 compatibility), if it will support remote streaming, the maximum size of file uploads and how it will respond to HTTP cache headers in requests.

Topics in this section:

- [handleSelect Element](#)
- [requestParsers Element](#)
- [httpCaching Element](#)

handleSelect Element

 `handleSelect` is for legacy back-compatibility; those new to Solr do not need to change anything about the way this is configured by default.

The first configurable item is the `handleSelect` attribute on the `<requestDispatcher>` element itself. This attribute can be set to one of two values, either "true" or "false". It governs how Solr responds to requests such as `/select?qt=xxx`. The default value "false" will ignore requests to `{/select}` if a requestHandler is not explicitly registered with the name `/select`. A value of "true" will route query requests to the parser defined with the `qt` value.

In recent versions of Solr, a `/select` requestHandler is defined by default, so a value of "false" will work fine. See the section [RequestHandlers and SearchComponents in SolrConfig](#) for more information.

```
<requestDispatcher handleSelect="true" >
  ...
</requestDispatcher>
```

requestParsers Element

The `<requestParsers>` sub-element controls values related to parsing requests. This is an empty XML element that doesn't have any content, only attributes.

The attribute `enableRemoteStreaming` controls whether remote streaming of content is allowed. If set to `false`, streaming will not be allowed. Setting it to `true` (the default) lets you specify the location of content to be streamed using `stream.file` or `stream.url` parameters.

If you enable remote streaming, be sure that you have authentication enabled. Otherwise, someone could potentially gain access to your content by accessing arbitrary URLs. It's also a good idea to place Solr behind a firewall to prevent it being accessed from untrusted clients.

The attribute `multipartUploadLimitInKB` sets an upper limit in kilobytes on the size of a document that may be submitted in a multi-part HTTP POST request. The value specified is multiplied by 1024 to determine the size in bytes.

The attribute `formdataUploadLimitInKB` sets a limit in kilobytes on the size of form data (`application/x-www-form-urlencoded`) submitted in a HTTP POST request, which can be used to pass request parameters that will not fit in a URL.

```
<requestParsers enableRemoteStreaming="true"
                 multipartUploadLimitInKB="2048000"
                 formdataUploadLimitInKB="2048" />
```

httpCaching Element

The `<httpCaching>` element controls HTTP cache control headers. Do not confuse these settings with Solr's internal cache configuration. This element controls caching of HTTP responses as defined by the W3C HTTP specifications.

This element allows for three attributes and one sub-element. The attributes of the `<httpCaching>` element control whether a 304 response to a GET request is allowed, and if so, what sort of response it should be. When an HTTP client application issues a GET, it may optionally specify that a 304 response is acceptable if the resource has not been modified since the last time it was fetched.

Parameter	Description
<code>never304</code>	If present with the value <code>true</code> , then a GET request will never respond with a 304 code, even if the requested resource has not been modified. When this attribute is set to true, the next two attributes are ignored. Setting this to true is handy for development, as the 304 response can be confusing when tinkering with Solr responses through a web browser or other client that supports cache headers.

lastModFrom	This attribute may be set to either <code>openTime</code> (the default) or <code>dirLastMod</code> . The value <code>openTime</code> indicates that last modification times, as compared to the <code>If-Modified-Since</code> header sent by the client, should be calculated relative to the time the Searcher started. Use <code>dirLastMod</code> if you want times to exactly correspond to when the index was last updated on disk.
etagSeed	This value of this attribute is sent as the value of the <code>ETag</code> header. Changing this value can be helpful to force clients to re-fetch content even when the indexes have not changed---for example, when you've made some changes to the configuration.

```
<httpCaching never304="false"
    lastModFrom="openTime"
    etagSeed="Solr">
    <cacheControl>max-age=30, public</cacheControl>
</httpCaching>
```

cacheControl Element

In addition to these attributes, `<httpCaching>` accepts one child element: `<cacheControl>`. The content of this element will be sent as the value of the Cache-Control header on HTTP responses. This header is used to modify the default caching behavior of the requesting client. The possible values for the Cache-Control header are defined by the HTTP 1.1 specification in [Section 14.9](#).

Setting the `max-age` field controls how long a client may re-use a cached response before requesting it again from the server. This time interval should be set according to how often you update your index and whether or not it is acceptable for your application to use content that is somewhat out of date. Setting `must-revalidate` will tell the client to validate with the server that its cached copy is still good before re-using it. This will ensure that the most timely result is used, while avoiding a second fetch of the content if it isn't needed, at the cost of a request to the server to do the check.

RequestHandlers and SearchComponents in SolrConfig

After the `<query>` section, request handlers and search components are configured. These are often referred to as "requestHandler" and "searchComponent", which is how they are defined in `solrconfig.xml`.

A *request handler* processes requests coming to Solr. These might be query requests or index update requests. You will likely need several of these defined, depending on how you want Solr to handle the various requests you will make.

A *search component* is a feature of search, such as highlighting or facetting. The search component is defined in `solrconfig.xml` separate from the request handlers, and then registered with a request handler as needed.

Topics covered in this section:

- Request Handlers
 - SearchHandlers
 - UpdateRequestHandlers
 - ShardHandlers
 - Other Request Handlers
- Search Components
 - Default Components
 - First-Components and Last-Components
 - Other Useful Components
- Related Topics

Request Handlers

Every request handler is defined with a name and a class. The name of the request handler is referenced with the request to Solr. For example, a request to

`http://localhost:8983/solr/collection1` is the default address for Solr, which will likely bring up the Solr Admin UI. However, add "/select" to the end, you can make a query:

```
http://localhost:8983/solr/collection1/select?q=solr
```

This query will be processed by the request handler with the name "/select". We've only used the "q" parameter here, which includes our query term, a simple keyword of "solr". If the request handler has more parameters defined, those will be used with any query we send to this request handler unless they are over-ridden by the client (or user) in the query itself.

If you have another request handler defined, you would send your request with that name - for example, "/update" is a request handler that handles index updates like sending new documents to the index.

SearchHandlers

The primary request handler defined with Solr by default is the "SearchHandler", which handles search queries. The request handler is defined, and then a list of defaults for the handler are defined with a `defaults` list.

For example, in the default `solrconfig.xml`, the first request handler defined looks like this:

```
<requestHandler name="/select" class="solr.SearchHandler">
  <lst name="defaults">
    <str name="echoParams">explicit</str>
    <int name="rows">10</int>
    <str name="df">text</str>
  </lst>
</requestHandler>
```

This example defines the `rows` parameter, which defines how many search results to return, to "10". The default field to search is the "text" field, set with the `df` parameter. The `echoParams` parameter defines that the parameters defined in the query should be returned when debug information is returned. Note also that the way the defaults are defined in the list varies if the parameter is a string, an integer, or another type.

All of the parameters described in the section on [searching](#) can be defined as defaults for any of the SearchHandlers.

Besides `defaults`, there are other options for the SearchHandler, which are:

- `appends`: This allows definition of parameters that are added to the user query. These might be [filter queries](#), or other query rules that should be added to each query. There is no mechanism in Solr to allow a client to override these additions, so you should be absolutely sure you always want these parameters applied to queries.

```
<lst name="appends">
  <str name="fq">inStock:true</str>
</lst>
```

In this example, the filter query "inStock:true" will always be added to every query.

- **invariants:** This allows definition of parameters that cannot be overridden by a client. The values defined in an `invariants` section will always be used regardless of the values specified by the user, by the client, in `defaults` or in `appends`.

```
<lst name="invariants">
  <str name="facet.field">cat</str>
  <str name="facet.field">manu_exact</str>
  <str name="facet.query">price:[* TO 500]</str>
  <str name="facet.query">price:[500 TO *]</str>
</lst>
```

In this example, facet fields have been defined which limits the facets that will be returned by Solr. If the client requests facets, the facets defined with a configuration like this are the only facets they will see.

The final section of a request handler definition is `components`, which defines a list of search components that can be used with a request handler. They are only registered with the request handler. How to define a search component is discussed further on in the section on [Search Components](#). The `components` element can only be used with a request handler that is a `SearchHandler`.

The `solrconfig.xml` file includes many other examples of `SearchHandlers` that can be used or modified as needed.

UpdateRequestHandlers

The `UpdateRequestHandlers` are request handlers which process updates to the index.

In this guide, we've covered these handlers in detail in the section [Uploading Data with Index Handlers](#).

ShardHandlers

It is possible to configure a request handler to search across shards of a cluster, used with distributed search. More information about distributed search and how to configure the `shardHandler` is in the section [Distributed Search with Index Sharding](#).

Other Request Handlers

There are other request handlers defined in `solrconfig.xml`, covered in other sections of this guide:

- [RealTime Get](#)
- [Index Replication](#)
- [Ping](#)

Search Components

The search components define the logic that is used by the SearchHandler to perform queries for users.

Default Components

There are several default search components that work with all SearchHandlers without any additional configuration. If no components are defined, these are used by default.

Component Name	Class Name	More Information
query	solr.QueryComponent	Described in the section Query Syntax and Parsing .
facet	solr.FacetComponent	Described in the section Faceting .
mlt	solr.MoreLikeThisComponent	Described in the section MoreLikeThis .
highlight	solr.HighlightComponent	Described in the section Highlighting .
stats	solr.StatsComponent	Described in the section The Stats Component .
debug	solr.DebugComponent	Described in the section on Common Query Parameters .

If you register a new search component with one of these default names, the newly defined component will be used instead of the default.

First-Components and Last-Components

It's possible to define some components as being used before (with `first-components`) or after (with `last-components`) other named components. This would be useful if custom search components have been configured to process data before the regular components are used. This is used when registering the components with the request handler.

```
<arr name="first-components">
  <str>mycomponent</str>
</arr>

<arr name="components">
  <str>query</str>
  <str>facet</str>
  <str>mlt</str>
  <str>highlight</str>
  <str>spellcheck</str>
  <str>stats</str>
  <str>debug</str>
</arr>
```

Other Useful Components

Many of the other useful components are described in sections of this Guide for the features they support. These are:

- `SpellCheckComponent`, described in the section [Spell Checking](#).
- `TermVectorComponent`, described in the section [The Term Vector Component](#).
- `QueryElevationComponent`, described in the section [The Query Elevation Component](#).
- `TermsComponent`, described in the section [The Terms Component](#).

Related Topics

- [SolrRequestHandler](#) from the Solr Wiki.
- [SearchHandler](#) from the Solr Wiki.
- [SearchComponent](#) from the Solr Wiki.

Core Admin and Configuring solr.xml

Use `solr.xml` to configure your Solr core (a logical index and associated configuration files), or to configure multiple cores. You can find `solr.xml` in your Solr Home directory. The default `solr.xml` file looks like this:

```
<solr persistent="true">
  <cores adminPath="/admin/cores"
    defaultCoreName="collection1" host="${host:}"
    hostPort="${jetty.port:}"
    hostContext="${hostContext:}"
    zkClientTimeout="${zkClientTimeout:15000}">
    <core name="collection1" instanceDir="collection1" />
  </cores>
</solr>
```

Covered in this section:

- [Using Multiple SolrCores](#)
- [Solr.xml Parameters](#)
- [CoreAdminHandler](#)

For more information about core configuration and `solr.xml`, see <http://wiki.apache.org/solr/CoreAdmin>.

Using Multiple SolrCores

It is possible to segment Solr into multiple cores, each with its own configuration and indices. Cores may be dedicated to a single application or to very different ones, but all are administered through a common administration interface. You can create new Solr cores on the fly, shutdown cores, even replace one running core with another, all without ever stopping or restarting your servlet container.

Solr cores are configured by placing a file named `solr.xml` in your `solr.home` directory. A typical `solr.xml` looks like this:

```
<solr persistent="false">
  <cores adminPath="/admin/cores" host="${host:}" hostPort="${jetty.port:}">
    <core name="core0" instanceDir="core0" />
    <core name="core1" instanceDir="core1" />
  </cores>
</solr>
```

This sets up two Solr cores, named "core0" and "core1", and names the directories (relative to the Solr installation path) which will store the configuration and data sub-directories.

- ⓘ You can run Solr without configuring any cores.

Solr.xml Parameters

The <solr> Element

There are two attributes that you can specify on <solr>, which is the root element of `solr.xml`.

Attribute	Description
<code>persistent</code>	Indicates that changes made through the API or admin UI should be saved back to this <code>solr.xml</code> . If not <code>true</code> , any runtime changes will be lost on the next Solr restart. The servlet container running Solr must have sufficient permissions to replace <code>solr.xml</code> (file delete and create), or errors will result. Any comments in <code>solr.xml</code> are not preserved when the file is updated. The default is <code>true</code> .
<code>sharedLib</code>	Specifies the path to a common library directory that will be shared across all cores. Any JAR files in this directory will be added to the search path for Solr plugins. This path is relative to the top-level container's Solr Home.

- ⚠ If you set the `persistent` attribute to `true`, be sure that the Web server has permission to replace the file. If the permissions are set incorrectly, the server will generate 500 errors and throw IOExceptions. Also, note that any comments in the `solr.xml` file will be lost when the file is overwritten.

The <cores> Element

The <cores> element, which contains definitions for each Solr core, is a child of <solr> and accepts several attributes of its own.

Attribute	Description
<code>adminPath</code>	This is the relative URL path to access the SolrCore administration pages. For example, a value of <code>/admin/cores</code> means that you can access the CoreAdminHandler with a URL that looks like this: http://localhost:8983/solr/admin/cores . If this attribute is not present, then SolrCore administration will not be possible.
<code>host</code>	The hostname Solr uses to access cores.

hostPort	The port Solr uses to access cores. In the default <code>solr.xml</code> file, this is set to <code> \${jetty.port:}</code> , which will use the Solr port defined in Jetty.
hostContext	The servlet context path.
zkClientTimeout	A timeout for connection to a ZooKeeper server. It is used with SolrCloud .
defaultCoreName	The name of a core that will be used for requests that do not specify a core.
transientCacheSize	Defines how many cores with <code>transient=true</code> that can be loaded before swapping the least recently used core for a new core. New in Solr 4.1.
shareSchema	This attribute, when set to <code>true</code> , ensures that the multiple cores pointing to the same <code>schema.xml</code> will be referring to the same <code>IndexSchema</code> Object. Sharing the <code>IndexSchema</code> Object makes loading the core faster. If you use this feature, make sure that no core-specific property is used in your <code>schema.xml</code> .
adminHandler	If used, this attribute should be set to the FQN (Fully qualified name) of a class that inherits from <code>CoreAdminHandler</code> . For example, <code>adminHandler="com.myorg.MyAdminHandler"</code> would configure the custom admin handler (<code>MyAdminHandler</code>) to handle admin requests. If this attribute isn't set, Solr uses the default admin handler, <code>org.apache.solr.handler.admin.CoreAdminHandler</code> . For more information on this parameter, see the Solr Wiki at http://wiki.apache.org/solr/CoreAdmin#cores .

The `<core>` Element

There is one `<core>` element for each SolrCore you define. They are children of the `<cores>` element and each one accepts six attributes.

Attribute	Description
name	The name of the SolrCore. You'll use this name to reference the SolrCore when running commands with the <code>CoreAdminHandler</code> .
instanceDir	This relative path defines the Solr Home for the core.
config	The configuration file name for a given core. The default is <code>solrconfig.xml</code> .
schema	The schema file name for a given core. The default is <code>schema.xml</code>
dataDir	This relative path defines the Solr Home for the core.
properties	The name of the properties file for this core. The value can be an absolute pathname or a path relative to the value of <code>instanceDir</code> .

transient	If true , the core can be unloaded if Solr reaches the <code>transientCacheSize</code> . The default if not specified is false . Cores are unloaded in order of least recently used first. New in Solr 4.1.
loadOnStartup	If true , the default if it is not specified, the core will loaded when Solr starts. New in Solr 4.1.

Properties in solr.xml

You can define properties in `solr.xml` that you may then reference in `solrconfig.xml` and `schema.xml`. Properties are name/value pairs. The scope of a property depends on which element it occurs within.

If a property is declared under `<solr>` but outside a `<core>` element, then it will have container scope and will be visible to all cores. In the example above, `productname` is such a property.

If a property declaration occurs within a `<core>` element, then its scope is limited to that core and it will not be visible to other cores. A property at core scope will override one of the same name declared at container scope.

```

<solr persistent="true" sharedLib="lib">
  <property name="productname" value="Acme Online"/>
  <cores adminPath="/admin/cores">
    <core name="core0" instanceDir="core0">
      <property name="dataDir" value="/data/core0"/></core>
      <core name="core1" instanceDir="core1"/>
    </cores>
  </solr>

```

In addition to any properties you declare at the core level, there are several properties that Solr defines automatically for each core. These properties are described in the table below:

Property	Description
<code>solr.core.name</code>	The core's name, as defined by the "name" attribute.
<code>solr.core.instanceDir</code>	The core's instance directory under which that its <code>conf/</code> and <code>data/</code> directories are located, derived from the core's <code>instanceDir</code> attribute.
<code>solr.core.dataDir</code>	The core's data directory, <code> \${solr.core.instanceDir}/data</code> by default.
<code>solr.core.configName</code>	The name of the core's configuration file, <code>solrconfig.xml</code> by default.
<code>solr.core.schemaName</code>	The name of the core's schema file, <code>schema.xml</code> by default.

Any of the above properties can be referenced by name in `schema.xml` or `solrconfig.xml`.

When defining properties, you can assign a property a default value that will be used if another value isn't specified. For example:

Without a default value, result will be empty if the property is not defined:

```
 ${productname}
```

With a default value:

```
 ${productname:SearchCo_MegaIndex}
```

CoreAdminHandler

The CoreAdminHandler is a special SolrRequestHandler that is used to manage Solr cores. Unlike normal SolrRequestHandlers, the CoreAdminHandler is not attached to a single core. Instead, it manages all the cores running in a single Solr instance. Only one CoreAdminHandler exists for each top-level Solr instance.

To use the CoreAdminHandler, make sure that the `adminPath` attribute is defined on the `<cores>` element; otherwise you will not be able to make HTTP requests to perform Solr core administration.

The CoreAdminHandler supports seven different actions that may be invoked on the `adminPath` URL. The action to perform is named by the HTTP request parameter "action", with arguments for a specific action being provided as additional parameters.

All action names are uppercase. The action names are:

- [STATUS](#)
- [CREATE](#)
- [RELOAD](#)
- [RENAME](#)
- [ALIAS](#)
- [SWAP](#)
- [UNLOAD](#)

These actions are described in detail in the sections below.

STATUS

The STATUS action returns the status of all running Solr cores, or status for only the named core.

```
http://localhost:8983/solr/admin/cores?action=STATUS
```

```
http://localhost:8983/solr/admin/cores?action=STATUS&core=core0
```

The STATUS action accepts one optional parameter:

Parameter	Description
core	(Optional) The name of a core, as listed in the "name" attribute of a <core> element in solr.xml.
indexInfo	If false , information about the index will not be returned with a core STATUS request. In Solr implementations with a large number of cores (i.e., more than hundreds), retrieving the index information for each core can take a lot of time and isn't always required. New in Solr 4.1.

CREATE

The CREATE action creates a new core and registers it. If persistence is enabled (`persistent="true"` on the `<solr>` element), the updated configuration for this new core will be saved in `solr.xml`. If a Solr core with the given name already exists, it will continue to handle requests while the new core is initializing. When the new core is ready, it will take new requests and the old core will be unloaded.

```
http://localhost:8983/solr/admin/cores?action=CREATE
    &name=coreX&instanceDir=path/to/dir
    &config=config_file_name.xml&schema=schem_file_name.xml&dataDir=data
```

The CREATE accepts the two mandatory parameters, as well as five optional parameters.

Parameter	Description
name	The name of the new core. Same as "name" on the <code><core></code> element.
instanceDir	The directory where files for this SolrCore should be stored. Same as <code>instanceDir</code> on the <code><core></code> element.
config	(Optional) Name of the config file (<code>solrconfig.xml</code>) relative to <code>instanceDir</code> .
schema	(Optional) Name of the schema file (<code>schema.xml</code>) relative to <code>instanceDir</code> .
datadir	(Optional) Name of the data directory relative to <code>instanceDir</code> .
collection	(Optional) The name of the collection to which this core belongs. The default is the name of the core. <code>collection.<param>=<value></code> causes a property of <code><param>=<value></code> to be set if a new collection is being created. Use <code>collection.configName=<configname></code> to point to the configuration for a new collection.
shard	(Optional) The shard id this core represents. Normally you want to be auto-assigned a shard id.

Use {{ collection.configName=<configname>}} to point to the config for a new collection.

For example: curl

```
'http://localhost:8983/solr/admin/cores?action=CREATE&name=mycore&collection=collect
```

RELOAD

The RELOAD action loads a new core from the configuration of an existing, registered Solr core.

