

Apache Solr



Yonik Seeley
yonik@apache.org

29 June 2006
Dublin, Ireland

ApacheCon
Europe 06



History

- Search for a replacement search platform
 - commercial: high license fees
 - open-source: no full solutions
- CNET grants code to Apache, Solr enters Incubator 17 Jan 2006
- Solr is a Lucene sub-project
- Users: CNET Reviews, CNET Channel, shopper.com, news.com, nines.org, krugle.com, oodle.com, booklooker.de

Lucene Refresher

- Lucene is a full-text search library
- Add documents to an index via IndexWriter
 - A document is a collection of fields
 - No config files, dynamic field typing
 - Flexible text analysis - tokenizers, filters
- Search for documents via IndexSearcher
Hits = search(Query, Filter, Sort, topN)
- Scoring: $tf * idf * lengthNorm$

2

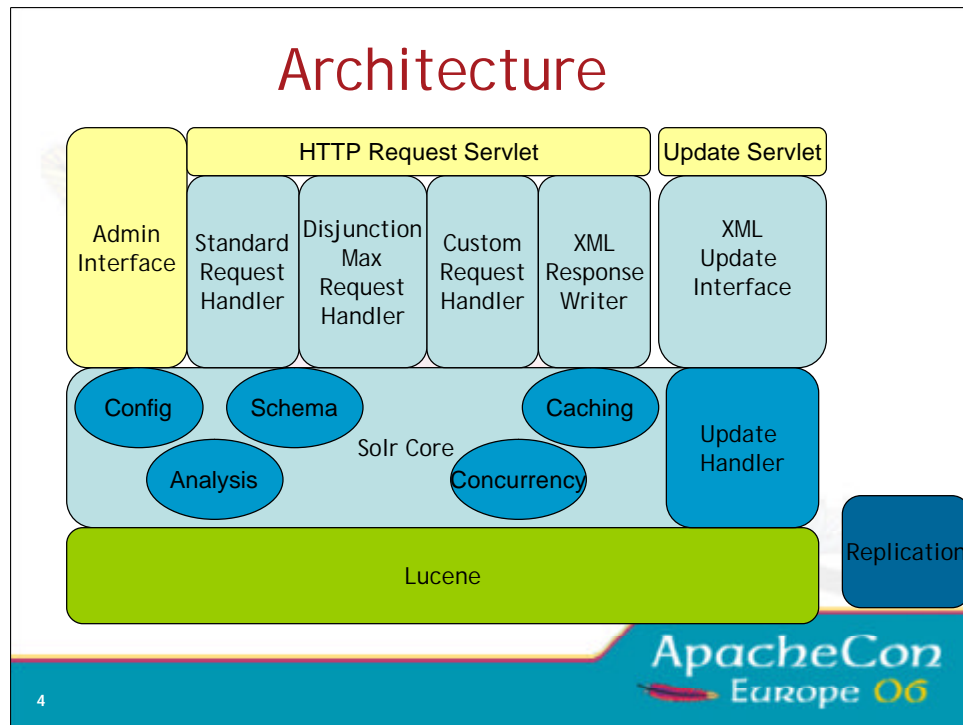
ApacheCon
Europe 06

Scoring:

- **tf** = term frequency in the document – a document that contains more instances of the search term will score higher
- **idf** = inverse document frequency – rarer terms in the corpus are considered more important and score higher
- **lengthNorm** = matches in shorter fields score higher (balances out the tf factor)

What Is Solr

- A full text search server based on Lucene
- XML/HTTP Interfaces
- Loose Schema to define types and fields
- Web Administration Interface
- Extensive Caching
- Index Replication
- Extensible Open Architecture
- Written in Java5, deployable as a WAR



- Yellow components rely on a servlet container, all other components can run outside of a container or embedding or easier testing.
- Configuration is done via XML files
- Request handlers are pluggable, allowing one to do custom server-side processing in a more efficient manner
- Pluggable response writers for different output formats (different XML, JSON, etc)
- Majority of replication implementation is outside the JVM – implemented as shell scripts, rsync, etc

Adding Documents

HTTP POST to /update

```
<add><doc boost="2">  
  <field name="article">05991</field>  
  <field name="title">Apache Solr</field>  
  <field name="subject">An intro...</field>  
  <field name="category">search</field>  
  <field name="category">lucene</field>  
  <field name="body">Solr is a full...</field>  
</doc></add>
```

ApacheCon
Europe 06

5

- Adds are done via an HTTP POST of XML documents
- The “article” field is the unique id field, defined in the schema.
- An add will overwrite other documents with the same id.
- There may be boosts per document or per field – a boost makes a document score higher when querying against it
- Multi-valued fields like category are permissible... just add the field twice
- Adds are only concerned with content, not analysis... notice no indexed, stored, modifiers on the fields – that’s the job of the schema.