While the new core is initializing, the existing one will continue to handle requests. When the new Solr core is ready, it takes over and the old core is unloaded.

This is useful when you've made changes to a Solr core's configuration on disk, such as adding new field definitions. Calling the RELOAD action lets you apply the new configuration without having to restart the Web container. However the Core Container does not persist the SolrCloud solr.xml parameters, such as solr/@zkHost and solr/cores/@hostPort, which are ignored.

```
http://localhost:8983/solr/admin/cores?action=RELOAD&core=core0
```

The RELOAD action accepts a single parameter

Parameter	Description
core	The name of the core to be reloaded.

RENAME

The RENAME action changes the name of a Solr core.

```
http://localhost:8983/solr/admin/cores?action=RENAME  
&core=core0&other=core5
```

The RENAME action requires the following two parameter:

Parameter	Description
core	The name of the Solr core to be renamed.
other	The new name for the Solr core. If the persistent attribute of <solr> is true, the new name will be written to solr.xml as the name attribute of the <core> attribute.

ALIAS

The ALIAS action establishes an additional name by which a SolrCore may be referenced. Subsequent actions may use the Solr core's original name or any of its aliases.

```
http://localhost:8983/solr/admin/cores?action=ALIAS&core=coreX&other=coreY
```

The ALIAS action requires two parameters:

Parameter	Description
core	The name or alias of an existing core.
other	The additional name by which this core should be known.

SWAP

SWAP atomically swaps the names used to access two existing Solr cores. This can be used to swap new content into production. The prior core remains available and can be swapped back, if necessary. Each core will be known by the name of the other, after the swap.

`http://localhost:8983/solr/admin/cores?action=SWAP&core=core1&other=core0`

The SWAP action requires two parameters, which are described in the table below.

Parameter	Description
core	The name of one of the cores to be swapped.
other	The name of one of the cores to be swapped.

UNLOAD

The UNLOAD action removes a core from Solr. Active requests will continue to be processed, but no new requests will be sent to the named core. If a core is registered under more than one name, only the given name is removed.

`http://localhost:8983/solr/admin/cores?action=UNLOAD&core=core0`

The UNLOAD action requires a parameter (`core`) identifying the core to be removed. If the persistent attribute of `<solr>` is set to true, the `<core>` element with this name attribute will be removed from `solr.xml`.



⚠️ Unloading all cores in a SolrCloud collection causes the removal of that collection's metadata from ZooKeeper.

Solr Plugins

Solr allows you to load custom code to perform a variety of tasks within Solr, from custom Request Handlers to process your searches, to custom Analyzers and Token Filters for your text field. You can even load custom Field Types. These pieces of custom code are called plugins.

Not everyone will need to create plugins for their Solr instances - what's provided is usually enough for most applications. However, if there's something that you need, you may want to review the Solr Wiki documentation on plugins at [SolrPlugins](#).

JVM Settings

Configuring your JVM can be a complex topic. A full discussion is beyond the scope of this document. Luckily, most modern JVMs are quite good at making the best use of available resources with default settings. The following sections contain a few tips that may be helpful when the defaults are not optimal for your situation.

For more general information about improving Solr performance, see
<https://wiki.apache.org/solr/SolrPerformanceFactors>.

Choosing Memory Heap Settings

The most important JVM configuration settings are those that determine the amount of memory it is allowed to allocate. There are two primary command-line options that set memory limits for the JVM. These are `-Xms`, which sets the initial size of the JVM's memory heap, and `-Xmx`, which sets the maximum size to which the heap is allowed to grow.

If your Solr application requires more heap space than you specify with the `-Xms` option, the heap will grow automatically. It's quite reasonable to not specify an initial size and let the heap grow as needed. The only downside is a somewhat slower startup time since the application will take longer to initialize. Setting the initial heap size higher than the default may avoid a series of heap expansions, which often results in objects being shuffled around within the heap, as the application spins up.

The maximum heap size, set with `-Xmx`, is more critical. If the memory heap grows to this size, object creation may begin to fail and throw `OutOfMemoryException`. Setting this limit too low can cause spurious errors in your application, but setting it too high can be detrimental as well.

It doesn't always cause an error when the heap reaches the maximum size. Before an error is raised, the JVM will first try to reclaim any available space that already exists in the heap. Only if all garbage collection attempts fail will your application see an exception. As long as the maximum is big enough, your app will run without error, but it may run more slowly if forced garbage collection kicks in frequently.

The larger the heap the longer it takes to do garbage collection. This can mean minor, random pauses or, in extreme cases, "freeze the world" pauses of a minute or more. As a practical matter, this can become a serious problem for heap sizes that exceed about two gigabytes, even if far more physical memory is available. On robust hardware, you may get better results running multiple JVMs, rather than just one with a large memory heap. Some specialized JVM implementations may have customized garbage collection algorithms that do better with large heaps. Also, Java 7 is expected to have a redesigned GC that should handle very large heaps efficiently. Consult your JVM vendor's documentation.

When setting the maximum heap size, be careful not to let the JVM consume all available physical memory. If the JVM process space grows too large, the operating system will start swapping it, which will severely impact performance. In addition, the operating system uses memory space not allocated to processes for file system cache and other purposes. This is especially important for I/O-intensive applications, like Lucene/Solr. The larger your indexes, the more you will benefit from filesystem caching by the OS. It may require some experimentation to determine the optimal tradeoff between heap space for the JVM and memory space for the OS to use.

On systems with many CPUs/cores, it can also be beneficial to tune the layout of the heap and/or the behavior of the garbage collector. Adjusting the relative sizes of the generational pools in the heap can affect how often GC sweeps occur and whether they run concurrently. Configuring the various settings of how the garbage collector should behave can greatly reduce the overall performance impact when it does run. There is a lot of good information on this topic available on Sun's website. A good place to start is here: <http://java.sun.com/javase/technologies/hotspot/gc/>.

Use the Server HotSpot VM

If you are using Sun's JVM, add the `-server` command-line option when you start Solr. This tells the JVM that it should optimize for a long running, server process. If the Java runtime on your system is a JRE, rather than a full JDK distribution (including `javac` and other development tools), then it is possible that it may not support the `-server` JVM option. Test this by running `java -help` and look for `-server` as an available option in the displayed usage message.

Checking JVM Settings

A great way to see what JVM settings your server is using, along with other useful information, is to use the admin RequestHandler, `solr/admin/system`. This request handler will display a wealth of server statistics and settings.

You can also use any of the tools that are compatible with the Java Management Extensions (JMX). See the section *Using JMX with Solr* in [Managing Solr](#) for more information.

Managing Solr

This section describes how to run Solr and how to look at Solr when it is running. It contains the following sections:

[Running Solr on Jetty](#): Describes how to run Solr in the Jetty web application container. The Solr example included in this distribution runs in a Jetty web application container.

[Running Solr on Tomcat](#): Describes how to run Solr in the Tomcat web application container.

[Configuring Logging](#): Describes how to configure logging for Solr.

[Backing Up](#): Describes backup strategies for your Solr indexes.

[Using JMX with Solr](#): Describes how to use Java Management Extensions with Solr.

For information on running Solr in a variety of Java application containers, see the [basic installation instructions](#) on the Solr wiki.

Running Solr on Tomcat

Solr comes with an example schema and scripts for running on [Jetty](#). The next section describes some of the details of how things work "under the hood," and covers running multiple Solr instances and deploying Solr using the Tomcat application manager.

For more information about running Solr on Tomcat, see the [basic installation instructions](#) and the [Solr Tomcat](#) page on the Solr wiki.

How Solr Works with Tomcat

The two basic steps for running Solr in any Web application container are as follows:

1. Make the Solr classes available to the container. In many cases, the Solr Web application archive (WAR) file can be placed into a special directory of the application container. In the case of Tomcat, you need to place the Solr WAR file in Tomcat's `webapps` directory. If you installed Tomcat with Solr, take a look in `tomcat/webapps`: you'll see the `solr.war` file is already there.
2. Point Solr to the Solr home directory that contains `conf/solrconfig.xml` and `conf/schema.xml`. There are a few ways to get this done. One of the best is to define the `solr.solr.home` Java system property. With Tomcat, the best way to do this is via a shell environment variable, `JAVA_OPTS`. Tomcat puts the value of this variable on the command line upon startup. Here is an example:

```
export JAVA_OPTS="-Dsolr.solr.home=/Users/jonathan/Desktop/solr"
```

Port 8983 is the default Solr listening port. If you are using Tomcat and wish to change this port, edit the file `tomcat/conf/server.xml` in the Solr distribution. You'll find the port in this part of the file:

```
<Connector port="8983" protocol="HTTP/1.1" connectionTimeout="20000"
redirectPort="8443" />
```

Modify the port number as desired and restart Tomcat if it is already running.



Modifying the port number will leave some of the samples and help file links pointing to the default port. It is out of the scope of this reference guide to provide full details of how to change all of the examples and other resources to the new port.

Running Multiple Solr Instances

The standard way to deploy multiple Solr index instances in a single Web application is to use the multicore API described in [Using Multiple SolrCores](#).

An alternative approach, which provides more code isolation, uses Tomcat context fragments. A context fragment is a file that contains a single `<context>` element and any subelements required for your application. The file omits all other XML elements.

Each context fragment specifies where to find the Solr WAR and the path to the solr home directory. The name of the context fragment file determines the URL used to access that instance of Solr. For example, a context fragment named `harvey.xml` would deploy Solr to be accessed at `http://localhost:8983/harvey`.

In Tomcat's `conf/Catalina/localhost` directory, store one context fragment per instance of Solr. If the `conf/Catalina/localhost` directory doesn't exist, go ahead and create it.

Using Tomcat context fragments, you could run multiple instances of Solr on the same server, each with its own schema and configuration. For full details and examples of context fragments, take a look at the Solr Wiki: <http://wiki.apache.org/solr/SolrTomcat>.

Here are examples of context fragments which would set up two Solr instances, each with its own `solr.home`:

```
<Context docBase="/some/path/solr.war" debug="0" crossContext="true" >
    <Environment name="solr/home" type="java.lang.String" value="/some/path/solr1home"
override="true" />
</Context>
<Context docBase="/some/path/solr.war" debug="0" crossContext="true" >
    <Environment name="solr/home" type="java.lang.String"
value="/some/path/solr2home" override="true" />
</Context>
```

Deploying Solr with the Tomcat Manager

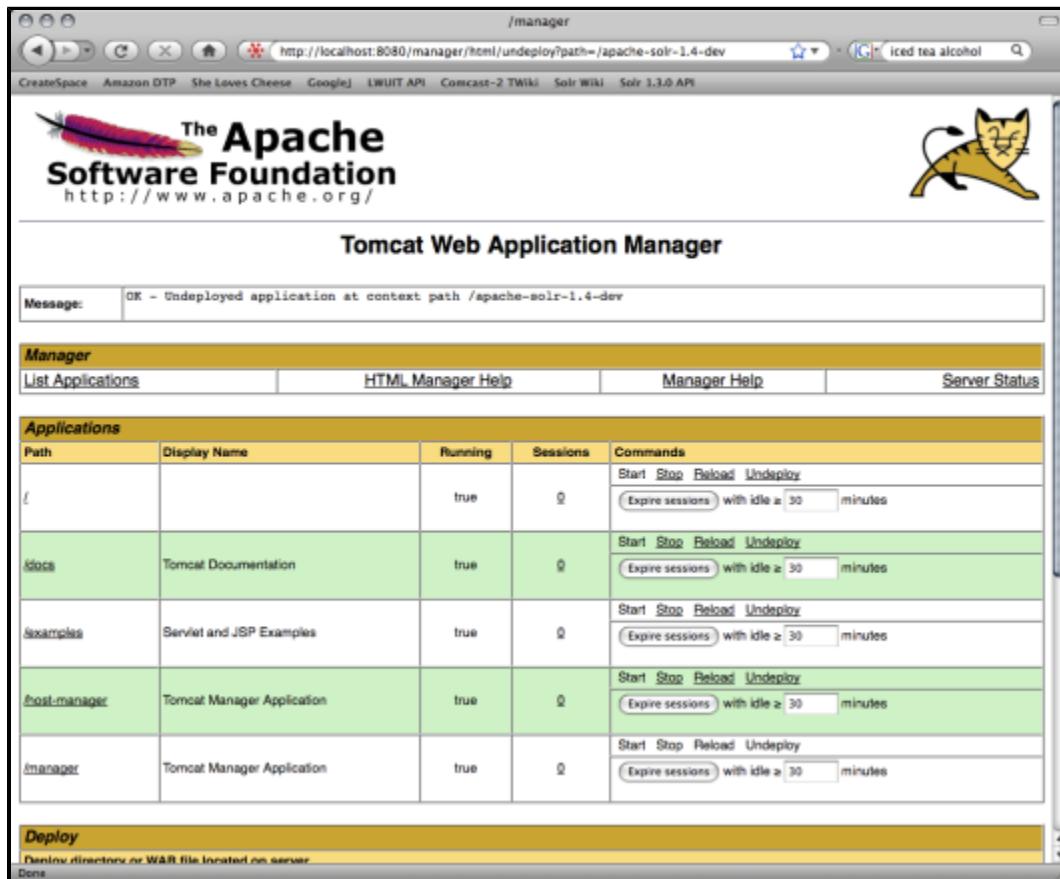
If your instance of Tomcat is running the Tomcat Web Application Manager, you can use its browser interface to deploy Solr.

Just as before, you have to tell Solr where to find the solr home directory. You can do this by setting `JAVA_OPTS` before starting Tomcat.

Once Tomcat is running, navigate to the Web application manager, probably available at a URL like this:

`http://localhost:8983/manager/html`

You will see the main screen of the manager.



The screenshot shows the Tomcat Web Application Manager interface. At the top, there's a banner for "The Apache Software Foundation" and a cartoon cat logo. Below the banner, the title "Tomcat Web Application Manager" is displayed. A message box shows "OK - Undeployed application at context path /apache-solr-1.4-dev". The main area is titled "Manager" and contains tabs for "List Applications", "HTML Manager Help", "Manager Help", and "Server Status". The "List Applications" tab is selected, showing a table of applications:

Path	Display Name	Running	Sessions	Commands
/		true	0	Start Stop Reload Undeploy <input type="button" value="Expire sessions"/> with idle ≥ 30 minutes
/docs	Tomcat Documentation	true	0	Start Stop Reload Undeploy <input type="button" value="Expire sessions"/> with idle ≥ 30 minutes
/examples	Servlet and JSP Examples	true	0	Start Stop Reload Undeploy <input type="button" value="Expire sessions"/> with idle ≥ 30 minutes
/host-manager	Tomcat Manager Application	true	0	Start Stop Reload Undeploy <input type="button" value="Expire sessions"/> with idle ≥ 30 minutes
/manager	Tomcat Manager Application	true	0	Start Stop Reload Undeploy <input type="button" value="Expire sessions"/> with idle ≥ 30 minutes

Below the table, there's a "Deploy" section with a note: "Deploying directory or WAR file located on server." and a "Done" button.

To add Solr, scroll down to the **Deploy** section, specifically **WAR file to deploy**. Click **Browse...** and find the Solr WAR file, usually something like `dist/apache-solr-3.x.0.war` within your Solr installation. Click **Deploy**. Tomcat will load the WAR file and start running it. Click the link in the application path column of the manager to see Solr. You won't see much, just a welcome screen, but it contains a link for the Admin Console.

Tomcat's manager screen, in its application list, has links so you can stop, start, reload, or undeploy the Solr application.

Running Solr on Jetty

Solr comes with an example schema and scripts for running on [Jetty](#), along with a working installation, in the `/example` directory. The included version of Jetty works well for small installations, but for more heavy-duty use, we recommend that you download the [full Jetty package](#), which includes additional modules ("JettyPlus").

For more information about the Jetty example installation, see the [Solr Tutorial](#) and the [basic installation instructions](#) on the Solr wiki.

For detailed information about running Solr on Jetty or JettyPlus, see <http://wiki.apache.org/solr/SolrJetty>.

Changing the Solr Listening Port

Port 8983 is the default port for Solr. If you are using Jetty and wish to change the port number, edit the file `jetty/etc/jetty.xml` in the Solr distribution. You'll find the port in this part of the file:

```
<New class="org.mortbay.jetty.bio.SocketConnector">
    <Set name="port">
        <SystemProperty name="jetty.port" default="8983"/>
    </Set>
    <Set name="maxIdleTime">50000</Set>
    <Set name="lowResourceMaxIdleTime">1500</Set>
</New>
```

Modify the port number as desired and restart Jetty if it is already running.



Modifying the port number will leave some of the samples and help file links pointing to the wrong port. It is out of the scope of this reference guide to provide full details of how to change all of the examples and other resources to the new port.

Configuring Logging

Logging is the practice of writing informative messages where they can be reviewed later. System administrators or developers can read logs to learn information about a system. If an application dies unexpectedly, the key to its demise might be written in a log somewhere. A canny developer can examine a log to understand what went wrong, much like a detective can examine the scene of a crime to find out what happened.

Solr uses the SLF4J Logging API (<http://www.slf4j.org>). If you want to see the log output in Tomcat, look in `solr/tomcat/logs`. You will find a file named something like `catalina.2011-05-01.log`, except with the current date.

For further information about Solr logging, see [Logging in Default Jetty Setup](#).

Temporary Logging Settings

You can control the amount of logging output in Solr by using the Admin Web interface. Select the **LOGGING** link. Note that this page only lets you change settings in the running system and is not saved for the next run. (For more information about the Admin Web interface, see [Using the Solr Administration User Interface](#).)

The screenshot shows the Apache Solr Admin Web interface. The left sidebar has a dark grey background with white text and icons. It includes links for Dashboard, Logging (which is highlighted in light grey), Level, Core Admin, Java Properties, Thread Dump, and a collection named "collection1". The main content area has a white background. At the top, it says "JUL (org.slf4j.impl.JDK14LoggerFactory)". Below that is a table with four columns: Time, Level, Logger, and Message. There is one row in the table: "15:58:19 WARNING SolrCore New index directory detected: old=null new=solr/collection1/data/index/". At the bottom of the content area, there is a small note: "Last Check: 16:01:06".

The Logging screen.

This part of the Admin Web interface allows you to set the logging level for many different log categories. Fortunately, any categories that are **unset** will have the logging level of its parent. This makes it possible to change many categories at once by adjusting the logging level of their parent.

When you select **Level**, you see the following menu:

Category	Current Level
root	UNSET
global	UNSET
javax	UNSET
management	UNSET
mbeanserver	UNSET
org	UNSET
apache	UNSET
http	UNSET
impl	UNSET
client	INFO
DefaultHttpClient	null
conn	UNSET
DefaultClientConnectionOperator	INFO
IdleConnectionHandler	INFO
tscmc	null
ConnPoolByRoute	INFO
ThreadSafeClientConnManager	INFO

The Log Level Menu.

Directories are shown with their current logging levels. The Log Level Menu floats over these. To set a log level for a particular directory, select it and click the appropriate log level button.

Log levels settings are as follows:

Level	Result
FINEST	Reports everything.
FINE	Reports everything but the least important messages.
CONFIG	Reports configuration errors.
INFO	Reports everything but normal status.
WARNING	Reports all warnings.
SEVERE	Reports only the most severe warnings.
OFF	Turns off logging.

UNSET	Removes the previous log setting.
-------	-----------------------------------

Multiple settings at one time are allowed.

Permanent Logging Settings

Making permanent changes to the JDK Logging API configuration is a matter of creating or editing a properties file.

Tomcat Logging Settings

Tomcat offers a choice between settings for all applications or settings specifically for the Solr application.

To change logging settings for Solr only, edit

`tomcat/webapps/solr/WEB-INF/classes/logging.properties`. You will need to create the classes directory and the logging.properties file. You can set levels from FINEST to SEVERE for a class or an entire package. Here are a couple of examples:

```
org.apache.commons.digester.Digester.level = FINEST
org.apache.solr.level = WARNING
```

Alternately, if you wish to change Tomcat's JDK Logging API settings for every application in this instance of Tomcat, edit `tomcat/conf/logging.properties`.

See the documentation for the SLF4J Logging API for more information:

<http://slf4j.org/docs.html>

<http://java.sun.com/javase/6/docs/technotes/guides/logging/index.html>

Jetty Logging Settings

To change settings for the SLF4J Logging API in Jetty, you need to create a settings file and tell Jetty where to find it.

Begin by creating a file `jetty/logging.properties`. Use the example lines above as a guide.

To tell Jetty how to find the file, edit `start.sh`. Find the line which launches Jetty, which looks something like this, except it will have an absolute path to `start.jar`:

```
java -DSTOP.PORT=8079 -DSTOP.KEY=secret -jar start.jar

#Add the location of the logging properties file like this:

java -Djava.util.logging.config.file=logging.properties
     -DSTOP.PORT=8079 -DSTOP.KEY=secret -jar start.jar
```

The next time you launch Jetty, it will use the settings in the file.

Backing Up

If you are worried about data loss, and of course you *should* be, you need a way to back up your Solr indexes so that you can recover quickly in case of catastrophic failure.

Making Backups with the Solr Replication Handler

The easiest way to make back-ups in Solr is to take advantage of the Replication Handler, which is described in detail in [Index Replication](#). The Replication Handler's primary purpose is to replicate an index on slave servers for load-balancing, but the Replication Handler can be used to make a back-up copy of a server's index, even if no slave servers are in operation.

Once you have configured the Replication Handler in `solrconfig.xml`, you can trigger a back-up with an HTTP command like this:

```
http://master_host/solr/replication?command=backup
```

For details on configuring the Replication Handler, see [Legacy Scaling and Distribution](#).

Backup Scripts from Earlier Solr Releases

Solr also provides shell scripts in the bin directory that make copies of the indexes. However, these scripts only work with a Linux-style shell, and not everybody in the world runs Linux.

The scripts themselves are relatively simple. Look in the bin directory of your Solr home directory, for example `example/solr/bin`. In particular, `backup.sh` makes a copy of Solr's index directory and gives it a name based on the current date.

This scripts include the following:

Script Name	Description
abc	Atomic Backup post-Commit tells the Solr server to perform a commit. A snapshot of the index directory is made after the commit if the Solr server is configured to do so (by enabling the postCommit event listener in <code>solr/conf/solrconfig.xml</code>). A backup of the most recent snapshot directory is then made if the commit is successful. Backup directories are named <code>yyyymmddHHMMSS</code> where <code>yyyymmddHHMMSS</code> is the timestamp of when the snapshot was taken.

abo	Atomic Backup post-Optimize tells the Solr server to perform an optimize. A snapshot of the index directory is made after the optimize if the Solr server is configured to do so (by enabling the postCommit or postOptimize event listener in <code>solr/conf/solrconfig.xml</code>). A backup of the most recent snapshot directory is then made if the optimize is successful. Backup directories are named <code>backup. yyyy-mm-dd HHMMSS</code> where <code>yyyy-mm-dd HHMMSS</code> is the timestamp of when the snapshot was taken.
backup	Backs up the index directory using hard links. Backup directories are named <code>backup. yyyy-mm-dd HHMMSS</code> where <code>yyyy-mm-dd HHMMSS</code> is the timestamp of when the backup was taken.
backupcleaner	Runs as a cron job to remove backups more than a configurable number of days old or all backups except for the most recent n number of backups. Also can be run manually.

For more details about backup scripts, see the Solr Wiki page

<http://wiki.apache.org/solr/SolrOperationsTools>.

Using JMX with Solr

Java Management Extensions (JMX) is a technology that makes it possible for complex systems to be controlled by tools without the systems and tools having any previous knowledge of each other. In essence, it is a standard interface by which complex systems can be viewed and manipulated.

Solr, like any other good citizen of the Java universe, can be controlled via a JMX interface. You can enable JMX support by adding lines to `solrconfig.xml`. You can use a JMX client, like `jconsole`, to connect with Solr. Check out the Wiki page <http://wiki.apache.org/solr/SolrJmx> for more information. You may also find the following overview of JMX to be useful: <http://java.sun.com/j2se/1.5.0/docs/guide/management/agent.html>.

Configuring JMX

JMX configuration is provided in `solrconfig.xml`. Please see the [JMX Technology Home Page](#) for more details.

A `rootName` attribute can be used when configuring `<jmx />` in `solrconfig.xml`. If this attribute is set, Solr uses it as the root name for all the MBeans that Solr exposes via JMX. The default name is "solr" followed by the core name.

 Enabling/disabling JMX and securing access to MBeanServers is left up to the user by specifying appropriate JVM parameters and configuration. Please explore the [JMX Technology Home Page](#) for more details.

Configuring an Existing MBeanServer

The command:

```
<jmx />
```

enables JMX support in Solr if and only if an existing MBeanServer is found. Use this if you want to configure JMX with JVM parameters. Remove this to disable exposing Solr configuration and statistics to JMX. If this is specified, Solr will try to list all available MBeanServers and use the first one to register MBeans.

Configuring an Existing MBeanServer with agentId

The command:

```
<jmx agentId="myMBeanServer" />
```

enables JMX support in Solr if and only if an existing MBeanServer is found matching the given agentId. If multiple servers are found, the first one is used. If none is found, an exception is raised and depending on the configuration, Solr may refuse to start.

Configuring a New MBeanServer

The command:

```
<jmx serviceUrl="service:jmx:rmi:///jndi/rmi://localhost:9999/solrjmx" />
```

creates a new MBeanServer exposed for remote monitoring at the specific service URL. If the JMXConnectorServer can't be started (probably because the serviceUrl is bad), an exception is thrown.

Example

Using the example jetty setup provided with Solr installation, the JMX support works like this in jconsole.png.

1. Run "ant example" to build the example war file.
2. Go to the example folder in the Solr installation and run the following command:

```
java -Dcom.sun.management.jmxremote -jar start.jar
```

3. Start jconsole (provided with the Sun JDK in the bin directory).
4. Connect to the "start.jar" shown in the list of local processes.
5. Switch to the "MBeans" tab. You should be able to see "solr" listed there.

Configuring a Remote Connection to Solr JMX

If you want to connect to Solr remotely, you need to pass in some extra parameters, documented here:

<http://docs.oracle.com/javase/1.5.0/docs/guide/management/agent.html>

If you are not able to connect from a remote machine, you may also need to specify the hostname of the Solr host by adding the following property as well:

 Making JMX connections into machines running behind NATs (e.g. Amazon's EC2 service) is not a simple task. The `java.rmi.server.hostname` system property may help, but running `jconsole` on the server itself and using a remote desktop is often the simplest solution. See <http://jmsbrdy.com/monitoring-java-applications-running-on-ec2-i>.

SolrCloud

Apache Solr includes the ability to set up a cluster of Solr servers that combines fault tolerance and high availability. Called **SolrCloud**, these capabilities provide distributed indexing and search capabilities, supporting the following features:

- Central configuration for the entire cluster
- Automatic load balancing and fail-over for queries
- ZooKeeper integration for cluster coordination and configuration.

SolrCloud is flexible distributed search and indexing, without a master node to allocate nodes, shards and replicas. Instead, Solr uses ZooKeeper to manage these locations, depending on configuration files and schemas. Documents can be sent to any server and ZooKeeper will figure it out.

In this section, we'll cover everything you need to know about using Solr in SolrCloud mode. We've split up the details into the following topics:

- [Getting Started with SolrCloud](#)
- [How SolrCloud Works](#)
 - [Shards and Indexing Data in SolrCloud](#)
 - [Distributed Requests](#)
 - [Read and Write Side Fault Tolerance](#)
 - [NRT, Replication, and Disaster Recovery with SolrCloud](#)
- [SolrCloud Configuration and Parameters](#)
 - [Using ZooKeeper to Manage Configuration Files](#)
 - [Managing Collections via the Collections API](#)
 - [Parameter Reference](#)
 - [Command Line Utilities](#)
 - [SolrCloud with Legacy Configuration Files](#)
- [SolrCloud Glossary](#)

You can also find more information on the [Solr wiki page on SolrCloud](#).

-  If upgrading an existing Solr 4.1 instance running with SolrCloud, be aware that the way the `name_node` parameter is defined has changed. This may cause a situation where the `name_node` uses the IP address of the machine instead of the server name, and thus SolrCloud is not aware of the existing node. If this happens, you can manually edit the `host` parameter in `solr.xml` to refer to the server name, or set the `host` in your system environment variables (since by default `solr.xml` is configured to inherit the `host` name from the environment variables). See also the section [Core Admin and Configuring solr.xml](#) for more information about the `host` parameter.

Getting Started with SolrCloud

SolrCloud is designed to provide a highly available, fault tolerant environment that can index your data for searching. It's a system in which data is organized into multiple pieces, or shards, that can be housed on multiple machines, with replicas providing redundancy for both scalability and fault tolerance, and a ZooKeeper server that helps manage the overall structure so that both indexing and search requests can be routed properly.

This section explains SolrCloud and its inner workings in detail, but before you dive in, it's best to have an idea of what it is you're trying to accomplish. This page provides a simple tutorial that explains how SolrCloud works on a practical level, and how to take advantage of its capabilities.

In this section you will learn:

- How to distribute data over multiple instances by using ZooKeeper and creating shards
- How to create redundancy for shards by using replicas
- How to create redundancy for the overall cluster by running multiple ZooKeeper instances

 This tutorial assumes that you're already familiar with the basics of using Solr. If you need a refresher, please visit the Solr Tutorial to get a grounding in Solr concepts.

Creating multiple shards

Creating a cluster with multiple shards involves two steps:

 Make sure to run Solr from the example directory in non-SolrCloud mode at least once before beginning; this process unpacks the jar files necessary to run SolrCloud. On the other hand, make sure also that there are no documents in the example directory before making copies.

1. Start the "overseer" node, which includes a ZooKeeper server to keep track of your cluster.
2. Start any remaining shard nodes and point them to the running ZooKeeper.

In this example, you'll create two separate Solr instances. Do this by making two copies of the example directory that is part of the Solr distribution:

```
cd <SOLR_DIST_HOME>
cp -r example shard1
cp -r example shard2
```

Next, start the first Solr instance, including the `-DzkRun` parameter, which also starts a local ZooKeeper instance:

```
cd shard1
java -DzkRun -DnumShards=2 -Dbootstrap_confdir=./solr/collection1/conf \
-Dcollection.configName=myconf -jar start.jar
```

Let's look at each of these parameters:

-DzkRun Starts up a ZooKeeper server embedded within Solr. This server will manage the cluster configuration.

-Dnumshards Determines how many pieces you're going to break your index into. In this case we're going to break the index into two pieces, or shards, so we're setting this value to 2. Note that once you start up a cluster, you cannot change this value. So if you expect to need more shards later on, build them into your configuration now (you can do this by starting all of your shards on the same server, then migrating them to different servers later).

-Dbootstrap_confdir ZooKeeper needs to get a copy of the cluster configuration, so this parameter tells it where to find that information.

-Dcollection.configName This parameter determines the name under which that configuration information is stored by ZooKeeper.

- ✓ The `-DnumShards`, `-Dbootstrap_confdir`, and `-Dcollection.configName` parameters need only be specified once, the first time you start Solr in SolrCloud mode. They load your configurations into ZooKeeper; if you run them again at a later time, they will re-load your configurations and may wipe out changes you have made.

OK, so at this point you have one sever running, but it represents only half the shards, so you will need to start the second one before you have a fully functional cluster. To do that, start the second instance in another window as follows:

```
cd shard2
java -Djetty.port=7574 -DzkHost=localhost:9983 -jar start.jar
```

Because this node isn't the overseer, the parameters are a bit less complex:

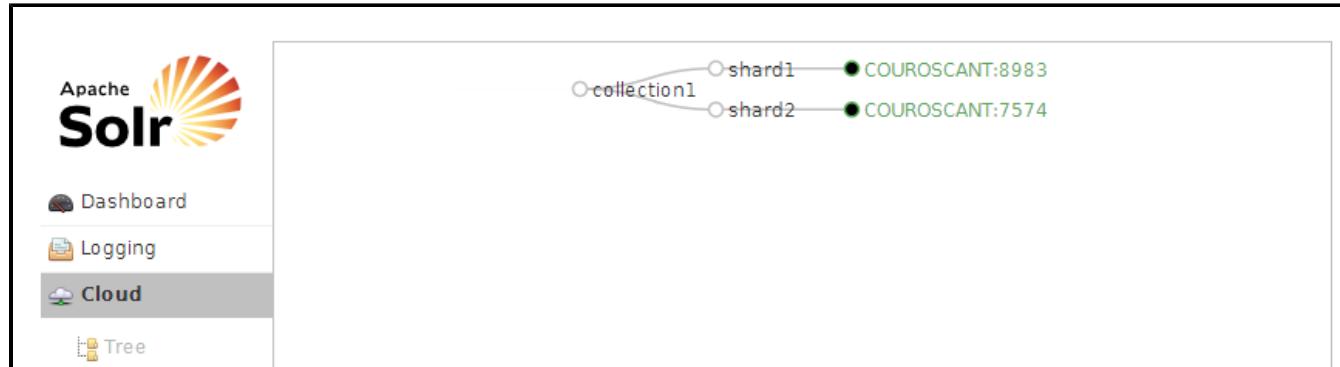
-Djetty.port The only reason we even have to set this parameter is because we're running both servers on the same machine, so they can't both use Jetty's default port. In this case we're choosing an arbitrary number that's different from the default.

-DzkHost This parameter tells Solr where to find the ZooKeeper server so that it can "report for duty". By default, the ZooKeeper server operates on the Solr port plus 1000. (Note that if you were running an external ZooKeeper server, you'd simply point to that.)

At this point you should have two Solr windows running, both being managed by ZooKeeper. To verify that, open the ZooKeeper administration page in your browser by clicking:

<http://localhost:8983/solr/#/~cloud>

You should see both shard1 and shard2, as in:



Now it's time to see the cluster in action. Start by indexing some data to one or both shards. You can do this any way you like, but the easiest way is to use the `exampledocs`, along with curl so that you can control which port (and thereby which server) gets the updates:

```
cd example/exampledocs
curl http://localhost:8983/solr/update?commit=true -H "Content-Type: text/xml" -d
"@mem.xml"
curl http://localhost:7574/solr/update?commit=true -H "Content-Type: text/xml" -d
"@monitor2.xml"
```

At this point each shard contains a subset of the data, but a search directed at either server should span both shards. For example, the following searches should both return the identical set of all results:

http://localhost:8983/solr/collection1/select?q=**

http://localhost:7574/solr/collection1/select?q=**

The reason that this works is that each shard knows about the other shards, so the search is carried out on all cores, then the results are combined and returned by the called server.

In this way you can have two cores or two hundred, with each containing a separate portion of the data.

- ✓ If you want to check the number of documents on each shard, you could add `distrib=false` to each query and your search would not span all shards.

But what about providing high availability, even if one of these servers goes down? To do that, you'll need to look at replicas.

Creating shard replicas

In order to provide high availability, you can create replicas, or copies of each shard that run in parallel with the main core for that shard. The architecture consists of the original shards, which are called the leaders, and their replicas, which contain the same data but let the leader handle all of the administrative tasks such as making sure data goes to all of the places it should go. This way, if one copy of the shard goes down, the data is still available and the cluster can continue to function.

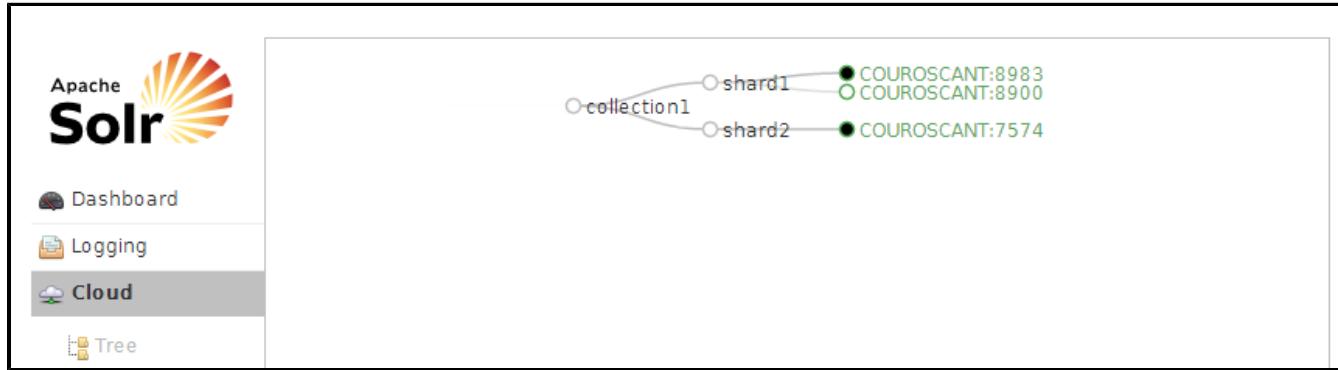
Start by creating two more fresh copies of the example directory:

```
cd <SOLR_DIST_HOME>
cp -r example shard1replica
cp -r example shard2replica
```

If you don't already have the two instances you created in the previous section up and running, go ahead and restart them. From there, it's simply a matter of adding additional instances. Start by adding `shard1replica`:

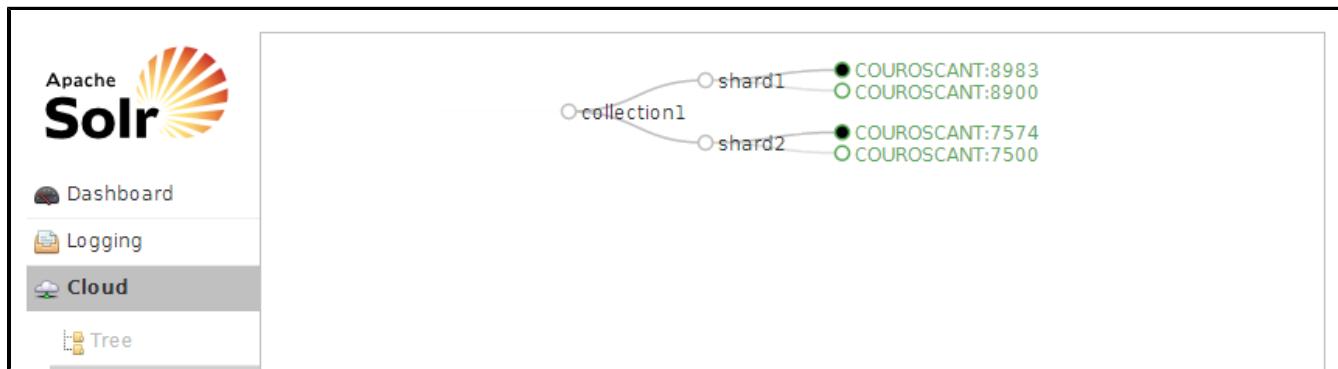
```
cd shard1replica
java -Djetty.port=8900 -DzkHost=localhost:9983 -jar start.jar
```

Notice that the parameters are exactly the same as they were for starting the second shard; you're simply pointing a new instance at the original ZooKeeper. But if you look at the SolrCloud admin page, you'll see that it was added not as a third shard, but as a replica for the first:



This is because the cluster already knew that there were only two shards and they were already accounted for, so new instances are added as replicas. Similarly, when you add the fourth instance, it's added as a replica for the second shard:

```
cd shard2replica  
java -Djetty.port=7500 -DzkHost=localhost:9983 -jar start.jar
```



If you were to add additional instances, the cluster would continue this round-robin, adding replicas as necessary.

So what does all of this mean?

It means that you now have four servers to handle your data. If you were to send data to a replica, as in:

```
curl http://localhost:7500/solr/update?commit=true -H "Content-Type: text/xml" -d  
"@money.xml"
```

the course of events goes like this:

1. Replica (in this case the server on port 7500) gets the request.
2. Replica forwards request to its leader (in this case the server on port 7574).

3. The leader processes the request, and makes sure that all of its replicas process the request as well.

In this way, the data is available via a request to any of the running instances, as you can see by requests to:

```
http://localhost:8983/solr/collection1/select?q=*&#42;
```

```
http://localhost:7574/solr/collection1/select?q=*&#42;
```

```
http://localhost:8900/solr/collection1/select?q=*&#42;
```

```
http://localhost:7500/solr/collection1/select?q=*&#42;
```

But how does this help provide high availability? Simply put, a cluster must have at least one server running for each shard in order to function. To test this, shut down the server on port 7574, and then check the other servers:

```
http://localhost:8983/solr/collection1/select?q=*&#42;
```

```
http://localhost:8900/solr/collection1/select?q=*&#42;
```

```
http://localhost:7500/solr/collection1/select?q=*&#42;
```

You should continue to see the full set of data, even though one of the servers is missing. In fact, you can have multiple servers down, and as long as at least one instance for each shard is running, the cluster will continue to function. If the leader goes down – as in this example – a new leader will be "elected" from among the remaining replicas.

Note that when we talk about servers going down, in this example it's crucial that one particular server stays up, and that's the one running on port 8983. That's because it's our overseer – the instance running ZooKeeper. If that goes down, the cluster can continue to function under some circumstances, but it won't be able to adapt to any servers that come up or go down.

That kind of single point of failure is obviously unacceptable. Fortunately, there is a solution for this problem: multiple ZooKeepers.