Deleting Documents

- Delete by Id

```
<delete><id>05591</id></delete>
```

- Delete by Query (multiple documents)

```
<delete>
```

```
  <query>manufacturer:microsoft</query>
```

```
</delete>
```

Deletes are also done simply by posting an XML document to the /update URL.

Commit

- <commit/> makes changes visible
 - closes IndexWriter
 - removes duplicates
 - opens new IndexSearcher
 - newSearcher/firstSearcher events
 - cache warming
 - “register” the new IndexSearcher
- <optimize/> same as commit, merges all index segments.

7

ApacheCon
Europe 06

commit and optimize are rather heavyweight in Solr/Lucene, so updates should be batched.

After an index searcher becomes “registered”, it can accept live search requests. A single searcher normally handles all live requests concurrently.

Default Query Syntax

Lucene Query Syntax [; sort specification]

1. mission impossible; releaseDate desc
2. +mission +impossible -actor:cruise
3. "mission impossible" -actor:cruise
4. title:spiderman^10 description:spiderman
5. description:"spiderman movie" ~10
6. +HDTV +weight:[0 TO 100]
7. Wildcard queries: te?t, te*t, test*

8

ApacheCon
Europe 06

1. OR query: finds all documents with the word "mission" OR the word "impossible", and sorts descending by the releaseDate field.
2. AND query: finds all documents with BOTH "mission" AND "impossible"
3. phrase query: finds all documents with "mission" directly followed by "impossible"
4. boosting: queries across two different fields, but boosts matches in the title field
5. phrase query slop: finds all documents where "spiderman" occurs within 10 words of "movie"
6. range query: finds HDTV with weight between 0 and 100 inclusive.

Default Parameters

Query Arguments for HTTP GET/POST to /select

param	default	description
q		The query
start	0	Offset into the list of matches
rows	10	Number of documents to return
fl	*	Stored fields to return
qt	standard	Query type; maps to query handler
df	(schema)	Default field to search

Search Results

`http://localhost:8983/solr/select?q=video&start=0&rows=2&fl=name,price`

```
<response><responseHeader><status>0</status>
<QTime>1</QTime></responseHeader>
<result numFound="16173" start="0">
  <doc>
    <str name="name">Apple 60 GB iPod with Video</str>
    <float name="price">399.0</float>
  </doc>
  <doc>
    <str name="name">ASUS Extreme N7800GTX/2DHTV</str>
    <float name="price">479.95</float>
  </doc>
</result>
</response>
```

10

ApacheCon
Europe 06

- status, always 0 in a successful response
- QTime, the server-side query time in milliseconds
- numFound, the total number of documents matching the query
- start, the offset into the ordered list of results
- field types in <doc> include str, boolean, int, long, float, double, date, lst, arr
 - lst is a named list <lst><int name="foo">33</int><int name="bar">42</int></lst>
 - arr is an array <arr><int>33</int><int>42</int></arr>
- multivalued fields are returned in an <arr> element.

Caching

IndexSearcher's view of an index is fixed

- Aggressive caching possible
- Consistency for multi-query requests

filterCache - unordered set of document ids
matching a query

resultCache - ordered subset of document ids
matching a query

documentCache - the stored fields of documents

userCaches - application specific, custom query
handlers

11

ApacheCon
Europe 06

filterCache: key=Query, val=DocSet

resultCache: key=(Query,Sort,Filter), val=DocList

documentCache: key=docid, val=Document

userCache: key=Object, val=Object

Warming for Speed

- Lucene IndexReader warming
 - field norms, FieldCache, tii - the term index
- Static Cache warming
 - Configurable static requests to warm new Searchers
- Smart Cache Warming (autowarming)
 - Using MRU items in the current cache to pre-populate the new cache
- Warming in parallel with live requests