Using Multiple ZooKeepers in an Ensemble

 To simplify setup for this example we're using the internal ZooKeeper server that comes with Solr, but in a production environment, you will likely be using an external ZooKeeper. The concepts are the same, however. You can find instructions on setting up an external ZooKeeper server here: <http://zookeeper.apache.org/doc/r3.3.4/zookeeperStarted.html>

To truly provide high availability, we need to make sure that not only do we also have at least one shard server running at all times, but also that the cluster also has a ZooKeeper running to manage it. To do that, you can set up a cluster to use multiple ZooKeepers. This is called using a ZooKeeper ensemble.

A ZooKeeper ensemble can keep running as long as more than half of its servers are up and running, so at least two servers in a three ZooKeeper ensemble, 3 servers in a 5 server ensemble, and so on, must be running at any given time. These required servers are called a quorum.

In this example, you're going to set up the same two-shard cluster you were using before, but instead of a single ZooKeeper, you'll run a ZooKeeper server on three of the instances. Start by cleaning up any ZooKeeper data from the previous example:

```
cd <SOLR_DIST_DIR>
rm -r shard*/solr/zoo_data
```

Next you're going to restart the Solr servers, but this time, rather than having them all point to a single ZooKeeper instance, each will run ZooKeeper **and** listen to the rest of the ensemble for instructions.

You're using the same ports as before – 8983, 7574, 8900 and 7500 – so any ZooKeeper instances would run on ports 9983, 8574, 9900 and 8500. You don't actually need to run ZooKeeper on every single instance, however, so assuming you run ZooKeeper on 9983, 8574, and 9900, the ensemble would have an address of:

```
localhost:9983,localhost:8574,localhost:9900
```

This means that when you start the first instance, you'll do it like this:

```
cd shard1
java -DzkRun -DnumShards=2 -Dbootstrap_confdir=./solr/collection1/conf \
-Dcollection.configName=myconf
-DzkHost=localhost:9983,localhost:8574,localhost:9900 \
-jar start.jar
```

- Note that the order of the parameters matters. Make sure to specify the `-DzkHost` parameter after the other ZooKeeper-related parameters.

You'll notice a lot of error messages scrolling past; this is because the ensemble doesn't yet have a quorum of ZooKeepers running.

Notice also, that this step takes care of uploading the cluster's configuration information to ZooKeeper, so starting the next server is more straightforward:

```
cd shard2
java -Djetty.port=7574 -DzkRun -DnumShards=2 \
-DzkHost=localhost:9983,localhost:8574,localhost:9900 -jar start.jar
```

Once you start this instance, you should see the errors begin to disappear on both instances, as the ZooKeepers begin to update each other, even though you only have two of the three ZooKeepers in the ensemble running.

Next start the last ZooKeeper:

```
cd shard1replica
java -Djetty.port=8900 -DzkRun -DnumShards=2 \
-DzkHost=localhost:9983,localhost:8574,localhost:9900 -jar start.jar
```

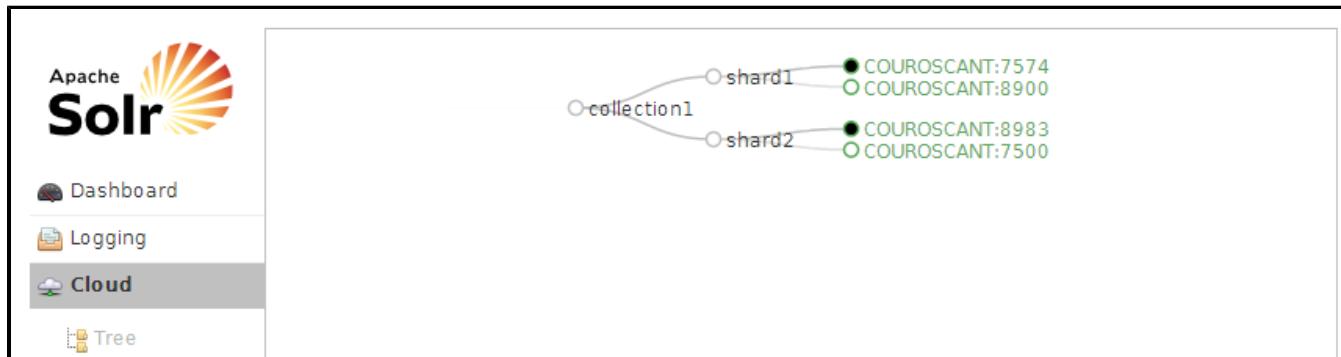
Finally, start the last replica, which doesn't itself run ZooKeeper, but references the ensemble:

```
cd shard2replica
java -Djetty.port=7500 -DzkHost=localhost:9983,localhost:8574,localhost:9900 \
-jar start.jar
```

Just to make sure everything's working properly, run a query:

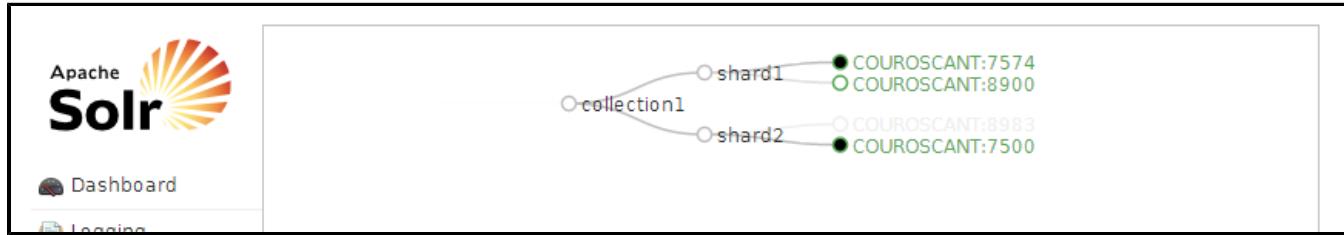
http://localhost:8983/solr/collection1/select?q=*&*

and check the SolrCloud admin page:



Now you can go ahead and kill the server on 8983, but ZooKeeper will still work, because you have more than half of the original servers still running. To verify, open the SolrCloud admin page on another server, such as:

<http://localhost:8900/solr/#/~cloud>



How SolrCloud Works

In this section, we'll discuss generally how SolrCloud works, covering these topics:

- Nodes, Cores and Clusters
- Shards and Indexing Data in SolrCloud
- Distributed Requests
- Read and Write Side Fault Tolerance
- NRT, Replication, and Disaster Recovery with SolrCloud

If you are already familiar with SolrCloud concepts and functionality, you can skip to the section covering [SolrCloud Configuration and Parameters](#).

Basic SolrCloud Concepts

On a single node, Solr has a **core** that is essentially a single **index**. If you want multiple indexes, you create multiple cores. With SolrCloud, a single index can span multiple Solr instances. This means that a single index can be made up of multiple cores on different machines.

The cores that make up one logical index are called a **collection**. A collection is essentially a single index that can span many cores, both for index scaling as well as redundancy. If, for instance, you wanted to move your two-core Solr setup to SolrCloud, you would have 2 collections, each made up of multiple individual cores.

In SolrCloud you can have multiple collections. Collections can be divided into slices. Each slice can exist in multiple copies; these copies of the same slice are called **shards**. One of the shards within a slice is the **leader**, designated by a leader-election process. Each shard is a physical index, so one shard corresponds to one core.

It is important to understand the distinction between a core and a collection. In classic single node Solr, a core is basically equivalent to a collection in that it presents one logical index. In SolrCloud, the cores on multiple nodes form a collection. This is still just one logical index, but multiple cores host different shards of the full collection. So a core encapsulates a single physical index on an instance. A collection is a combination of all of the cores that together provide a logical index that is distributed across many nodes.

Differences Between Solr 3.x-style Scaling and SolrCloud

In Solr 3.x, Solr included following features:

- The index and all changes to it are replicated to another Solr instance.
- In distributed searches, queries are sent to multiple Solr instances and the results are combined into a single output.

- Documents are available only after committing, which may be expensive and not very timely.
- Sharding must be done manually, usually through SolrJ or a similar utility, and there is no distributed indexing: your index code must understand your sharding schema.
- Replication must be manually configured and can slow down access to recent content because the system needs to wait for a commit and the replication to be triggered and to complete.
- Failure recovery may result in the loss of your ability to index, and make recovering your indexing process difficult.

With SolrCloud, some capabilities are distributed:

- SolrCloud automatically distributes index updates to the appropriate shard, distributes searches across multiple shards, and assigns replicas to shards when replicas are available.
- Near Real Time searching is supported, and if configured, documents are available after a "soft" commit.
- Indexing accesses your sharding schema automatically.
- Replication is automatic for backup purposes.
- Recovery is robust and automatic.
- ZooKeeper serves as a repository for cluster state.

Nodes, Cores and Clusters

Nodes and Cores

In SolrCloud, a node is Java Virtual Machine instance running Solr, commonly called a server. Each Solr core can also be considered a node. Any node can contain both an instance of Solr and various kinds of data.

A Solr core is basically an index of the text and fields found in documents. A single Solr instance can contain multiple "cores", which are separate from each other based on local criteria. It might be that they are going to provide different search interfaces to users (customers in the US and customers in Canada, for example), or they have security concerns (some users cannot have access to some documents), or the documents are really different and just won't mix well in the same index (a shoe database and a dvd database).

When you start a new core in SolrCloud mode, it registers itself with ZooKeeper. This involves creating an Ephemeral node that will go away if the Solr instance goes down, as well as registering information about the core and how to contact it (such as the base Solr URL, core name, etc). Smart clients and nodes in the cluster can use this information to determine who they need to talk to in order to fulfill a request.

New Solr cores may also be created and associated with a collection via [CoreAdmin](#). Additional cloud-related parameters are discussed in the [Parameter Reference](#) page. Terms used for the CREATE action are:

- **collection:** the name of the collection to which this core belongs. Default is the name of the core.
- **shard:** the shard id this core represents. (Optional: normally you want to be auto assigned a shard id.)
- **collection.<param>=<value>:** causes a property of <param>=<value> to be set if a new collection is being created. For example, use `collection.configName=<configname>` to point to the config for a new collection.

For example:

```
curl 'http://localhost:8983/solr/admin/cores?  
action=CREATE&name=mycore&collection=collection1&shard=shard2'
```

Clusters

A cluster is set of Solr nodes managed by ZooKeeper as a single unit. When you have a cluster, you can always make requests to the cluster and if the request is acknowledged, you can be sure that it will be managed as a unit and be durable, i.e., you won't lose data. Updates can be seen right after they are made and the cluster can be expanded or contracted.

Creating a Cluster

A cluster is created as soon as you have more than one Solr instance registered with ZooKeeper. The section [Getting Started with SolrCloud](#) reviews how to set up a simple cluster.

Resizing a Cluster

Clusters contain a settable number of shards. You set the number of shards for a new cluster by passing a system property, `numShards`, when you start up Solr. The `numShards` parameter must be passed on the first startup of any Solr node, and is used to auto-assign which shard each instance should be part of. Once you have started up more Solr nodes than `numShards`, the nodes will create replicas for each shard, distributing them evenly across the node, as long as they all belong to the same collection.

To add more cores to your collection, simply start the new core. You can do this at any time and the new core will sync its data with the current replicas in the shard before becoming active.

You can also avoid `numShards` and manually assign a core a shard ID if you choose.

The number of shards determines how the data in your index is broken up, so you cannot change the number of shards of the index after initially setting up the cluster.

However, you do have the option of breaking your index into multiple shards to start with, even if you are only using a single machine. You can then expand to multiple machines later. To do that, follow these steps:

1. Set up your collection by hosting multiple cores on a single physical machine (or group of machines). Each of these shards will be a leader for that shard.
2. When you're ready, you can migrate shards onto new machines by starting up a new replica for a given shard on each new machine.
3. Remove the shard from the original machine. ZooKeeper will promote the replica to the leader for that shard.

Shards and Indexing Data in SolrCloud

When your data is too large for one node, you can break it up and store it in sections by creating one or more **shards**. Each is a portion of the logical index, or core, and it's the set of all nodes containing that section of the index.

A shard is a way of splitting a core over a number of "servers" (in quotes because the server can be real servers in a rack or multiple virtual machines on a single physical box). For example, you might have a shard for data that represents each state, or different categories that are likely to be searched independently, but are often combined.

Before SolrCloud, Solr supported Distributed Search, which allowed one query to be executed across multiple shards, so the query was executed against the entire Solr index and no documents would be missed from the search results. So splitting the core across shards is not exclusively a SolrCloud concept. There were, however, several problems with the distributed approach that necessitated improvement with SolrCloud:

1. Splitting of the core into shards was somewhat manual.
2. There was no support for distributed indexing, which meant that you needed to explicitly send documents to a specific shard; Solr couldn't figure out on its own what shards to send documents to.
3. There was no load balancing or failover, so if you got a high number of queries, you needed to figure out where to send them and if one shard died it was just gone.

SolrCloud fixes all those problems. There is support for distributing both the index process and the queries automatically, and ZooKeeper provides failover and load balancing. Additionally, every shard can also have multiple replicas for additional robustness.

Unlike Solr 3.x, in SolrCloud there are no masters or slaves. Instead, there are leaders and replicas. Leaders are automatically elected, initially on a first-come-first-served basis, and then based on the Zookeeper process described at

http://zookeeper.apache.org/doc/trunk/recipes.html#sc_leaderElection..

If a leader goes down, one of its replicas is automatically elected as the new leader. As each node is started, it's assigned to the shard with the fewest replicas. When there's a tie, it's assigned to the shard with the lowest shard ID.

When a document is sent to a machine for indexing, the system first determines if the machine is a replica or a leader.

- If the machine is a replica, the document is forwarded to the leader for processing.
- If the machine is a leader, SolrCloud determines which shard the document should go to, forwards the document to the leader for that shard, indexes the document for this shard, and forwards the index notation to itself and any replicas.

Document Routing

New in Solr 4.1 is the ability to co-locate documents to improve query performance.

First, if you specify `numShards` when you create a collection, you will use the "compositeId" router by default, which will then allow you to send documents with a prefix in the document ID. The prefix will be used to calculate the hash Solr uses to determine the shard a document is sent to for indexing. The prefix can be anything you'd like it to be (it doesn't have to be the shard name, for example), but it must be consistent so Solr behaves consistently. For example, if you wanted to co-locate documents for a customer, you could use the customer name or ID as the prefix. If your customer is "IBM", for example, with a document with the ID "12345", you would insert the prefix into the document id field: "IBM!12345". The exclamation mark ('!') is critical here, as it defines the shard to direct the document to.

Then at query time, you include the prefix(es) into your query with the `shard.keys` parameter (i.e., `q=solr&shard.keys=IBM!`). In some situations, this may improve query performance because it overcomes network latency when querying all the shards.

If you do not want to influence how documents are stored, you don't need to specify a prefix in your document ID.

If you did not create the collection with the `numShards` parameter, you will be using the "implicit" router by default. In this case, you could use the `_shard_` parameter or a field to name a specific shard.

Distributed Requests

One of the advantages of using SolrCloud is the ability to distribute requests among various shards that may or may not contain the data that you're looking for. You have the option of searching over all of your data or just parts of it.

Querying all shards for a collection should look familiar; it's as though SolrCloud didn't even come into play:

```
http://localhost:8983/solr/collection1/select?q=*:*
```

If, on the other hand, you wanted to search just one shard, you can specify that shard, as in:

```
http://localhost:8983/solr/collection1/select?q=*:*&shards=localhost:7574/solr
```

If you want to search a group of shards, you can specify them together:

```
http://localhost:8983/solr/collection1/select?q=*:*&shards=localhost:7574/solr,localhost:8
```

Or you can specify a list of servers to choose from for load balancing purposes by using the pipe symbol (|):

```
http://localhost:8983/solr/collection1/select?q=*:*&shards=localhost:7574/solr|localhost:7
```

(If you have explicitly created your shards using ZooKeeper and have shard IDs, you can use those IDs rather than server addresses.)

You also have the option of searching multiple collections. For example:

```
http://localhost:8983/solr/collection1/select?collection=collection1,collection2,collectio
```

Read and Write Side Fault Tolerance

Read Side Fault Tolerance

With earlier versions of Solr, you had to set up your own load balancer. Now each individual node load balances requests across the replicas in a cluster. You still need a load balancer on the 'outside' that talks to the cluster, or you need a smart client. (Solr provides a smart Java SolrJ client called CloudSolrServer.)

A smart client understands how to read and interact with ZooKeeper and only requests the ZooKeeper ensembles' address to start discovering to which nodes it should send requests.

Write Side Fault Tolerance

SolrCloud supports near real-time actions, elasticity, high availability, and fault tolerance. What this means, basically, is that when you have a large cluster, you can always make requests to the cluster, and if a request is acknowledged you are sure it will be durable; i.e., you won't lose data. Updates can be seen right after they are made and the cluster can be expanded or contracted.

Recovery

A Transaction Log is created for each node so that every change to content or organization is noted. The log is used to determine which content in the node should be included in a replica. When a new replica is created, it refers to the Leader and the Transaction Log to know which content to include. If it fails, it retries.

Since the Transaction Log consists of a record of updates, it allows for more robust indexing because it includes redoing the uncommitted updates if indexing is interrupted.

If a leader goes down, it may have sent requests to some replicas and not others. So when a new potential leader is identified, it runs a synch process against the other replicas. If this is successful, everything should be consistent, the leader registers as active, and normal actions proceed. If the a replica is too far out of synch, the system asks for a full replication/replay-based recovery.

If an update fails because cores are reloading schemas and some have finished but others have not, the leader tells the nodes that the update failed and starts the recovery procedure.

NRT, Replication, and Disaster Recovery with SolrCloud

SolrCloud and Replication

Replication ensures redundancy for your data, and enables you to send an update request to any node in the shard. If that node is a replica, it will forward the request to the leader, which then forwards it to all existing replicas, using versioning to make sure every replica has the most up-to-date version. This architecture enables you to be certain that your data can be recovered in the event of a disaster, even if you are using Near Real Time searching.

Near Real Time Searching

If you want to use the [NearRealtimeSearch](#) support, enable auto soft commits in your `solrconfig.xml` file before storing it into Zookeeper. Otherwise you can send explicit soft commits to the cluster as you need.

SolrCloud doesn't work very well with separated data clusters connected by an expensive pipe. The root problem is that SolrCloud's architecture sends documents to all the nodes in the cluster (on a per-shard basis), and that architecture is really dictated by the NRT functionality.

Imagine that you have a set of servers in China and one in the US that are aware of each other. Assuming 5 replicas, a single update to a shard may make multiple trips over the expensive pipe before it's all done, probably slowing indexing speed unacceptably.

So the SolrCloud recommendation for this situation is to maintain these clusters separately; nodes in China don't even know that nodes exist in the US and vice-versa. When indexing, you send the update request to one node in the US and one in China and all the node-routing after that is local to the separate clusters. Requests can go to any node in either country and maintain a consistent view of the data.

However, if your US cluster goes down, you have to re-synchronize the down cluster with up-to-date information from China. The process requires you to replicate the index from China to the repaired US installation and then get everything back up and working.

Disaster Recovery for an NRT system

Use of Near Real Time (NRT) searching affects the way that systems using SolrCloud behave during disaster recovery.

The procedure outlined below assumes that you are maintaining separate clusters, as described above. Consider, for example, an event in which the US cluster goes down (say, because of a hurricane), but the China cluster is intact. Disaster recovery consists of creating the new system and letting the intact cluster create a replicate for each shard on it, then promoting those replicas to be leaders of the newly created US cluster.

Here are the steps to take:

1. Take the downed system offline to all end users.
2. Take the indexing process offline.
3. Repair the system.
4. Bring up one machine per shard in the repaired system as part of the ZooKeeper cluster on the good system, and wait for replication to happen, creating a replica on that machine.
(SoftCommits will not be repeated, but data will be pulled from the transaction logs if necessary.)

 SolrCloud will automatically use old-style replication for the bulk load. By temporarily having only one replica, you'll minimize data transfer across a slow connection.

5. Bring the machines of the repaired cluster down, and reconfigure them to be a separate Zookeeper cluster again, optionally adding more replicas for each shard.
6. Make the repaired system visible to end users again.
7. Start the indexing program again, delivering updates to both systems.

SolrCloud Configuration and Parameters

In this section, we'll cover the various configuration options for SolrCloud.

In general, with a new Solr 4 instance, the required configuration is in the sample `schema.xml` and `solrconfig.xml` files. However, there may be reasons to change default settings or configure the cloud elements manually.

The following topics are covered in these pages:

- [Using ZooKeeper to Manage Configuration Files](#)
- [Managing Collections via the Collections API](#)
- [Parameter Reference](#)
- [Command Line Utilities](#)
- [SolrCloud with Legacy Configuration Files](#)

Using ZooKeeper to Manage Configuration Files

With SolrCloud your configuration files (particularly `solrconfig.xml` and `schema.xml`) are kept in ZooKeeper. These files are uploaded when you first start Solr in SolrCloud mode.

Startup Bootstrap Parameters

There are two different ways you can use system properties to upload your initial configuration files to ZooKeeper the first time you start Solr. Remember that these are meant to be used only on first startup or when overwriting configuration files. Every time you start Solr with these system properties, any current configuration files in ZooKeeper may be overwritten when `conf.set` names match.

The first way is to look at `solr.xml` and upload the `conf` for each core found. The `config.set` name will be the collection name for that core, and collections will use the `config.set` that has a matching name. One parameter is used with this approach, `bootstrap_conf`. If you pass `-Dbootstrap_conf=true` on startup, each core you have configured will have its configuration files automatically uploaded and linked to the collection containing the core.

An alternate approach is to upload the given directory as a `config.set` with the given name. No linking of collection to `config.set` is done. However, if only one `conf.set` exists, a collection will autolink to it. Two parameters are used with this approach:

Parameter	Default value	Description
<code>bootstrap_confdir</code>	No default	If you pass <code>-bootstrap_confdir=<directory></code> on startup, that specific directory of configuration files will be uploaded to ZooKeeper with a <code>conf.set</code> name defined by the system property below, <code>collection.configName</code> .
<code>collection.configName</code>	Defaults to <code>configuration1</code>	Determines the name of the <code>conf.set</code> pointed to by <code>bootstrap_confdir</code> .

Using the [ZooKeeper Command Line Interface \(zkCLI\)](#), you can download and re-upload these configuration files.



It's important to keep these files under version control.

Managing Your SolrCloud Configuration Files

To update or change your SolrCloud configuration files:

1. Download the lastest configuration files from ZooKeeper, using the source control checkout process.
2. Make your changes.
3. Commit your changed file to source control.
4. Push the changes back to ZooKeeper.
5. Reload the collection so that the changes will be in effect.

There are some scripts available with the ZooKeeper Command Line Utility to help manage changes to configuration files, discussed in the section on [Command Line Utilities](#).

Managing Collections via the Collections API

The Collections API is generally used to enable you to create, remove, or reload collections, but in the context of SolrCloud you can also use it to create collections with a specific number of shards and replicas. For example, to create a new collection called `mycollection` with 5 shards and 10 replicas, you could use the command:

```
http://localhost:8983/solr/admin/collections?action=CREATE&name=mycollection&numShards=5&r
```

To delete that collection, you could use the command:

```
http://localhost:8983/solr/admin/collections?action=DELETE&name=mycollection
```

Or to reload it, you could use:

```
http://localhost:8983/solr/admin/collections?action=RELOAD&name=mycollection
```

There is nothing magical about these commands; they simply aggregate all of the manual steps you would normally use to manage your collections.

Collection API Parameters

The collection API has several possible parameters.

Parameter	Description
action	The CoreAdmin operation to be performed across all cores of a collection. It can currently support: <ul style="list-style-type: none">• CREATE• RELOAD• DELETE
name	The name of the collection. In the case of CREATE action, this would be the new name for the collection. This is required for all API calls.
numShards	The number of shards to be created as part of the collection (used when creating a collection).

replicationFactor	<p>The number of replicas to be created for each shard.</p> <p>Info In Solr 4.1, the meaning of <code>replicationFactor</code> has changed. In Solr 4.1, it means the total number of replicas; in Solr 4.0, it referred to the number of <i>additional</i> copies.</p>
maxShardsPerNode	<p>When creating collections, the shards and/or replicas are spread across all available (i.e., live) nodes, and two replicas of the same shard will never be on the same node. If a node is not live when the CREATE operation is called, it will not get any parts of the new collection, which could lead to too many replicas being created on a single (live) node. Defining the <code>maxShardsPerNode</code> sets a limit on the number of replicas CREATE will spread to each node. The default is 1. If the entire collection can not be fit into the live nodes, no collection will be created at all.</p>
createNodeSet	<p>Allows defining the nodes to spread the new collection across. If not provided, the CREATE operation will create shard-replica spread across all of your live Solr nodes. The format is a comma-separated list of <code>node_names</code>, such as <code>localhost:8983_solr,localhost:8984_solr,localhost:8985_solr</code>.</p>
collection.configName	

Parameter Reference

Cluster Parameters

numShards	Defaults to 1	The number of shards to hash documents to. There must be one leader per shard and each leader can have N replicas.	
-----------	---------------	--	--

SolrCloud Instance Parameters

These are set in `solr.xml`, but by default they are set up to also work with system properties.

host	Defaults to the first local host address found	If the wrong host address is found automatically, you can override the host address with this parameter.	
hostPort	Defaults to the jetty.port system property	The port that Solr is running on. By default this is found by looking at the <code>jetty.port</code> system property.	
hostContext	Defaults to solr	The context path for the Solr web application.	

SolrCloud Instance ZooKeeper Parameters

zkRun	Defaults to <code>localhost:<solrPort+1001></code>	Causes Solr to run an embedded version of ZooKeeper. Set to the address of ZooKeeper on this node; this allows us to know who you are in the list of addresses in the <code>zkHost</code> connect string. Use <code>-DzkRun</code> to get the default value.
zkHost	No default	The host address for ZooKeeper. Usually this is a comma-separated list of addresses to each node in your ZooKeeper ensemble.
zkClientTimeout	Defaults to 15000	The time a client is allowed to not talk to ZooKeeper before its session expires.

`zkRun` and `zkHost` are set up using system properties. `zkClientTimeout` is set up in `solr.xml` by default, but can also be set using a system property.

SolrCloud Core Parameters

shardId	Defaults to being automatically assigned based on numShards	Allows you to specify the id used to group cores into shards.
---------	---	---

shardId can be configured in `solr.xml` for each core element as an attribute.

Additional cloud related parameters are discussed in [Core Admin and Configuring solr.xml](#).

Command Line Utilities

Solr's Administration page (found by default at <http://hostname:8983/solr/admin>), provides a section with menu items for monitoring indexing and performance statistics, information about index distribution and replication, and information on all threads running in the JVM at the time. There is also a section where you can run queries, and an assistance area.

In addition, SolrCloud provides its own administration page (found by default at <http://localhost:8983/solr/#/~cloud>), as well as a few tools available via ZooKeeper's Command Line Utility (CLI). The CLI lets you upload configuration information to ZooKeeper, in the same two ways that were shown in the examples in [Parameter Reference](#). It also provides a few other commands that let you link collection sets to collections, make ZooKeeper paths or clear them, and download configurations from ZooKeeper to the local filesystem.

Using The ZooKeeper CLI

ZooKeeper has a utility that lets you pass command line parameters: `zkcli.bat` (for Windows environments) and `zkcli.sh` (for Unix environments).

zkcli Parameters

The `zkcli` command takes the form `-cmd <arg>` where the command to run can be `bootstrap`, `upconfig`, `downconfig`, `linkconfig`, `makepath`, or `clear`.

Parameter	Usage	Meaning
<code>-c</code>	<code>-collection <arg></code>	For <code>linkconfig</code> : name of the collection.
<code>-d</code>	<code>--confdir <arg></code>	For <code>upconfig</code> : a directory of configuration files.
<code>-h</code>	<code>--help</code>	Bring up the help page.
<code>-n</code>	<code>--confname <arg></code>	For <code>upconfig</code> , <code>linkconfig</code> : name of the configuration set.
<code>-r</code>	<code>--runzk <arg></code>	Run ZooKeeper internally by passing the Solr run port; only for clusters on one machine.
<code>-s</code>	<code>--solrhome <arg></code>	For <code>bootstrap</code> , <code>runzk</code> : <code>solrhome <location></code> .
<code>-z</code>	<code>--zkhost <arg></code>	ZooKeeper host address. This parameter is mandatory for all CLI commands.

ZooKeeper CLI Examples

Below are some examples of using the zkcli CLI:

Uploading a Configuration Directory

```
java -classpath example/solr-webapp/WEB-INF/lib/*  
      org.apache.solr.cloud.ZkCLI -cmd upconfig -zkhost 127.0.0.1:9983  
      -confdir example/solr/collection1/conf -confname conf1 -solrhome example/solr
```

Linking a Collection to a Configuration Set

```
java -classpath example/solr-webapp/webapp/WEB-INF/lib/*  
      org.apache.solr.cloud.ZkCLI -cmd linkconfig -zkhost 127.0.0.1:9983  
      -collection collection1 -confname conf1 -solrhome example/solr
```

Bootstrapping All the Configuration Directories in solr.xml

```
java -classpath example/solr-webapp/webapp/WEB-INF/lib/*  
      org.apache.solr.cloud.ZkCLI -cmd bootstrap -zkhost 127.0.0.1:9983  
      -solrhome example/solr
```

Scripts

There are scripts in example/cloud-scripts that handle the classpath and class name for you if you are using Solr out of the box with Jetty. Commands then become:

```
sh zkcli.sh -cmd linkconfig -zkhost 127.0.0.1:9983  
      -collection collection1 -confname conf1 -solrhome example/solr
```

SolrCloud with Legacy Configuration Files

All of the required configuration is already set up in the sample configurations shipped with Solr. You only need to add the following if you are migrating old configuration files. Do not remove these files and parameters from a new Solr instance if you intend to use Solr in SolrCloud mode.

These properties exist in 3 files: schema.xml, solrconfig.xml, and solr.xml.

1. In schema.xml, you must have a _version_ field defined:

```
<field name="_version_" type="long" indexed="true" stored="true" multiValued="false"/>
```

2. In solrconfig.xml, you must have an UpdateLog defined. This should be defined in the updateHandler section.

```
<updateHandler>
...
<updateLog>
    <str name="dir">${solr.data.dir:}</str>
</updateLog>
...
</updateHandler>
```

3. You must have a replication handler called /replication defined:

```
<requestHandler name="/replication" startup="lazy" />
```

There are several parameters available for this handler, discussed in the section [Index Replication](#).

4. You must have a Realtime Get handler called "/get" defined:

```
<requestHandler name="/get">
    <lst name="defaults">
        <str name="omitHeader">true</str>
    </lst>
</requestHandler>
```

The parameters for this handler are discussed in the section [RealTime Get](#).

5. You must have the admin handlers defined:

```
<requestHandler name="/admin/" class="solr.admin.AdminHandlers" />
```

6. And you must leave the admin path in `solr.xml` as the default:

```
<cores adminPath="/admin/cores" />
```

7. The [DistributedUpdateProcessor](#) is part of the default update chain and is automatically injected into any of your custom update chains, so you don't actually need to make any changes for this capability. However, should you wish to add it explicitly, you can still add it to the `solrconfig.xml` file as part of an `updateRequestProcessorChain`. For example:

```
<updateRequestProcessorChain name="sample">
  <processor class="solr.LogUpdateProcessorFactory" />
  <processor class="solr.DistributedUpdateProcessorFactory"/>
  <processor class="my.package.UpdateFactory" />
  <processor class="solr.RunUpdateProcessorFactory" />
</updateRequestProcessorChain>
```

If you do not want the [DistributedUpdateProcessFactory](#) auto-injected into your chain (for example, if you want to use SolrCloud functionality, but you want to distribute updates yourself) then specify the `NoOpDistributingUpdateProcessorFactory` update processor factory in your chain:

```
<updateRequestProcessorChain name="sample">
  <processor class="solr.LogUpdateProcessorFactory" />
  <processor class="solr.NoOpDistributingUpdateProcessorFactory" />
  <processor class="my.package.MyDistributedUpdateFactory" />
  <processor class="solr.RunUpdateProcessorFactory" />
</updateRequestProcessorChain>
```

In the update process, Solr skips updating processors that have already been run on other nodes.

Setting Up an External ZooKeeper Ensemble

Although Solr comes bundled with Apache ZooKeeper, you should consider yourself discouraged from using this internal ZooKeeper in production, because shutting down a redundant Solr instance will also shut down its ZooKeeper server, which might not be quite so redundant. Because a ZooKeeper ensemble must have a quorum of more than half its servers running at any given time, this can be a problem.

The solution to this problem is to set up an external ZooKeeper ensemble. Fortunately, while this process can seem intimidating due to the number of powerful options, setting up a simple ensemble is actually quite straightforward. The basic steps are as follows:

Download Apache ZooKeeper

The first step in setting up Apache ZooKeeper is, of course, to download the software. It's available from <http://zookeeper.apache.org/releases.html>.

- ⚠ When using stand-alone ZooKeeper, you need to take care to keep your version of ZooKeeper updated with the latest version distributed with Solr. Since you are using it as a stand-alone application, it does not get upgraded when you upgrade Solr.
 - Solr 4.0 uses Apache ZooKeeper v3.3.6.
 - Solr 4.1 uses Apache ZooKeeper v3.4.5.

Create the instance

Creating the instance is a simple matter of extracting the files into a specific target directory. The actual directory itself doesn't matter, as long as you know where it is, and where you'd like to have ZooKeeper store its internal data.

Configure the instance

The next step is to configure your ZooKeeper instance. To do that, create the following file:

```
<ZOOKEEPER_HOME>/conf/zoo.cfg
```

and add the following information:

```
tickTime=2000  
dataDir=/var/lib/zookeeper  
clientPort=2181
```

The parameters are as follows:

tickTime: Part of what ZooKeeper does is to determine which servers are up and running at any given time, and the minimum session time out is defined as two "ticks". The `tickTime` parameter specifies, in milliseconds, how long each tick should be.

dataDir: This is the directory in which ZooKeeper will store data about the cluster. This directory should start out empty.

clientPort: This is the port on which Solr will access ZooKeeper.

Once this file is in place, you're ready to start the ZooKeeper instance.

Run the instance

To run the instance, you can simply use the script provided:

```
bin/zkServer.sh start
```

Again, ZooKeeper provides a great deal of power through additional configurations, but delving into them is beyond the scope of this tutorial. For more information, see the ZooKeeper [Getting Started](#) page. For this example, however, the defaults are fine.

Point Solr at the instance

Pointing Solr at the ZooKeeper instance you've created is a simple matter of using the `-DzkHost` parameter. For example, in the [Getting Started](#) example you learned how to point to the internal ZooKeeper. In this example, you would point to the ZooKeeper you've started on port 2181:

```
cd shard1
java -DnumShards=2 -Dbootstrap_confdir=./solr/collection1/conf \
      -Dcollection.configName=myconf -DzkHost=localhost:2181 -jar start.jar
cd shard2java
java -Djetty.port=7574 -DzkHost=localhost:2181 -jar start.jar
```

As before, you must first upload the configuration information, and then you can connect a second instance.

Shut down ZooKeeper

To shut down ZooKeeper, use the `zkServer` script:

```
bin/zkServer.sh stop
```

Setting up a ZooKeeper ensemble

In the Getting Started example, using a ZooKeeper ensemble was a simple matter of starting multiple instances and pointing to them. With an external ZooKeeper ensemble, you need to set things up just a little more carefully.

The difference is that rather than simply starting up the servers, you need to configure them to know about and talk to each other first. So your original `zoo.cfg` file might look like this:

```
dataDir=/var/lib/zookeeperdata/1
clientPort=2181
initLimit=5
syncLimit=2
server.1=localhost:2888:3888
server.2=localhost:2889:3889
server.3=localhost:2890:3890
```

Here you see three new parameters:

initLimit: The time, in ticks, the server allows for connecting to the leader. In this case, you have 5 ticks, each of which is 2000 milliseconds long, so the server will wait as long as 10 seconds to connect.

syncLimit: The time, in ticks, the server will wait before updating itself from the leader.

server.X: These are the locations of all servers in the ensemble, and the ports on which they communicate with each other. The server id, stored in the `<dataDir>/myid` file, identifies each server, so in the case of this first instance, you would create the file `/var/lib/zookeeperdata/1/myid` with the content

```
1
```

Now, whereas with Solr you need to create entirely new directories to run multiple instances, all you need for a new ZooKeeper instance, even if it's on the same machine for testing purposes, is a new configuration file. To complete the example you'll create two more configuration files.

The `<ZOOKEEPER_HOME>/conf/zoo2.cfg` file should have the content:

```
tickTime=2000
dataDir=c:/sw/zookeeperdata/2
clientPort=2182
initLimit=5
syncLimit=2
server.1=localhost:2888:3888
server.2=localhost:2889:3889
server.3=localhost:2890:3890
```

You'll also need to create <ZOOKEEPER_HOME>/conf/zoo3.cfg:

```
tickTime=2000
dataDir=c:/sw/zookeeperdata/3
clientPort=2183
initLimit=5
syncLimit=2
server.1=localhost:2888:3888
server.2=localhost:2889:3889
server.3=localhost:2890:3890
```

Finally, create your `myid` files in each of the `dataDir` directories so that each server knows which instance it is.

To start the servers, you can simply explicitly reference the configuration files:

```
cd <ZOOKEEPER_HOME>
bin/zkServer.sh start zoo.cfg
bin/zkServer.sh start zoo2.cfg
bin/zkServer.sh start zoo3.cfg
```

Once these servers are running, you can reference them from Solr just as you did before:

```
java -DnumShards=2 -Dbootstrap_confdir=./solr/collection1/conf \
-Dcollection.configName=myconf
-DzkHost=localhost:2181,localhost:2182,localhost:2183 -jar start.jar
```

For more information on getting the most power from your ZooKeeper installation, check out the [ZooKeeper Administrator's Guide](#).

SolrCloud Glossary

SolrCloud Glossary

- **Node** : A JVM instance running Solr; a server.
- **Cluster** : A cluster is a set of Solr nodes managed as a single unit.
- **Core** : An individual Solr instance (represents a logical index). Multiple cores can run on a single node.
- **Shard** : In Solr, a logical section of a single collection. This may be spread across multiple nodes of the cluster. Each shard can have as many replicas as needed.
- **Leader** : Each shard has one node identified as its leader. All the writes for documents belonging to a partition are routed through the leader.
- **Collection (Solr)** : Multiple documents that make up one logical index. A collection must have a single schema, but can be spread over multiple cores. A cluster can have multiple collections. When you create a collection, you specify its number of shards. You may have many collections in the same set of servers, each with a different number of shards.
- **ZooKeeper** : Apache ZooKeeper keeps track of configuration and naming, among other things, for a cluster of Solr nodes. A ZooKeeper cluster is used as (1) the central configuration store for the cluster, (2) a coordinator for operations requiring distributed synchronization, and (3) the system of record for cluster topology.
- **Ensemble** : Multiple ZooKeeper instances running simultaneously.
- **Collection (ZooKeeper)** : A group of cores managed together as part of a SolrCloud installation.
- **Overseer** : The Overseer coordinates the clusters. It keeps track of the existing nodes and shards and assigns shards to nodes.
- **Transaction Log** : An append-only log of write operations maintained by each node.
- **Document** : A group of fields and their values. Documents are the basic unit of storage, and their specific locations are found using an index. Documents are assigned to shards using standard hashing, or by specifically assigning a shard within the document ID. Documents are versioned after each write operation.

Related Pages

- [Solr Glossary](#)

Legacy Scaling and Distribution

This section describes how to set up distribution and replication in Solr. It is considered "legacy" behavior, since while it is still supported in Solr, the SolrCloud functionality described in the previous chapter is where the current development is headed. However, if you don't need all that SolrCloud delivers, search distribution and index replication may be sufficient.

This section covers the following topics:

[Introduction to Scaling and Distribution](#): Conceptual information about distribution and replication in Solr.

[Distributed Search with Index Sharding](#): Detailed information about implementing distributed searching in Solr.

[Index Replication](#): Detailed information about replicating your Solr indexes.

[Combining Distribution and Replication](#): Detailed information about replicating shards in a distributed index.

[Merging Indexes](#): Information about combining separate indexes in Solr.

Introduction to Scaling and Distribution

Both Lucene and Solr were designed to scale to support large implementations with minimal custom coding. This section covers:

- [distributing](#) an index across multiple servers
- [replicating](#) an index on multiple servers
- [merging indexes](#)

If you need full scale distribution of indexes and queries, as well as replication, load balancing and failover, you may want to use SolrCloud. Full details on configuring and using SolrCloud is available in the section [SolrCloud](#).

What Problem Does Distribution Solve?

If searches are taking too long or the index is approaching the physical limitations of its machine, you should consider distributing the index across two or more Solr servers.

To distribute an index, you divide the index into partitions called shards, each of which runs on a separate machine. Solr then partitions searches into sub-searches, which run on the individual shards, reporting results collectively. The architectural details underlying index sharding are invisible to end users, who simply experience faster performance on queries against very large indexes.

What Problem Does Replication Solve?