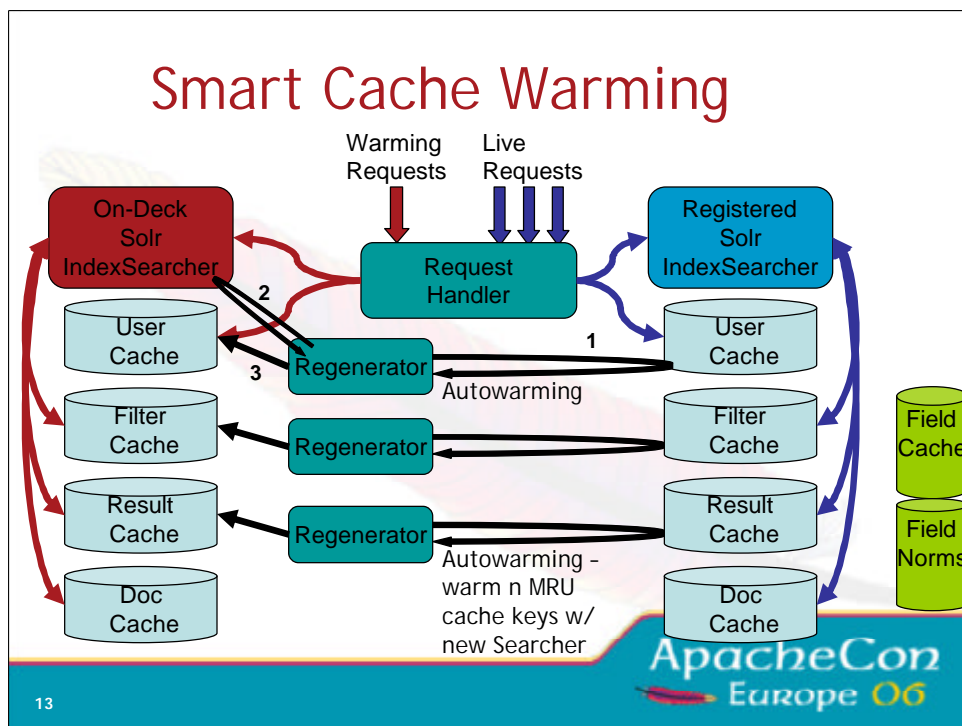
12

ApacheCon
Europe 06

Lucene Internal Caches

- per-field norms are read the first time a field is searched on
- FieldCache entry is populated the first time each indexed field is used to sort
- tii – the term index is loaded on the first random access search

Warming of a new searcher happens in parallel with the current searcher handling live requests



- Static warming requests configured from solrconfig.xml, triggered by events (newSearcher or firstSearcher)
- Autowarming: The top keys from the old (current) cache are re-queried using the new IndexSearcher to pre-populate the new cache(s).
- Cache specific regenerators are used that take keys from old caches and use the new Searcher to pre-populate the new caches.
- The docCache does not have autowarming done since document ids change from one searcher to the next.
- Lucene also has some internal caches (FieldCache and field norms) than benefit from warming.
- After all warming is completed, the new IndexSearcher is registered, and starts serving live requests
- The old index searcher hangs around until all of it's requests have completed, then it is closed.

Schema

- Lucene has no notion of a schema
 - Sorting - string vs. numeric
 - Ranges - val:42 included in val:[1 TO 5] ?
 - Lucene QueryParser has date-range support, but must guess.
- Defines fields, their types, properties
- Defines unique key field, default search field, Similarity implementation

14

ApacheCon
Europe 06

- Fields may still be sparse... just because a field is defined in the schema does not mean it needs to have a value for every document!

Field Definitions

- Field Attributes: name, type, indexed, stored, multiValued, omitNorms

```
<field name="id"      type="string"    indexed="true" stored="true"/>
<field name="sku"     type="textTight" indexed="true" stored="true"/>
<field name="name"    type="text"      indexed="true" stored="true"/>
<field name="reviews" type="text"      indexed="true" stored="false"/>
<field name="category" type="text_ws"  indexed="true" stored="true"
  multiValued="true"/>
```

- Dynamic Fields, in the spirit of Lucene!

```
<dynamicField name="*_i" type="sint"    indexed="true" stored="true"/>
<dynamicField name="*_s" type="string"  indexed="true" stored="true"/>
<dynamicField name="*_t" type="text"    indexed="true" stored="true"/>
```