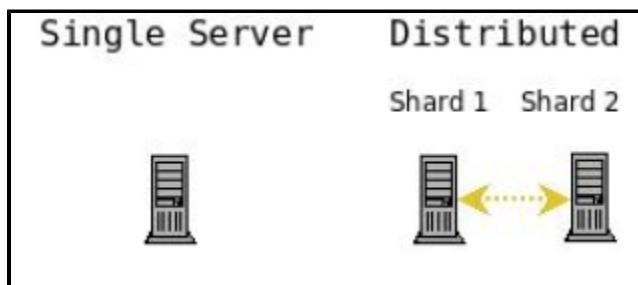
Replicating an index is useful when:

- You have a large search volume which one machine cannot handle, so you need to distribute searches across multiple read-only copies of the index.
- There is a high volume/high rate of indexing which consumes machine resources and reduces search performance on the indexing machine, so you need to separate indexing and searching.
- You want to make a backup of the index (see [Backing Up](#)).

Distributed Search with Index Sharding

When an index becomes too large to fit on a single system, or when a query takes too long to execute, an index can be split into multiple shards, and Solr can query and merge results across those shards. A single shard receives the query, distributes the query to other shards, and integrates the results. You can find additional information about distributed search on the Solr wiki: <http://wiki.apache.org/solr/DistributedSearch>.

The figure below compares a single server to a distributed configuration with two shards.



- ⓘ If single queries are currently fast enough and if one simply wants to expand the capacity (queries/sec) of the search system, then standard index replication (replicating the entire index on multiple servers) should be used instead of index sharding.

Update commands may be sent to any server with distributed indexing configured correctly. Document adds and deletes are forwarded to the appropriate server/shard based on a hash of the unique document id. **commit** commands and **deleteByQuery** commands are sent to every server in shards.

Update reorders (i.e., replica A may see update X then Y, and replica B may see update Y then X). **deleteByQuery** also handles reorders the same way, to ensure replicas are consistent. All replicas of a shard are consistent, even if the updates arrive in a different order on different replicas.

Distributed Support for Date and Numeric Range Faceting

You can now use range facetting for anything that uses date math (both date and numeric ranges). In addition, you can now use NOW by including `FacetParams.FACET_DATE_NOW` in the original request to sync remote shard requests to a common 'now' time. For example, using range facetting is a convenient way to keep the rest of your request the same, but check how the current date affects your date boosting strategies.

FacetParams.FACET_DATE_NOW takes as a parameter a (stringified) long that is the number of milliseconds from 1 Jan 1970 00:00, i.e., the returned value from a System.currentTimeMillis() call. This delineates it from a 'searchable' time and avoids superfluous date parsing.

 NOTE: This parameter affects date facet timing only. If there are other areas of a query that rely on 'NOW', these will not interpret this value.

For distributed facet_dates, Solr steps through each date facet, adding and merging results from the current shard.

Any time and/or time zone differences are NOT taken into account here. The issue of time zone/skew on distributed shards is currently handled by passing a facet.date.now=<epochtime> parameter in the search query. This is then used by the participating shards to use as 'now'.

If you use the first encountered shard's facet_dates as the basis for subsequent shards' data to be merged in, if subsequent shards' facet_dates are skewed in relation to the first by a >1 'gap', these 'earlier' or 'later' facets will not be merged in.

There are two reasons for this:

1. Performance: It is of course faster to check facet_date lists against a single map's data, rather than against each other.
2. If 'earlier' and/or 'later' facet_dates are added in, this makes the time range larger than that which was requested (e.g. a request for one hour's worth of facets could bring back 2, 3, or more hours of data).

Distributing Documents across Shards

It is up to you to get all your documents indexed on each shard of your server farm. Solr does not include out-of-the-box support for distributed indexing, but your method can be as simple as a round robin technique. Just index each document to the next server in the circle. (For more information about indexing, see [Indexing and Basic Data Operations](#).)

A simple hashing system would also work. The following should serve as an adequate hashing function.

```
uniqueId.hashCode() % numServers
```

One advantage of this approach is that it is easy to know where a document is if you need to update it or delete. In contrast, if you are moving documents around in a round-robin fashion, you may not know where a document actually is.

Solr does not calculate universal term/doc frequencies. For most large-scale implementations, it is not likely to matter that Solr calculates TD/IDF at the shard level. However, if your collection is heavily skewed in its distribution across servers, you may find misleading relevancy results in your searches. In general, it is probably best to randomly distribute documents to your shards.

You can directly configure aspects of the concurrency and thread-pooling used within distributed search in Solr. This allows for finer grained control and you can tune it to target your own specific requirements. The default configuration favors throughput over latency.

To configure the standard handler, provide a configuration like this:

```
<requestHandler name="standard" class="solr.SearchHandler" default="true">
    <!-- other params go here -->
    <shardHandlerFactory class="HttpShardHandlerFactory">
        <int name="socketTimeOut">1000</int>
        <int name="connTimeOut">5000</int>
    </shardHandler>
</requestHandler>
```

The parameters that can be specified are as follows:

Parameter	Default	Explanation
socketTimeout	0 (use OS default)	The amount of time in ms that a socket is allowed to wait.
connTimeout	0 (use OS default)	The amount of time in ms that is accepted for binding / connecting a socket
maxConnectionsPerHost	20	The maximum number of connections that is made to each individual shard in a distributed search.
corePoolSize	0	The retained lowest limit on the number of threads used in coordinating distributed search.
maximumPoolSize	Integer.MAX_VALUE	The maximum number of threads used for coordinating distributed search.
maxThreadIdleTime	5 seconds	The amount of time to wait for before threads are scaled back in response to a reduction in load.

sizeOfQueue	-1	If specified, the thread pool will use a backing queue instead of a direct handoff buffer. High throughput systems will want to configure this to be a direct hand off (with -1). Systems that desire better latency will want to configure a reasonable size of queue to handle variations in requests.
fairnessPolicy	false	Chooses the JVM specifics dealing with fair policy queuing, if enabled distributed searches will be handled in a First in First out fashion at a cost to throughput. If disabled throughput will be favored over latency.

Executing Distributed Searches with the shards Parameter

If a query request includes the `shards` parameter, the Solr server distributes the request across all the shards listed as arguments to the parameter. The `shards` parameter uses this syntax:

```
host:port/base_url[,host:port/base_url]*
```

For example, the `shards` parameter below causes the search to be distributed across two Solr servers: **solr1** and **solr2**, both of which are running on port 8983:

```
http://localhost:8983/solr/select?
```

```
shards=solr1:8983/solr,solr2:8983/solr&indent=true&q=ipod+solr
```

Rather than require users to include the `shards` parameter explicitly, it is usually preferred to configure this parameter as a default in the `RequestHandler` section of `solrconfig.xml`.

 Do not add the `shards` parameter to the standard `requestHandler`; otherwise, search queries may enter an infinite loop. Instead, define a new `requestHandler` that uses the `shards` parameter, and pass distributed search requests to that handler.

Currently, only query requests are distributed. This includes requests to the standard `requestHandler` (and subclasses such as the `DisMax RequestHandler`), and any other handler (`org.apache.solr.handler.component.SearchHandler`) using standard components that support distributed search.

Where `shards.info=true`, distributed responses will include information about the shard (where each shard represents a logically different index or physical location), such as the following:

```
<lst name="shards.info">
  <lst name="localhost:7777/solr">
    <long name="numFound">1333</long>
    <float name="maxScore">1.0</float>
    <long name="time">686</long>
  </lst>
  <lst name="localhost:8888/solr">
    <long name="numFound">342</long>
    <float name="maxScore">1.0</float>
    <long name="time">602</long>
  </lst>
</lst>
```

The following components support distributed search:

- The **Query** component, which returns documents matching a query
- The **Facet** component, which processes facet.query and facet.field requests where facets are sorted by count (the default).
- The **Highlighting** component, which enables Solr to include "highlighted" matches in field values.
- The **Stats** component, which returns simple statistics for numeric fields within the DocSet.
- The **Debug** component, which helps with debugging.

Limitations to Distributed Search

Distributed searching in Solr has the following limitations:

- Each document indexed must have a unique key.
- If Solr discovers duplicate document IDs, Solr selects the first document and discards subsequent ones.
- Inverse-document frequency (IDF) calculations cannot be distributed.
- Distributed searching does not support the QueryElevationComponent, which configures the top results for a given query regardless of Lucene's scoring. For more information, see <http://wiki.apache.org/solr/QueryElevationComponent>.
- The index for distributed searching may become momentarily out of sync if a commit happens between the first and second phase of the distributed search. This might cause a situation where a document that once matched a query and was subsequently changed may no longer match the query but will still be retrieved. This situation is expected to be quite rare, however, and is only possible for a single query request.
- Distributed searching supports only sorted-field faceting, not date facetting
- The number of shards is limited by number of characters allowed for GET method's URI; most Web servers generally support at least 4000 characters, but many servers limit URI length to reduce their vulnerability to Denial of Service (DoS) attacks.

- TF/IDF computations are per shard. This may not matter if content is well (randomly) distributed.
- Shard information can be returned with each document in a distributed search by including `fl=id, [shard]` in the search request. This returns the shard URL.
- In a distributed search, the data directory from the core descriptor overrides any data directory in `solrconfig.xml`.
- Update commands may be sent to any server with distributed indexing configured correctly. Document adds and deletes are forwarded to the appropriate server/shard based on a hash of the unique document id. **commit** commands and **deleteByQuery** commands are sent to every server in shards.

Avoiding Distributed Deadlock

Each shard may also serve top-level query requests and then make sub-requests to all of the other shards. In this configuration, care should be taken to ensure that the max number of threads serving HTTP requests in the servlet container is greater than the possible number of requests from both top-level clients and other shards. If this is not the case, the configuration may result in a distributed deadlock.

For example, a deadlock might occur in the case of two shards, each with just a single thread to service HTTP requests. Both threads could receive a top-level request concurrently, and make sub-requests to each other. Because there are no more remaining threads to service requests, the servlet containers will block the incoming requests until the other pending requests are finished, but they will not finish since they are waiting for the sub-requests. By ensuring that the servlets are configured to handle a sufficient number of threads, you can avoid deadlock situations like this.

Testing Index Sharding on Two Local Servers

For simple functionality testing, it's easiest to just set up two local Solr servers on different ports. (In a production environment, of course, these servers would be deployed on separate machines.)

1. Make a copy of the solr example directory:

```
cd solr  
cp -r example example7574
```

2. Change the port number:

```
perl -pi -e s/8983/7574/g example7574/etc/jetty.xml  
example7574/exampledocs/post.sh
```

3. In the first window, start up the server on port 8983:

```
cd examplejava -server -jar start.jar
```

4. In the second window, start up the server on port 7574:

```
cd example7574/java -server -jar start.jar
```

5. In the third window, index some example documents to each server:

```
cd example/exampledocs./post.sh \[a-m\]\*.xmlcd  
.../example7574/exampledocs./post.sh \[n-z\]\*.xml
```

6. Now do a distributed search across both servers with your browser or curl:

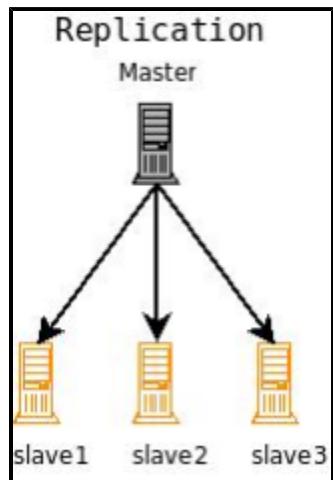
```
curl  
'http://localhost:8983/solr/select?shards=localhost:8983/solr,localhost:7574/solr&inc
```

Index Replication

⚠ The Lucene index format has changed with Solr 4. As a result, once you upgrade, previous versions of Solr will no longer be able to read the rest of your indices. In a master/slave configuration, all searchers/slaves should be upgraded before the master. If the master is updated first, the older searchers will not be able to read the new index format.

Index Replication distributes complete copies of a master index to one or more slave servers. The master server continues to manage updates to the index. All querying is handled by the slaves. This division of labor enables Solr to scale to provide adequate responsiveness to queries against large search volumes.

The figure below shows a Solr configuration using index replication. The master server's index is replicated on the slaves.



A Solr index can be replicated across multiple slave servers, which then process requests.

Topics covered in this section:

- [Index Replication in Solr](#)
- [Replication Terminology](#)
- [Configuring the Replication RequestHandler on a Master Server](#)
- [Configuring the Replication RequestHandler on a Slave Server](#)
- [Setting Up a Repeater with the ReplicationHandler](#)
- [Commit and Optimize Operations](#)
- [Slave Replication](#)
- [Index Replication using ssh and rsync](#)
- [The Snapshot and Distribution Process](#)
- [Commit and Optimization](#)
- [Distribution and Optimization](#)

Index Replication in Solr

Solr includes a Java implementation of index replication that works over HTTP.

For information on the ssh/rsync based replication, see [Index Replication using ssh and rsync](#).

The Java-based implementation of index replication offers these benefits:

- Replication without requiring external scripts
- The configuration affecting replication is controlled by a single file, `solrconfig.xml`
- Supports the replication of configuration files as well as index files
- Works across platforms with same configuration
- No reliance on OS-dependent hard links
- Tightly integrated with Solr; an admin page offers fine-grained control of each aspect of replication
- The Java-based replication feature is implemented as a RequestHandler. Configuring replication is therefore similar to any normal RequestHandler.

Replication Terminology

The table below defines the key terms associated with Solr replication.

Term	Definition
Collection	A Lucene collection is a directory of files. These files make up the indexed and returnable data of a Solr search repository.
Distribution	The copying of a collection from the master server to all slaves. The distribution process takes advantage of Lucene's index file structure.
Inserts and Deletes	As inserts and deletes occur in the collection, the directory remains unchanged. Documents are always inserted into newly created files. Documents that are deleted are not removed from the files. They are flagged in the file, deletable, and are not removed from the files until the collection is optimized.
Master and Slave	The Solr distribution model uses the master/slave model. The master is the service which receives all updates initially and keeps everything organized. Solr uses a single update master server coupled with multiple query slave servers. All changes (such as inserts, updates, deletes, etc.) are made against the single master server. Changes made on the master are distributed to all the slave servers which service all query requests from the clients.
Update	An update is a single change request against a single Solr instance. It may be a request to delete a document, add a new document, change a document, delete all documents matching a query, etc. Updates are handled synchronously within an individual Solr instance.

Optimization	A process that compacts the index and merges segments in order to improve query performance. New secondary segment(s) are created to contain documents inserted into the collection after it has been optimized. A Lucene collection must be optimized periodically to maintain satisfactory query performance. Optimization is run on the master server only. An optimized index will give you a performance gain at query time of at least 10%. This gain may be more on an index that has become fragmented over a period of time with many updates and no optimizations. Optimizations require a much longer time than does the distribution of an optimized collection to all slaves.
Segments	The number of files in a collection.
mergeFactor	A parameter that controls the number of files (segments) in a collection. For example, when mergeFactor is set to 3, Solr will fill one segment with documents until the limit maxBufferedDocs is met, then it will start a new segment. When the number of segments specified by mergeFactor is reached--in this example, 3--then Solr will merge all the segments into a single index file, then begin writing new documents to a new segment.
Snapshot	A directory containing hard links to the data files. Snapshots are distributed from the master server when the slaves pull them, "smartcopying" the snapshot directory that contains the hard links to the most recent collection data files.

Configuring the Replication RequestHandler on a Master Server

Before running a replication, you should set the following parameters on initialization of the handler:

Name	Description
replicateAfter	String specifying action after which replication should occur. Valid values are commit, optimize, or startup. There can be multiple values for this parameter. If you use "startup", you need to have a "commit" and/or "optimize" entry also if you want to trigger replication on futures commits or optimizes.
backupAfter	String specifying action after which a backup should occur. Valid values are commit, optimize, or startup. There can be multiple values for this parameter. It is not required for replication, it just makes a backup.
maxNumberOfBackups	Integer specifying how many backups to keep. This can be used to delete all but the most recent N backups.
confFiles	The configuration files to replicate, separated by a comma.

commitReserveDuration	If your commits are very frequent and your network is slow, you can tweak this parameter to increase the amount of time taken to download 5Mb from the master to a slave. The default is 10 seconds.
-----------------------	--

The example below shows how to configure the Replication RequestHandler on a master server.

```
<requestHandler name="/replication" class="solr.ReplicationHandler" >
  <lst name="master">
    <str name="replicateAfter">optimize</str>
    <str name="backupAfter">optimize</str>
    <int name="maxNumberOfBackups">2</int> -->
    <str name="confFiles">schema.xml,stopwords.txt,elevate.xml</str>
    <str name="commitReserveDuration">00:00:10</str>
  </lst>
</requestHandler>
```

Replicating solrconfig.xml

In the configuration file on the master server, include a line like the following:

```
<str name="confFiles">solrconfig_slave.xml:solrconfig.xml,x.xml,y.xml</str>
```

This ensures that the local configuration `solrconfig_slave.xml` will be saved as `solrconfig.xml` on the slave. All other files will be saved with their original names.

On the master server, the file name of the slave configuration file can be anything, as long as the name is correctly identified in the `confFiles` string; then it will be saved as whatever file name appears after the colon ':'.

Configuring the Replication RequestHandler on a Slave Server

The code below shows how to configure a ReplicationHandler on a slave.

```
<requestHandler name="/replication" class="solr.ReplicationHandler" >
  <lst name="slave">

    <!--fully qualified url for the replication handler of master. It is possible
    to pass on this as
      a request param for the fetchindex command-->

    <str name="masterUrl">http://remote_host:port/solr/corename/replication</str>

    <!--Interval in which the slave should poll master .Format is HH:mm:ss . If
```

```
this is absent slave does not
poll automatically.

But a fetchindex can be triggered from the admin or the http API -->

<str name="pollInterval">00:00:20</str>

<!-- THE FOLLOWING PARAMETERS ARE USUALLY NOT REQUIRED-->

<!--to use compression while transferring the index files. The possible values
are internal|external
    if the value is 'external' make sure that your master Solr has the settings to
    honor the
        accept-encoding header.
    See here for details: http://wiki.apache.org/solr/SolrHttpCompression
    If it is 'internal' everything will be taken care of automatically.
    USE THIS ONLY IF YOUR BANDWIDTH IS LOW . THIS CAN ACTUALLY SLOWDOWN
    REPLICATION IN A LAN-->

<str name="compression">internal</str>

<!--The following values are used when the slave connects to the master to
download the index files.
    Default values implicitly set as 5000ms and 10000ms respectively. The user
DOES NOT need to specify
    these unless the bandwidth is extremely low or if there is an extremely high
latency-->

<str name="httpConnTimeout">5000</str>
<str name="httpReadTimeout">10000</str>

<!-- If HTTP Basic authentication is enabled on the master, then the slave can
be
    configured with the following -->

<str name="httpBasicAuthUser">username</str>
<str name="httpBasicAuthPassword">password</str>
</lst>
</requestHandler>
```

 If you are not using cores, then you simply omit the corename parameter above in the masterUrl. To ensure that the URL is correct, just hit the URL with a browser. You must get a status OK response.

Setting Up a Repeater with the ReplicationHandler

A master may be able to serve only so many slaves without affecting performance. Some organizations have deployed slave servers across multiple data centers. If each slave downloads the index from a remote data center, the resulting download may consume too much network bandwidth. To avoid performance degradation in cases like this, you can configure one or more slaves as repeaters. A repeater is simply a node that acts as both a master and a slave.

- To configure a server as a repeater, the definition of the Replication requestHandler in the solrconfig.xml file must include file lists of use for both masters and slaves.
- Be sure to set the replicateAfter parameter to commit, even if replicateAfter is set to optimize on the main master. This is because on a repeater (or any slave), a commit is called only after the index is downloaded. The optimize command is never called on slaves.
- Optionally, one can configure the repeater to fetch compressed files from the master through the compression parameter to reduce the index download time.

Here is an example of a ReplicationHandler configuration for a repeater:

```
<requestHandler name="/replication" class="solr.ReplicationHandler">
    <lst name="master">
        <str name="replicateAfter">commit</str>
        <str name="confFiles">schema.xml,stopwords.txt,synonyms.txt</str>
    </lst>
    <lst name="slave">
        <str name="masterUrl">http://master.solr.company.com:8983/solr/replication</str>
        <str name="pollInterval">00:00:60</str>
    </lst>
</requestHandler>
```

Commit and Optimize Operations

When a commit or optimize operation is performed on the master, the RequestHandler reads the list of file names which are associated with each commit point. This relies on the replicateAfter parameter in the configuration to decide which types of events should trigger replication.

Setting on the Master	Description
commit	Triggers replication whenever a commit is performed on the master index.
optimize	Triggers replication whenever the master index is optimized.
startup	Triggers replication whenever the master index starts up.