15

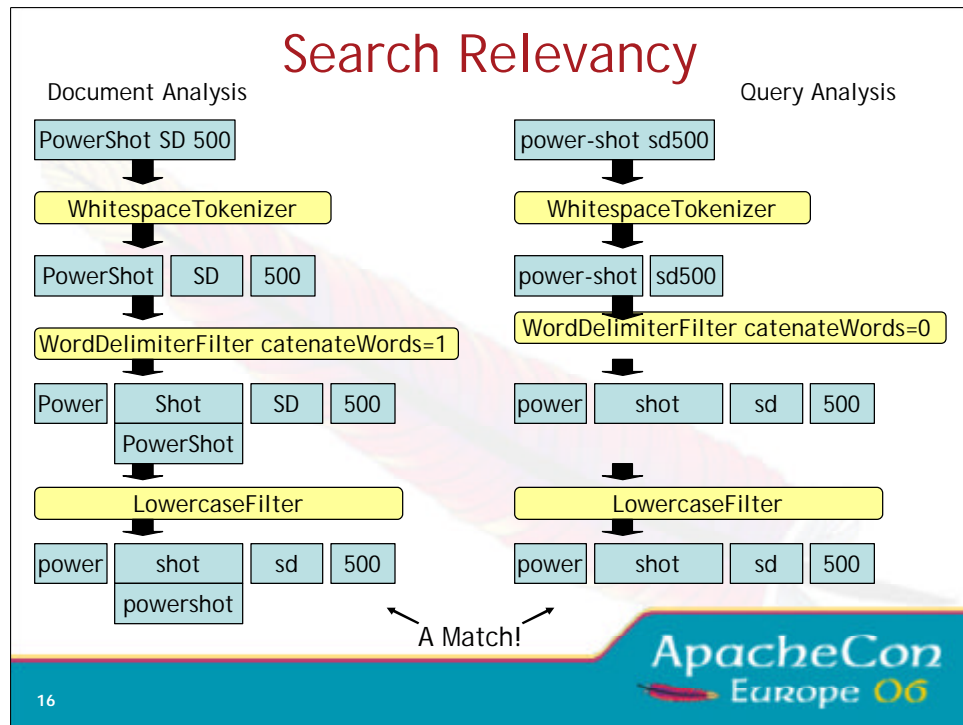
ApacheCon
Europe 06

type: the field type, defined as a fieldtype declaration in the schema

indexed: a field must be indexed if you want to search on it or sort by it

stored: a field must be stored if you want it returned in search results

omitNorms: saves memory, good for sparse fields and when lengthNorm and index-time boosting is not needed



LHS: index analyzer, RHS: query analyzer

"PowerShot SD 500" is the value of a particular field, while "power-shot sd500" is what someone typed into a search box.

When extra tokens can be injected (WordDelimiterFilter or SynonymFilter), you often want different analyzers for indexing and querying.

More than one token may be indexed at the same position ("shot" and "powershot" for example)

Configuring Relevancy

```
<fieldtype name="text" class="solr.TextField">
  <analyzer>
    <tokenizer class="solr.WhitespaceTokenizerFactory"/>
    <filter class="solr.LowerCaseFilterFactory"/>
    <filter class="solr.SynonymFilterFactory"
      synonyms="synonyms.txt" />
    <filter class="solr.StopFilterFactory"
      words="stopwords.txt" />
    <filter class="solr.EnglishPorterFilterFactory"
      protected="protwords.txt"/>
  </analyzer>
</fieldtype>
```

17

ApacheCon
Europe 06

- fieldtype name is arbitrary... the behavior is governed by the implementing class
- one can set defaults for fields on the fieldtype such as indexed, stored, multiValued, omitNorms
- different analyzers may be configured for a TextField... an “index” analyzer, and a “query” analyzer.
- SynonymFilter may inject synonyms or “normalize” by replacing words.
- StopFilter – common words that don’t increase relevancy “and, or, a”, there is a default english stopword list, or one may configure via external file.
- EnglishPorterFilter – Porter stemmer - reduces related words, different tenses, to a common root.

copyField

- Copies one field to another at **index time**
- Usecase: Analyze same field different ways
 - copy into a field with a different analyzer
 - boost exact-case, exact-punctuation matches
 - language translations, thesaurus, soundex

```
<field name="title" type="text" />  
<field name="title_exact" type="text_exact" stored="false" />  
<copyField source="title" dest="title_exact"/>
```

- Usecase: Index multiple fields into single searchable field

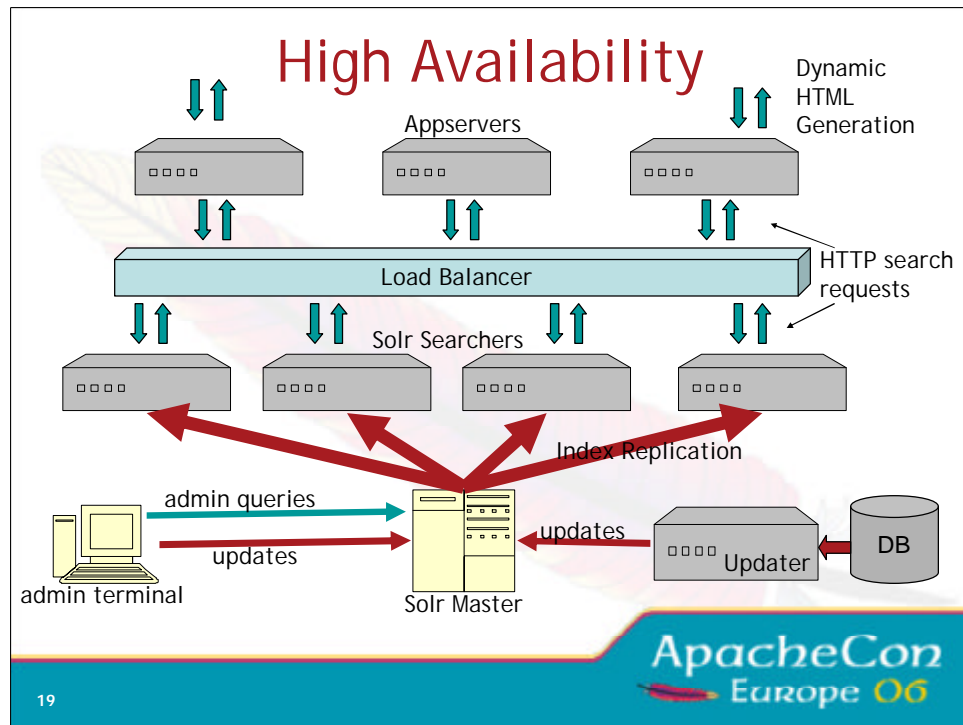
ApacheCon
Europe 06

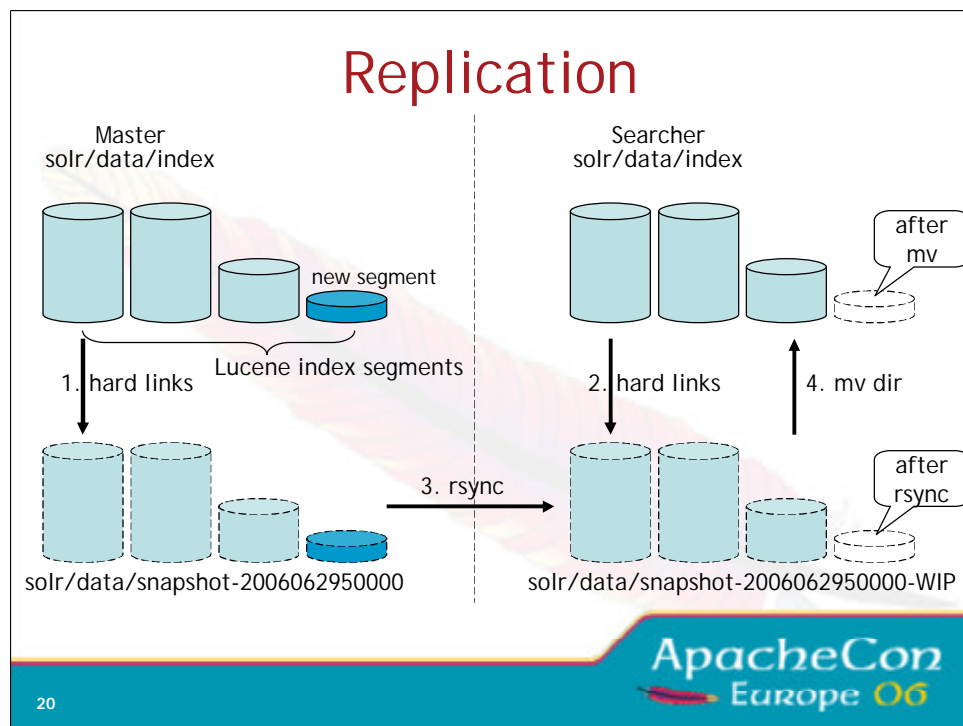
18

copyField decouples data providers (updaters) from details of text analysis and indexing... you can add new text analysis features in the future without having to change clients updating the index.

copyField only takes effect when adding a new document to the index. It causes one field in the input to be copied to another.