The replicateAfter parameter can accept multiple arguments. For example:

```
<str name="replicateAfter">startup</str>
<str name="replicateAfter">commit</str>
<str name="replicateAfter">optimize</str>
```

Slave Replication

The master is totally unaware of the slaves. The slave continuously keeps polling the master (depending on the pollInterval parameter) to check the current index version the master. If the slave finds out that the master has a newer version of the index it initiates a replication process. The steps are as follows:

- The slave issues a filelist command to get the list of the files. This command returns the names of the files as well as some metadata (for example, size, a lastmodified timestamp, an alias if any).
- The slave checks with its own index if it has any of those files in the local index. It then runs the filecontent command to download the missing files. This uses a custom format (akin to the HTTP chunked encoding) to download the full content or a part of each file. If the connection breaks in between, the download resumes from the point it failed. At any point, the slave tries 5 times before giving up a replication altogether.
- The files are downloaded into a temp directory, so that if either the slave or the master crashes during the download process, no files will be corrupted. Instead, the current replication will simply abort.
- After the download completes, all the new files are 'mov'ed to the live index directory and the file's timestamp is same as its counterpart in on the master master.
- A commit command is issued on the slave by the Slave's ReplicationHandler and the new index is loaded.

Replicating Configuration Files

To replicate configuration files, list them using using the confFiles parameter. Only files found in the conf directory of the master's Solr instance will be replicated

Solr replicates configuration files only when the index itself is replicated. That means even if a configuration file is changed on the master, that file will be replicated only after there is a new commit/optimize on master's index.

Unlike the index files, where the timestamp is good enough to figure out if they are identical, configuration files are compared against their checksum. The schema.xml files (on master and slave) are judged to be identical if their checksums are identical.

As a precaution when replicating configuration files, Solr copies configuration files to a temporary directory before moving them into their ultimate location in the conf directory. The old configuration files are then renamed and kept in the same conf/ directory. The ReplicationHandler does not automatically clean up these old files.

If a replication involved downloading of at least one configuration file, the ReplicationHandler issues a core-reload command instead of a commit command.

Resolving Corruption Issues on Slave Servers

If documents are added to the slave, then the slave is no longer in sync with its master. However, the slave will not undertake any action to put itself in sync, until the master has new index data. When a commit operation takes place on the master, the index version of the master becomes different from that of the slave. The slave then fetches the list of files and finds that some of the files present on the master are also present in the local index but with different sizes and timestamps. This means that the master and slave have incompatible indexes. To correct this problem, the slave then copies all the index files from master to a new index directory and asks the core to load the fresh index from the new directory.

HTTP API Commands for the ReplicationHandler

You can use the HTTP commands below to control the ReplicationHandler's operations.

Command	Description
http://_master_host_:_port_/solr/replication?command=enablereplication	Enables replication on the master for all its slaves.
http://_master_host_:_port_/solr/replication?command=disablereplication	Disables replication on the master for all its slaves.
http://_host_:_port_/solr/replication?command=indexversion	Returns the version of the latest replicatable index on the specified master or slave

http://_slave_host_:_port_/solr/replication?command=fetchindex	Forces the specified slave to fetch a copy of the index from its master. If you like, you can pass an extra attribute such as masterUrl or compression (or any other parameter which is specified in the <lst name="slave"> tag) to do a one time replication from a master. This obviates the need for hard-coding the master in the slave.
http://_slave_host_:_port_/solr/replication?command=abortfetch	Aborts copying an index from a master to the specified slave.
http://_slave_host_:_port_/solr/replication?command=enablepoll	Enables the specified slave to poll for changes on the master.
http://_slave_host_:_port_/solr/replication?command=disablepoll	Disables the specified slave from polling for changes on the master.
http://_slave_host_:_port_/solr/replication?command=details	Retrieves configuration details and current status.

<code>http://host:port/solr/replication?command=filelist&indexversion=<index-version-number></code>	Retrieves a list of Lucene files present in the specified host's index. You can discover the version number of the index by running the indexversion command.
<code>http://_master_host_:_port_/solr/replication?command=backup</code>	Creates a backup on master if there are committed index data in the server; otherwise, does nothing. This command is useful for making periodic backups. The numberToKeep request parameter can be used with the backup command unless the maxNumberOfBackups initialization parameter has been specified on the handler – in which case maxNumberOfBackups is always used and attempts to use the numberToKeep request parameter will cause an error.

Index Replication using ssh and rsync

Solr supports ssh/rsync-based replication. *This mechanism only works on systems that support removing open hard links.*

Solr distribution is similar in concept to database replication. All collection changes come to one master Solr server. All production queries are done against query slaves. Query slaves receive all their collection changes indirectly — as new versions of a collection which they pull from the master. These collection downloads are polled for on a cron'd basis.

A collection is a directory of many files. Collections are distributed to the slaves as snapshots of these files. Each snapshot is made up of hard links to the files so copying of the actual files is not necessary when snapshots are created. Lucene only *significantly* rewrites files following an optimization command. Generally, once a file is written, it will change very little, if at all. This makes the underlying transport of rsync very useful. Files that have already been transferred and have not changed do not need to be re-transferred with the new edition of a collection.

The Snapshot and Distribution Process

Here are the steps that Solr follows when replicating an index:

1. The **snapshooter** command takes snapshots of the collection on the master. It runs when invoked by Solr after it has done a commit or an optimize.
2. The **snappuller** command runs on the query slaves to pull the newest snapshot from the master. This is done via rsync in daemon mode running on the master for better performance and lower CPU utilization over rsync using a remote shell program as the transport.
3. The **snapinstaller** runs on the slave after a snapshot has been pulled from the master. This signals the local Solr server to open a new index reader, then auto-warming of the cache(s) begins (in the new reader), while other requests continue to be served by the original index reader. Once auto-warming is complete, Solr retires the old reader and directs all new queries to the newly cache-warmed reader.
4. All distribution activity is logged and written back to the master to be viewable on the distribution page of its GUI.
5. Old versions of the index are removed from the master and slave servers by a cron'd **snapcleaner**.

If you are building an index from scratch, distribution is the final step of the process.

Manual copying of index files is not recommended; however, running distribution commands manually (that is, not relying on crond to run them) is perfectly fine.

Snapshot Directories

Snapshots are stored in directories whose names follow this format: `snapshot.yyyymmddHHMMSS`

All the files in the index directory are hard links to the latest snapshot. This design offers these advantages:

- The Solr implementation can keep multiple snapshots on each host without needing to keep multiple copies of index files that have not changed.

- File copying from master to slave is very fast.
- Taking a snapshot is very fast as well.

Solr Distribution Scripts

For the Solr distribution scripts, the name of the index directory is defined either by the environment variable `data_dir` in the configuration file `solr/conf/scripts.conf` or the command line argument `-d`. It should match the value used by the Solr server which is defined in `solr/conf/solrconfig.xml`.

All Solr collection distribution scripts are bundled in a Solr release and reside in the directory `solr/src/scripts`. LucidWorks recommends that you install the scripts in a `solr/bin/` directory.

Collection distribution scripts create and prepare for distribution a snapshot of a search collection after each commit and optimize request if the `postCommit` and `postOptimize` event listener is configured in `solrconfig.xml` to execute **snapshotter**.

The **snapshotter** script creates a directory `snapshot.<ts>`, where `<ts>` is a timestamp in the format, `yyyymmddHHMMSS`. It contains hard links to the data files.

Snapshots are distributed from the master server when the slaves pull them, "smartcopying" the snapshot directory that contains the hard links to the most recent collection data files.

Name	Description
snapshotter	Creates a snapshot of a collection. Snapshotter is normally configured to run on the master Solr server when a commit or optimize happens. Snapshotter can also be run manually, but one must make sure that the index is in a consistent state, which can only be done by pausing indexing and issuing a commit.
snappuller	A shell script that runs as a cron job on a slave Solr server. The script looks for new snapshots on the master Solr server and pulls them.
snappuller-enable	Creates the file <code>solr/logs/snappuller-enabled</code> , whose presence enables snappuller.
snapinstaller	Installs the latest snapshot (determined by the timestamp) into the place, using hard links (similar to the process of taking a snapshot). Then <code>solr/logs/snapshot.current</code> is written and scp'd (secure copied) back to the master Solr server. snapinstaller then triggers the Solr server to open a new Searcher.
snapcleaner	Runs as a cron job to remove snapshots more than a configurable number of days old or all snapshots except for the most recent n number of snapshots. Also can be run manually.

rsyncd-start	Starts the rsyncd daemon on the master Solr server which handles collection distribution requests from the slaves.
rsyncd daemon	Efficiently synchronizes a collection--between master and slaves--by copying only the files that actually changed. In addition, rsync can optionally compress data before transmitting it.
rsyncd-stop	Stops the rsyncd daemon on the master Solr server. The stop script then makes sure that the daemon has in fact exited by trying to connect to it for up to 300 seconds. The stop script exits with error code 2 if it fails to stop the rsyncd daemon.
rsyncd-enable	Creates the file <code>solr/logs/rsyncd-enabled</code> , whose presence allows the rsyncd daemon to run, allowing replication to occur.
rsyncd-disable	Removes the file <code>solr/logs/rsyncd-enabled</code> , whose absence prevents the rsyncd daemon from running, preventing replication.

For more information about usage arguments and syntax see the [SolrCollectionDistributionScripts](#) page on the Solr Wiki.

Solr Distribution-related Cron Jobs

The distribution process is automated through the use of cron jobs. The cron jobs should run under the user ID that the Solr server is running under.

Cron Job	Description
snapcleaner	<p>The snapcleaner job should be run out of cron at the regular basis to clean up old snapshots. This should be done on both the master and slave Solr servers. For example, the following cron job runs everyday at midnight and cleans up snapshots 8 days and older:</p> <pre>0 0 * * * <solr.solr.home>/solr/bin/snapcleaner -D 7</pre> <p>Additional cleanup can always be performed on-demand by running snapcleaner manually.</p>
snappuller snapinstaller	On the slave Solr servers, snappuller should be run out of cron regularly to get the latest index from the master Solr server. It is a good idea to also run snapinstaller with snappuller back-to-back in the same crontab entry to install the latest index once it has been copied over to the slave Solr server.

For example, the following cron job runs every 5 minutes to keep the slave Solr server in sync with the master Solr server:

```
0,5,10,15,20,25,30,35,40,45,50,55 * * * *
<solr.solr.home>/solr/bin/snappuller;<solr.solr.home>/solr/bin/snapinstaller
```

 Modern cron allows this to be shortened to */5 * * * *....

Performance Tuning for Script-based Replication

Because fetching a master index uses the rsync utility, which transfers only the segments that have changed, replication is normally very fast. However, if the master server has been optimized, then rsync may take a long time, because many segments will have been changed in the process of optimization.

- If replicating to multiple slaves consumes too much network bandwidth, consider the use of a repeater.
- Make sure that slaves do not pull from the master so frequently that a previous replication is still running when a new one is started. In general, it's best to allow at least a minute for the replication process to complete. But in configurations with low network bandwidth or a very large index, even more time may be required.

Commit and Optimization

On a very large index, adding even a few documents and then running an optimize operation causes the complete index to be rewritten. This consumes a lot of disk I/O and impacts query performance. Optimizing a very large index may even involve copying the index twice and calling optimize at the beginning *and* at the end. If some documents have been deleted, the first optimize call will rewrite the index even before the second index is merged.

Optimization is an I/O intensive process, as the entire index is read and re-written in optimized form. Anecdotal data shows that optimizations on modest server hardware can take around 5 minutes per GB, although this obviously varies considerably with index fragmentation and hardware bottlenecks. We do not know what happens to query performance on a collection that has not been optimized for a long time. We *do* know that it will get worse as the collection becomes more fragmented, but how much worse is very dependent on the manner of updates and commits to the collection. The setting of the `mergeFactor` attribute affects performance as well. Dividing a large index with millions of documents into even as few as five segments may degrade search performance by as much as 15-20%.

While optimizing has many benefits, a rapidly changing index will not retain those benefits for long, and since optimization is an intensive process, it may be better to consider other options, such as lowering the merge factor (discussed in this Guide in the section on [Configuring the Lucene Index Writers](#)).

Distribution and Optimization

The time required to optimize a master index can vary dramatically. A small index may be optimized in minutes. A very large index may take hours. The variables include the size of the index and the speed of the hardware.

Distributing a newly optimized collection may take only a few minutes or up to an hour or more, again depending on the size of the index and the performance capabilities of network connections and disks. During optimization the machine is under load and does not process queries very well. Given a schedule of updates being driven a few times an hour to the slaves, we cannot run an optimize with every committed snapshot.

Copying an optimized collection means that the **entire** collection will need to be transferred during the next snappull. This is a large expense, but not nearly as huge as running the optimize everywhere. Consider this example: on a three-slave one-master configuration, distributing a newly-optimized collection takes approximately 80 seconds *total*. Rolling the change across a tier would require approximately ten minutes per machine (or machine group). If this optimize were rolled across the query tier, and if each collection being optimized were disabled and not receiving queries, a rollout would take at least twenty minutes and potentially as long as an hour and a half. Additionally, the files would need to be synchronized so that the *following* rsync, snappull would not think that the independently optimized files were different in any way. This would also leave the door open to independent corruption of collections instead of each being a perfect copy of the master.

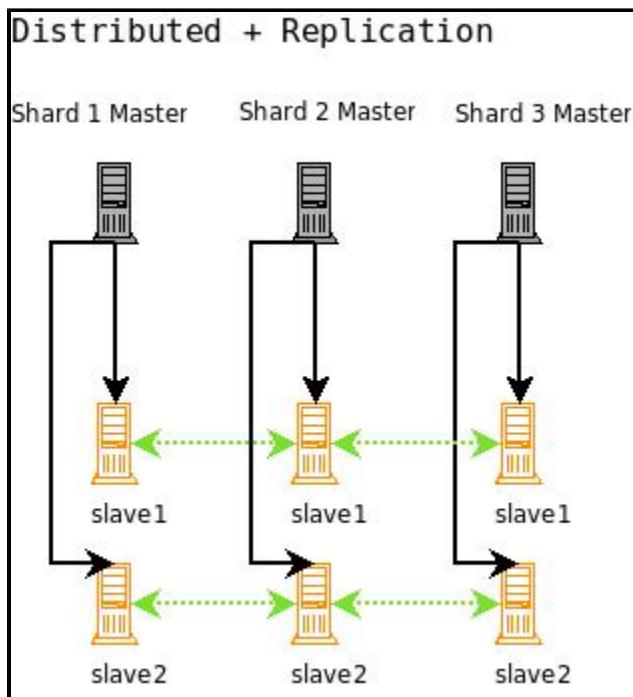
Optimizing on the master allows for a straight-forward optimization operation. No query slaves need to be taken out of service. The optimized collection can be distributed in the background as queries are being normally serviced. The optimization can occur at any time convenient to the application providing collection updates.

Combining Distribution and Replication

When your index is too large for a single machine and you have a query volume that single shards cannot keep up with, it's time to replicate each shard in your distributed search setup.

The idea is to combine distributed search with replication. As shown in the figure below, a combined distributed-replication configuration features a master server for each shard and then $1-n$ slaves that are replicated from the master. As in a standard replicated configuration, the master server handles updates and optimizations without adversely affecting query handling performance.

Query requests should be load balanced across each of the shard slaves. This gives you both increased query handling capacity and fail-over backup if a server goes down.



A Solr configuration combining both replication and master-slave distribution.

None of the master shards in this configuration know about each other. You index to each master, the index is replicated to each slave, and then searches are distributed across the slaves, using one slave from each master/slave shard.

For high availability you can use a load balancer to set up a virtual IP for each shard's set of slaves. If you are new to load balancing, HAProxy (<http://haproxy.1wt.eu/>) is a good open source software load-balancer. If a slave server goes down, a good load-balancer will detect the failure using some technique (generally a heartbeat system), and forward all requests to the remaining live slaves that served with the failed slave. A single virtual IP should then be set up so that requests can hit a single IP, and get load balanced to each of the virtual IPs for the search slaves.

With this configuration you will have a fully load balanced, search-side fault-tolerant system (Solr does not yet support fault-tolerant indexing). Incoming searches will be handed off to one of the functioning slaves, then the slave will distribute the search request across a slave for each of the shards in your configuration. The slave will issue a request to each of the virtual IPs for each shard, and the load balancer will choose one of the available slaves. Finally, the results will be combined into a single results set and returned. If any of the slaves go down, they will be taken out of rotation and the remaining slaves will be used. If a shard master goes down, searches can still be served from the slaves until you have corrected the problem and put the master back into production.

Merging Indexes

If you need to combine indexes from two different projects or from multiple servers previously used in a distributed configuration, you can use either the `IndexMergeTool` included in `lucene-misc` or the `CoreAdminHandler`.

To merge indexes, they must meet these requirements:

- The two indexes must be compatible: their schemas should include the same fields and they should analyze fields the same way.
- The indexes must not include duplicate data.

Optimally, the two indexes should be built using the same schema.

Using IndexMergeTool

To merge the indexes, do the following:

1. Find the lucene JAR file that your version of Solr is using. You can do this by copying your `solr.war` file somewhere and unpacking it (`jar xvf solr.war`). Your lucene JAR file should be in `WEB-INF/lib`. It is probably called something like `lucene-core-2007-05-20_00-04-53.jar`.
2. Copy it somewhere easy to find.
3. Download a copy of Lucene from <http://www.lucidimagination.com/downloads> and unpack it. The file you're interested in is `contrib/misc/lucene-misc-VERSION.jar`.
4. Make sure that both indexes you want to merge are closed.
5. Issue this command:

```
java -cp /path/to/lucene-core-VERSION.jar:/path/to/lucene-misc-VERSION.jar  
org/apache/lucene/misc/IndexMergeTool  
/path/to/newindex  
/path/to/index1  
/path/to/index2
```

This will create a new index at `/path/to/newindex` that contains both `index1` and `index2`.

6. Copy this new directory to the location of your application's solr index (move the old one aside first, of course) and start Solr.

For example:

```
java -cp /tmp/lucene-core-2007-05-20_00-04-53.jar:  
./lucene-2.2.0/contrib/misc/lucene-misc-2.2.0.jarorg/apache/lucene/misc/IndexMergeTo  
  
. /newindex  
. /app1/solr/data/index  
. /app2/solr/data/index
```

Using CoreAdmin

This method uses the [CoreAdminHandler](#) with either the `indexDir` or `srcCore` parameters.

The `indexDir` parameter is used to define the path to the indexes for the cores that should be merged, and merge them into a 3rd core that must already exist prior to initiation of the merge process. The indexes must exist on the disk of the Solr host, which may make using this in a distributed environment cumbersome. With the `indexDir` parameter, a commit should be called on the cores to be merged (so the IndexWriter will close), and no writes should be allowed on either core until the merge is complete. If writes are allowed, corruption may occur on the merged index. Once complete, a commit should be called on the merged core to make sure the changes are visible to searchers.

The following example shows how to construct the merge command with `indexDir`:

```
http://localhost:8983/solr/admin/cores?action=mergeindexes&core=core0&indexDir=/home
```

In this example, `core` is the new core that is created prior to calling the merge process.

The `srcCore` parameter is used to call the cores to be merged by name instead of defining the path. The cores do not need to exist on the same disk as the Solr host, and the merged core does not need to exist prior to issuing the command. `srcCore` also protects against corruption during creation of the merged core index, so writes are still possible while the merge occurs. However, `srcCore` can only merge Solr Cores - indexes built directly with Lucene should be merged with either the `IndexMergeTool` or the `indexDir` parameter.

The following example shows how to construct the merge command with `srcCore`:

```
http://localhost:8983/solr/admin/cores?action=mergeindexes&core=core0&srcCore=core1&
```

Client APIs

This section discusses the available client APIs for Solr. It covers the following topics:

[Introduction to Client APIs](#): A conceptual overview of Solr client APIs.

[Choosing an Output Format](#): Information about choosing a response format in Solr.

[Using JavaScript](#): Explains why a client API is not needed for JavaScript responses.

[Using Python](#): Information about Python and JSON responses.

[Client API Lineup](#): A list of all Solr Client APIs, with links.

[Using SolrJ](#): Detailed information about SolrJ, an API for working with Java applications.

[Using Solr From Ruby](#): Detailed information about using Solr with Ruby applications.

[MBean Request Handler](#): Describes the MBean request handler for programmatic access to Solr server statistics and information.

Introduction to Client APIs

At its heart, Solr is a Web application, but because it is built on open protocols, any type of client application can use Solr.

HTTP is the fundamental protocol used between client applications and Solr. The client makes a request and Solr does some work and provides a response. Clients use requests to ask Solr to do things like perform queries or index documents.