Since the fields being copied to contain redundant data, they should normally be unstored.





for performance reasons rsync is instructed to use timestamps & file sizes to determine which files have changed.

Step 1, taking an index snapshot, normally happens for every commit or optimize on the master. This is configurable.

A cron task on the searcher periodically checks the master for a new snapshot. If it exists, this kicks off steps 2-4.

Other replication strategies are possible!

Faceted Browsing Example

DESKTOPS

You found 1045 items for System type: [Budget desktop system](#)
 Too few results? Click a link above to remove that filter, or [remove all filters](#).

Find by price

- Less than \$400 (76)
- \$400 to \$699 (337)
- \$700 to \$999 (468)
- \$1000 to \$1299 (5)

Find by manufacturer

- Dell, Inc. (43)
- Lenovo (490)
- HP (342)
- Acer America Corp. (28)
- Cyberpower Inc (22)
- [See all manufacturers](#)

Find by processor manufacturer

- Intel (804)
- AMD (122)
- Motorola (1)

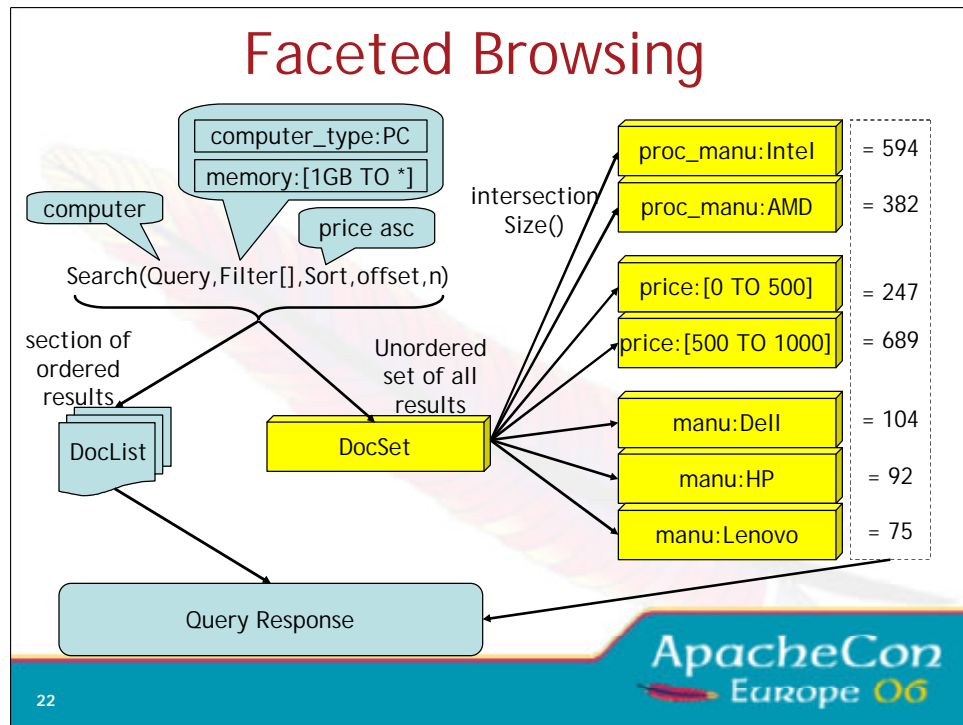
Or find by

- [Clock speed](#)
- [Graphics processor](#)
- [RAM installed](#)
- [Hard drive size](#)
- [OS provided](#)
- [See all](#)

Sort by: [Product name](#) | [Lowest price](#) | [Editors' rating](#) | [Review date](#) Check products to [Compare](#)

 <p>Reviewed on 05/05/2006</p>	<p>Dell Dimension B110 Desktop Computer for Home (Cel-D 2.53GHz/160GB/512MB)</p> <p>Dell's entry-level Dimension B110 series features aging technology and a dated design, but its members will suffice as second PCs for basic tasks.</p> <p>Specs: Celeron D (2.53 GHz), 512 MB, 160 GB, 15 in, Microsoft Windows XP Home Edition</p> <p>Add to my products New! What is this?</p>	<p>\$479 at 1 store</p> <p>Check prices</p>	COMPARE
	<p>Dell Dimension B110 Desktop Computer for Home (Cel-D 2.53GHz/80GB/256MB)</p> <p>Dell's entry-level Dimension B110 series features aging</p>	<p>\$349 at 1 store</p>	COMPARE

Faceted browsing: listing category or facet counts, and allowing the user to narrow search results by those facets.