Client applications can reach Solr by creating HTTP requests and parsing the HTTP responses.

Client APIs encapsulate much of the work of sending requests and parsing responses, which makes it much easier to write client applications.

Clients use Solr's five fundamental operations to work with Solr. The operations are query, index, delete, commit, and optimize.

Queries are executed by creating a URL that contains all the query parameters. Solr examines the request URL, performs the query, and returns the results. The other operations are similar, although in certain cases the HTTP request is a POST operation and contains information beyond whatever is included in the request URL. An index operation, for example, may contain a document in the body of the request.

Solr also features an `EmbeddedSolrServer` that offers a Java API without requiring an HTTP connection. For details, see [Using SolrJ](#).

Choosing an Output Format

Many programming environments are able to send HTTP requests and retrieve responses. Parsing the responses is a slightly more thorny problem. Fortunately, Solr makes it easy to choose an output format that will be easy to handle on the client side.

Specify a response format using the `wt` parameter in a query. The available response formats are documented in [Response Writers](#).

Most client APIs hide this detail for you, so for many types of client applications, you won't ever have to specify a `wt` parameter. In JavaScript, however, the interface to Solr is a little closer to the metal, so you will need to add this parameter yourself.

Using JavaScript

Using Solr from JavaScript clients is so straightforward that it deserves a special mention. In fact, it is so straightforward that there is no client API. You don't need to install any packages or configure anything.

HTTP requests can be sent to Solr using the standard `XMLHttpRequest` mechanism.

Out of the box, Solr can send [JavaScript Object Notation \(JSON\) responses](#), which are easily interpreted in JavaScript. Just add `wt=json` to the request URL to have responses sent as JSON.

For more information and an excellent example, read the SolJSON page on the Solr Wiki:

<http://wiki.apache.org/solr/SolJSON>

Using Python

Solr includes an output format specifically for [Python](#), but [JSON output](#) is a little more robust.

Simple Python

Making a query is a simple matter. First, tell Python you will need to make HTTP connections.

```
from urllib2 import *
```

Now open a connection to the server and get a response. The `wt` query parameter tells Solr to return results in a format that Python can understand.

```
connection = urlopen(
    'http://localhost:8983/solr/select?q=cheese&wt=python')
response = eval(connection.read())
```

Now interpreting the response is just a matter of pulling out the information that you need.

```
print response['response']['numFound'], "documents found."
# Print the name of each document.

for document in response['response']['docs']:
    print "  Name =", document['name']
```

Python with JSON

JSON is a more robust response format, but you will need to add a Python package in order to use it. At a command line, install the `simplejson` package like this:

```
$ sudo easy_install simplejson
```

Once that is done, making a query is nearly the same as before. However, notice that the `wt` query parameter is now `json`, and the response is now digested by `simplejson.load()`.

```
from urllib2 import *
import simplejson
connection = urlopen('http://localhost:8983/solr/select?q=cheese&wt=json')
response = simplejson.load(connection)
print response['response']['numFound'], "documents found."

# Print the name of each document.

for document in response['response']['docs']:
    print "  Name =", document['name']
```

Client API Lineup

The Solr Wiki contains a list of client APIs at <http://wiki.apache.org/solr/IntegratingSolr>.

Here is the list of client APIs, current at this writing (November 2011):

Name	Environment	URL
SolRuby	Ruby	http://wiki.apache.org/solr/SolRuby
DelSolr	Ruby	http://delsolr.rubyforge.org/
acts_as_solr	Rails	http://acts-as-solr.rubyforge.org/ , http://rubyforge.org/projects/background-solr/
Flare	Rails	http://wiki.apache.org/solr/Flare
SolPHP	PHP	http://wiki.apache.org/solr/SolPHP
SolrJ	Java	http://wiki.apache.org/solr/SolJava
Python API	Python	http://wiki.apache.org/solr/SolPython
PySolr	Python	http://code.google.com/p/pysolr/
SolPerl	Perl	http://wiki.apache.org/solr/SolPerl
Solr.pm	Perl	http://search.cpan.org/~garafola/Solr-0.03/lib/Solr.pm
SolrForrest	Forrest/Cocoon	http://wiki.apache.org/solr/SolrForrest
SolrSharp	C#	http://www.codeplex.com/solrsharp
SolColdfusion	ColdFusion	http://solcoldfusion.riaforge.org/
SolrNet	.NET	http://code.google.com/p/solrnet/
AJAX Solr	AJAX	http://github.com/evolvingweb/ajax-solr/wiki

Using SolrJ

SolrJ (also sometimes known as SolJava) is an API that makes it easy for Java applications to talk to Solr. SolrJ hides a lot of the details of connecting to Solr and allows your application to interact with Solr with simple high-level methods.

The center of SolrJ is the `org.apache.solr.client.solrj` package, which contains just five main classes. Begin by creating a `SolrServer`, which represents the Solr instance you want to use. Then send `SolrRequests` or `SolrQuerys` and get back `SolrResponses`.

`SolrServer` is abstract, so to connect to a remote Solr instance, you'll actually create an instance of `org.apache.solr.client.solrj.impl.HttpSolrServer`, which knows how to use HTTP to talk to Solr.

```
String urlString = "http://localhost:8983/solr";
SolrServer solr = new HttpSolrServer(urlString);
```

Creating a `SolrServer` does not make a network connection - that happens later when you perform a query or some other operation - but it will throw `MalformedURLException` if you give it a bad URL string.

Once you have a `SolrServer`, you can use it by calling methods like `query()`, `add()`, and `commit()`.

Building and Running SolrJ Applications

The SolrJ API is included with Solr, so you do not have to download or install anything else. However, in order to build and run applications that use SolrJ, you have to add some libraries to the classpath.

At build time, the examples presented with this section require the following libraries in the classpath (all paths are relative to the root of the Solr installation).

- `apache-solr-common-3.x.0.jar`
- `apache-solr-solrj-4.x.x.jar`

At run time, the examples in this section require the libraries found in the 'dist/solrj-lib' directory.

The Ant script bundled with this sections' examples includes the libraries as appropriate when building and running.

You can sidestep a lot of the messing around with the JAR files by using Maven instead of Ant. All you will need to do to include SolrJ in your application is to put the following dependency in the project's `pom.xml`:

```
<dependency>
  <groupId>org.apache.solr</groupId>
  <artifactId>solr-solrj</artifactId>
  <version>3.x.0</version>
</dependency>
```

If you are worried about the SolrJ libraries expanding the size of your client application, you can use a code obfuscator like [ProGuard](#) to remove APIs that you are not using.

Setting XMLResponseParser

SolrJ uses a binary format, rather than XML, as its default format. Users of earlier Solr releases who wish to continue working with XML must explicitly set the parser to the `XMLResponseParser`, like so:

```
server.setParser(new XMLResponseParser());
```

Performing Queries

Use `query()` to have Solr search for results. You have to pass a `SolrQuery` object that describes the query, and you will get back a `QueryResponse` (from the `org.apache.solr.client.solrj.response` package).

`SolrQuery` has methods that make it easy to add parameters to choose a request handler and send parameters to it. Here is a very simple example that uses the default request handler and sets the `q` parameter:

```
SolrQuery parameters = new SolrQuery();
parameters.set("q", mQueryString);
```

To choose a different request handler, for example, just set the `qt` parameter like this:

```
parameters.set("qt", "/spellCheckCompRH");
```

Once you have your `SolrQuery` set up, submit it with `query()`:

```
QueryResponse response = solr.query(parameters);
```

The client makes a network connection and sends the query. Solr processes the query, and the response is sent and parsed into a `QueryResponse`.

The `QueryResponse` is a collection of documents that satisfy the query parameters. You can retrieve the documents directly with `getResults()` and you can call other methods to find out information about highlighting or facets.

```
SolrDocumentList list = response.getResults();
```

Indexing Documents

Other operations are just as simple. To index (add) a document, all you need to do is create a `SolrInputDocument` and pass it along to the `SolrServer's add()` method.

```
String urlString = "http://localhost:8983/solr";
SolrServer solr = new HttpSolrServer(urlString);
SolrInputDocument document = new SolrInputDocument();
document.addField("id", "552199");
document.addField("name", "Gouda cheese wheel");
document.addField("price", "49.99");
UpdateResponse response = solr.add(document);

Remember to commit your changes!

solr.commit();
```

Uploading Content in XML or Binary Formats

SolrJ lets you upload content in XML and binary formats instead of the default XML format. Use the following to upload using binary format, which is the same format SolrJ uses to fetch results.

```
server.setRequestWriter(new BinaryRequestWriter());
```

EmbeddedSolrServer

The `EmbeddedSolrServer` provides the Java interface described above without requiring an HTTP connection. This is the recommended approach if you need to use Solr in an embedded application. This approach enables you to work with the same Java interface whether or not you have access to HTTP.

 **EmbeddedSolrServer** works only with handlers registered in `solrconfig.xml`. The Solr RequestHandler must be mapped to `/update` for a request to function. For more information about configuring handlers, see [Configuring solrconfig.xml](#).

Note that the following property could be set through JVM level arguments:

```
System.setProperty("solr.solr.home",
"/home/shalinsmangar/work/oss/branch-1.3/example/solr");
CoreContainer.Initializer initializer = new CoreContainer.Initializer();
CoreContainer coreContainer = initializer.initialize();
EmbeddedSolrServer server = new EmbeddedSolrServer(coreContainer, "");
```

If you want to use [MultiCore](#) features (which are described in [Core Admin and Configuring solr.xml](#)), then you should use this:

```
File home = new File( "/path/to/solr/home" );
File f = new File( home, "solr.xml" );
CoreContainer container = new CoreContainer();
container.load( "/path/to/solr/home", f );
EmbeddedSolrServer server = new EmbeddedSolrServer( container, "core name as defined in
solr.xml" );
...
```

Using the ConcurrentUpdateSolrServer

If you are working with Java, you can take advantage of the `ConcurrentUpdateSolrServer` to perform bulk updates at high speed. The `ConcurrentHttpSolrServer` buffers all added documents and writes them into open HTTP connections. This class is thread safe. Although any `SolrServer` request can be made with this implementation, it is only recommended to use the `ConcurrentUpdateSolrServer` for `/update` requests.

You can learn more about the `ConcurrentUpdateSolrServer` at
http://lucene.apache.org/solr/4_0_0/solr-solrj/org/apache/solr/client/solrj/impl/ConcurrentUpdateS

Related Topics

- [SolrJ API documentation](#)



- [Solr Wiki page on SolrJ](#)

- [Indexing and Basic Data Operations](#)

Using Solr From Ruby

For Ruby applications, the `solr-ruby` gem encapsulates the fundamental Solr operations.

At a command line, install `solr-ruby` as follows:

```
$ gem install solr-ruby
Bulk updating Gem source index for: http://gems.rubyforge.org
Successfully installed solr-ruby-0.0.8
1 gem installed
Installing ri documentation for solr-ruby-0.0.8...
Installing RDoc documentation for solr-ruby-0.0.8...
```

This gives you a `Solr::Connection` class that makes it easy to add documents, perform queries, and do other Solr stuff.

`Solr-ruby` takes advantage of Solr's Ruby response writer, which is a subclass of the JSON response writer. This response writer sends information from Solr to Ruby in a form that Ruby can understand and use directly.

Performing Queries

To perform queries, you just need to get a `Solr::Connection` and call its `query` method. Here is a script that looks for cheese. The return value from `query()` is an array of documents, which are dictionaries, so the script iterates through each document and prints out a few fields.

```
require 'rubygems'
require 'solr'
solr = Solr::Connection.new('http://localhost:8983/solr')
response = solr.query('cheese')
response.each do |hit|
  puts hit['id'] + ' ' + hit['name'] + ' ' + hit['price'].to_s
end
```

An example run looks like this:

```
$ ruby query.rb
551299 Gouda cheese wheel 49.99
123 Fresh mozzarella cheese
```

Indexing Documents

Indexing is just as simple. You have to get the `Solr::Connection` just as before. Then call the `add()` and `commit()` methods.

```
require 'rubygems'  
require 'solr'  
solr = Solr::Connection.new('http://localhost:8983/solr')  
solr.add(:id => 123, :name => 'Fresh mozzarella cheese')  
solr.commit()
```

More Information

For more information on solr-ruby, read the page at the Solr Wiki:

<http://wiki.apache.org/solr/solr-ruby>

MBean Request Handler

The MBean Request Handler offers programmatic access to the information provided on the [Statistics](#) and [Info](#) pages of the Admin UI. You can access the MBean Request Handler here: <http://localhost:8983/solr/admin/mbeans>.

The MBean Request Handler accepts the following parameters:

Parameter	Type	Default	Description
key	multivalued	all	Restricts results by object key.
cat	multivalued	all	Restricts results by category name.
stats	boolean	false	Specifies whether statistics are returned with results. You can override the stats parameter on a per-field basis.
wt	multivalued	xml	The output format. This operates the same as the wt parameter in a query .

Examples

To return information about the CACHE category only:

```
http://localhost:8983/solr/admin/mbeans?cat=CACHE
```

To return information and statistics about the CACHE category only:

```
http://localhost:8983/solr/admin/mbeans?stats=true&cat=CACHE
```

To return information for everything, and statistics for everything except the fieldCache:

```
http://localhost:8983/solr/admin/mbeans?stats=true&f.fieldCache.stats=false
```

To return information and statistics for the fieldCache only:

```
http://localhost:8983/solr/admin/mbeans?key=fieldCache&stats=true
```

Further Assistance

LucidWorks (formerly known as Lucid Imagination) is the trusted name in Search, Discovery and Analytics, delivering the only enterprise-grade embedded search development solution built on the power of the Apache Lucene/Solr open source search project. Founded in 2008, the company initially provided support, consulting services, documentation and training for the Apache Lucene/Solr open source search project.

Within a few years, the LucidWorks team realized the need to add value to the open source search platform by developing an extensive layer of services which made Lucene/Solr secure and easier to use and manage. The company shipped the first version of its flagship product, LucidWorks Search, in 2011, followed by LucidWorks Big Data in May 2012. LucidWorks continues to offer support, documentation, consulting services and training products for Lucene/Solr and remains committed to giving back to the Apache Lucene/Solr community.

In addition to providing this Reference Guide for Solr, LucidWorks offers other helpful documentation and tips on a Web site devoted to the Solr and Lucene community, [SearchHub](#). Visit the Web site for:

- Technical Notes on special topics
- White Papers about important search topics and methodologies
- Blog posts about the latest news and events of interest to the Lucene and Solr communities
- Podcasts presenting Lucene and Solr tutorials, as well as interviews with Lucene and Solr committers and customers

There is also a very active user community around Solr and Lucene. The solr-user mailing list is a great resource for questions. To view the archives or subscribe to the list, see http://mail-archives.apache.org/mod_mbox/lucene-solr-user/.

LucidWorks has created a search index for all things Lucene and Solr, called [LucidFind](#). Content to be found there includes: All messages to the mailing lists for Lucene and Solr, the full contents of the Lucene/Solr websites, the Lucene/Solr documentation wikis, all of the Lucid Imagination published content, and mailing lists, websites and wikis for a host of related Apache projects.

For more information about services or software offered by LucidWorks, [contact us online](#) or at:

LucidWorks
3800 Bridge Parkway, Suite 101
Redwood City, CA 94065

Tel: 650.353.4057

Fax: 650.620.9540

Solr Glossary

Where possible, terms are linked to relevant parts of the Solr Reference Guide for more information.

Jump to a letter:

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

B

Boolean Operators

These control the inclusion or exclusion of keywords in a query by using operators such as AND, OR, and NOT.

C

Cluster

In Solr, a cluster is a set of Solr nodes managed as a unit. They may contain many cores, collections, shards, and/or replicas. See also [SolrCloud](#).

Collection

In Solr, one or more documents grouped together in a single logical index. A collection must have a single schema, but can be spread across multiple cores.

In [ZooKeeper](#), a group of cores managed together as part of a SolrCloud installation.

Commit

To make document changes permanent in the index. In the case of added documents, they would be searchable after a *commit*.

Core

An individual Solr instance (represents a logical index). Multiple cores can run on a single node. See also [SolrCloud](#).

Core Reload

To re-initialize Solr after changes to `schema.xml`, `solrconfig.xml` or other configuration files.

D

Distributed Search

Distributed search is one where queries are processed across more than one [shard](#).

Document

One or more Fields and their values that are considered related for indexing. See also [Field](#).

E

Ensemble

A [ZooKeeper](#) term to indicate multiple ZooKeeper instances running simultaneously.

F

Facet

The arrangement of search results into categories based on indexed terms.

Field

The content to be indexed/searched along with metadata defining how the content should be processed by Solr.

I

Inverse Document Frequency (IDF)

A measure of the general importance of a term. It is calculated as the number of total Documents divided by the number of Documents that a particular word occurs in the collection. See <http://en.wikipedia.org/wiki/Tf-idf> and http://lucene.apache.org/core/old_versioned_docs/versions/3_5_0/scoring.html for more info on TF-IDF based scoring and Lucene scoring in particular. See also [Term Frequency](#).

Inverted Index

A way of creating a searchable index that lists every word and the documents that contain those words, similar to an index in the back of a book which lists words and the pages on which they can be found. When performing keyword searches, this method is considered more efficient than the alternative, which would be to create a list of documents paired with every word used in each document. Since users search using terms they expect to be in documents, finding the term before the document saves processing resources and time.

L

Leader

The main shard for each node that routes document adds, updates, or deletes to other shards on the same node. See also [SolrCloud](#).

M

Metadata

Literally, *data about data*. Metadata is information about a document, such as its title, author, or location.

N

Natural Language Query

A search that is entered as a user would normally speak or write, as in, "What is aspirin?"

Node

A JVM instance running Solr. Also known as a Solr server.

O

Overseer

The name of the SolrCloud process that coordinates the clusters. It keeps track of existing nodes and shards, and assigns shards to nodes. See also [SolrCloud](#).

Q

Query Parser

A query parser processes the terms entered by a user.

R

Recall

The ability of a search engine to retrieve *all* of the possible matches to a user's query.

Relevance

The appropriateness of a document to the search conducted by the user.

Replica

A copy of a shard or single logical index, for use in failover or load balancing.

Replication

A method of copying a master index from one server to one or more "slave" or "child" servers.

RequestHandler

Logic and configuration parameters that tell Solr how to handle incoming "requests", whether the requests are to return search results, to index documents, or to handle other custom situations.

S

SearchComponent

Logic and configuration parameters used by request handlers to process query requests. Examples of search components include faceting, highlighting, and "more like this" functionality.

Shard

In SolrCloud, a logical section of a single collection. This may be spread across multiple nodes. See also [SolrCloud](#).

SolrCloud

Umbrella term for a suite of functionality in Solr which allows managing a cluster of Solr servers for scalability, fault tolerance, and high availability.

Solr Schema (`schema.xml`)

The Apache Solr index schema. The schema defines the fields to be indexed and the type for the field (text, integers, etc.) The schema is stored in `schema.xml` and is located in the Solr home conf directory.

SolrConfig (`solrconfig.xml`)

The Apache Solr configuration file. Defines indexing options, RequestHandlers, highlighting, spellchecking and various other configurations. The file, `solrconfig.xml` is located in the Solr home conf directory.

Spell Check

The ability to suggest alternative spellings of search terms to a user, as a check against spelling errors causing few or zero results.

Stopwords

Generally, words that have little meaning to a user's search but which may have been entered as part of a [natural language](#) query. Stopwords are generally very small pronouns, conjunctions and prepositions (such as, "the", "with", or "and")

Suggester

Functionality in Solr that provides the ability to suggest possible query terms to users as they type.

Synonyms

Synonyms generally are terms which are near to each other in meaning and may substitute for one another. In a search engine implementation, synonyms may be abbreviations as well as words, or terms that are not consistently hyphenated. Examples of synonyms in this context would be "Inc." and "Incorporated" or "iPod" and "i-pod".

T

Term Frequency

The number of times a word occurs in a given document. See <http://en.wikipedia.org/wiki/Tf-idf> and http://lucene.apache.org/java/2_3_2/scoring.html for more info on TF-IDF based scoring and Lucene scoring in particular.

See also [Inverse Document Frequency \(IDF\)](#).

Transaction Log

An append-only log of write operations maintained by each node. This log is only required with SolrCloud implementations and is created and managed automatically by Solr.

W

Wildcard

A wildcard allows a substitution of one or more letters of a word to account for possible variations in spelling or tenses.

Z

ZooKeeper

Also known as [Apache ZooKeeper](#). The system used by SolrCloud to keep track of configuration files and node names for a cluster. A ZooKeeper cluster is used as the central configuration store for the cluster, a coordinator for operations requiring distributed synchronization, and the system of record for cluster topology. See also [SolrCloud](#).