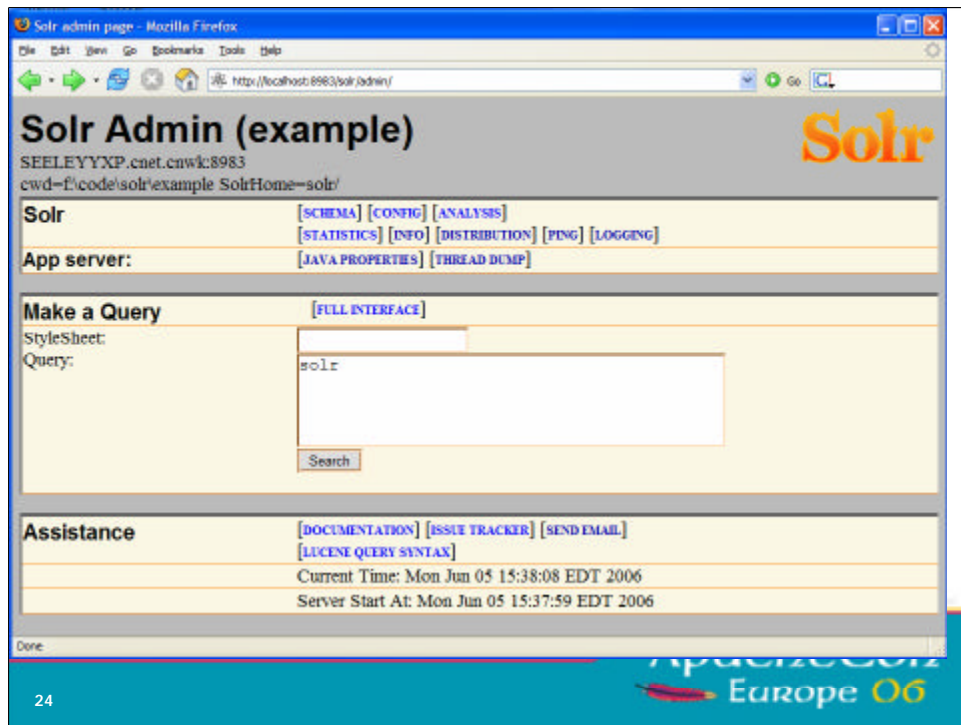
Faceted Browsing currently requires a custom query handler.

Solr Optimizations for Faceted Browsing:

- OpenBitSet – 3 to 4 times faster than `java.util.BitSet` for calculating intersection sizes
- HashDocSet – much faster and more memory efficient for smaller sets

Web Admin Interface

- Show Config, Schema, Distribution info
- Query Interface
- Statistics
 - Caches: lookups, hits, hitratio, inserts, evictions, size
 - RequestHandlers: requests, errors
 - UpdateHandler: adds, deletes, commits, optimizes
 - IndexReader, open-time, index-version, numDocs, maxDocs,
- Analysis Debugger
 - Shows tokens after each Analyzer stage
 - Shows token matches for query vs index



Selling Points

- Fast
- Powerful & Configurable
- High Relevancy
- Mature Product
- Same features as software costing \$\$\$
- Leverage Community
 - Lucene committers, IR experts
 - Free consulting: shared problems & solutions

Where are we going?

- OOTB Simple Faceted Browsing
- Automatic Database Indexing
- Federated Search
 - HA with failover
- Alternate output formats (JSON, Ruby)
- Highlighter integration
- Spellchecker
- Alternate APIs (Google Data, OpenSearch)

Resources

- WWW
 - <http://incubator.apache.org/solr>
 - <http://incubator.apache.org/solr/tutorial.html>
 - <http://wiki.apache.org/solr/>
- Mailing Lists
 - solr-user-subscribe@lucene.apache.org
 - solr-dev-subscribe@lucene.apache.org