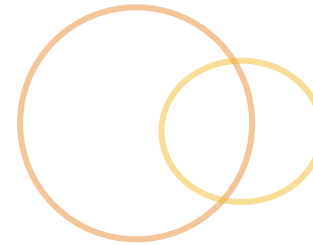
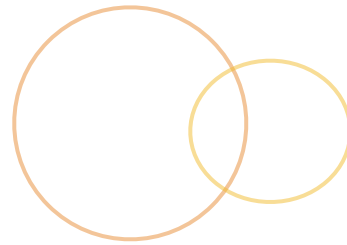


Solr Primer Day 1

Macy's

jeffnb@gmail.com





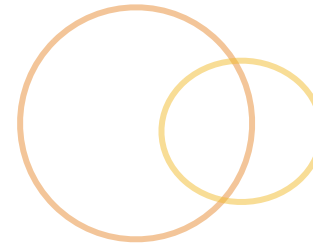
🕒 Download Solr:

🕒 <http://www.apache.org/dyn/closer.lua/lucene/solr/6.5.1>
Github Repo

🕒 <https://github.com/jeffnb/solr-lab-intro>

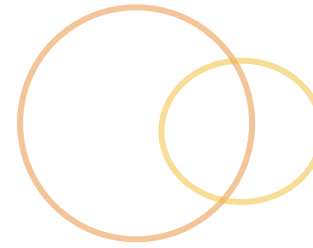
🕒 Run Through “Setup Complete Movies Demo”

Course Outline [Day 1]



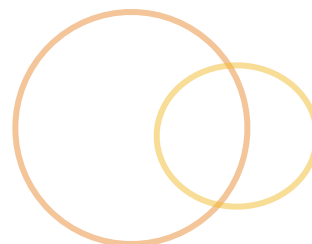
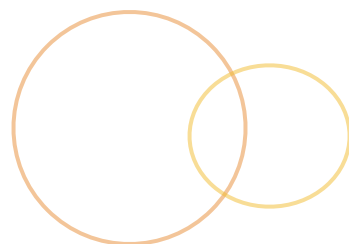
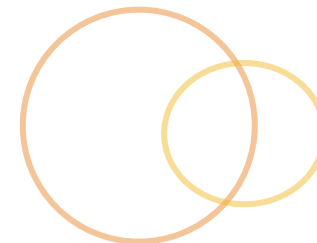
- ② Introduction to Solr
- ② Solr Overview
- ② Solr Admin
- ② Searching
- ② SolrJ Searching
- ② Schema

Course Outline [Day 2]



- 🕒 Schema Continued
- 🕒 Indexing
- 🕒 Configuration
- 🕒 Optimization
- 🕒 SolrCloud
- 🕒 NLP

Part 0: Solr



Solr: The Database Issue



🕒 Filtering simple things

```
SELECT * FROM movies WHERE Year = 1995;
```

🕒 A little more complicated looking in multiple fields

```
SELECT * FROM movies WHERE Title like '%Toy%'  
OR Plot like '%Toy%';
```

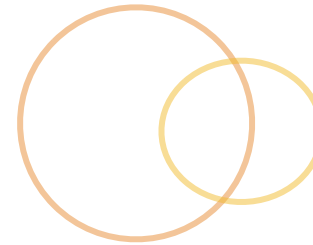
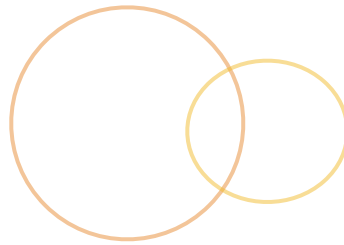
🕒 Even more so if we need to check many fields

```
SELECT * FROM movies WHERE Title like '%Tom%'  
OR Plot like '%Tom%' OR Actors like '%Tom%'  
OR Director like '%Tom%';
```

Solr: The Database Issue



- ⦿ Text field searching is slow
- ⦿ No “best” matches
- ⦿ “Try” != “Tries”

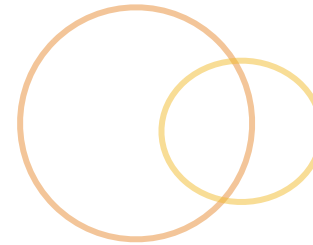
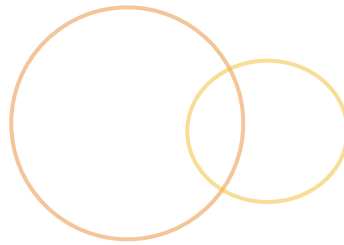


☉ What?

- ☉ Enterprise Level Search Engine

☉ Why?

- ☉ Extremely fast
- ☉ Extremely flexible
- ☉ Extremely powerful
- ☉ Extremely scalable



⦿ What?

- ⦿ Enterprise Level Search Engine

⦿ Why?

- ⦿ Extremely fast
- ⦿ Extremely flexible
- ⦿ Extremely powerful
- ⦿ Extremely scalable



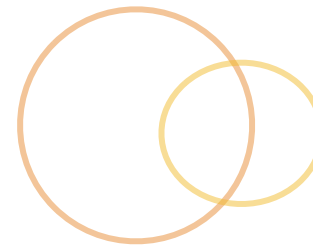
Solr: What Can it Do



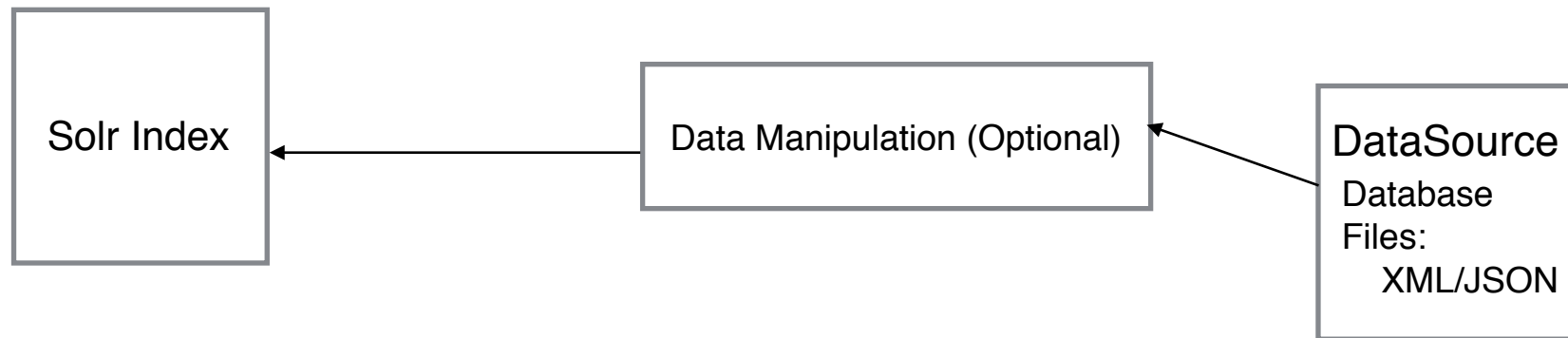
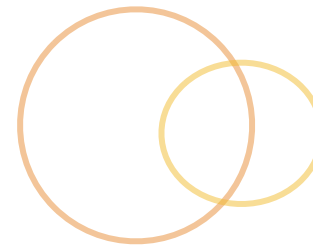
- 🔍 Searching
- 🔍 Faceting
- 🔍 Tokenizing
- 🔍 Spell Checking
- 🔍 Replicating
- 🔍 Similarity-ing

Solr: History

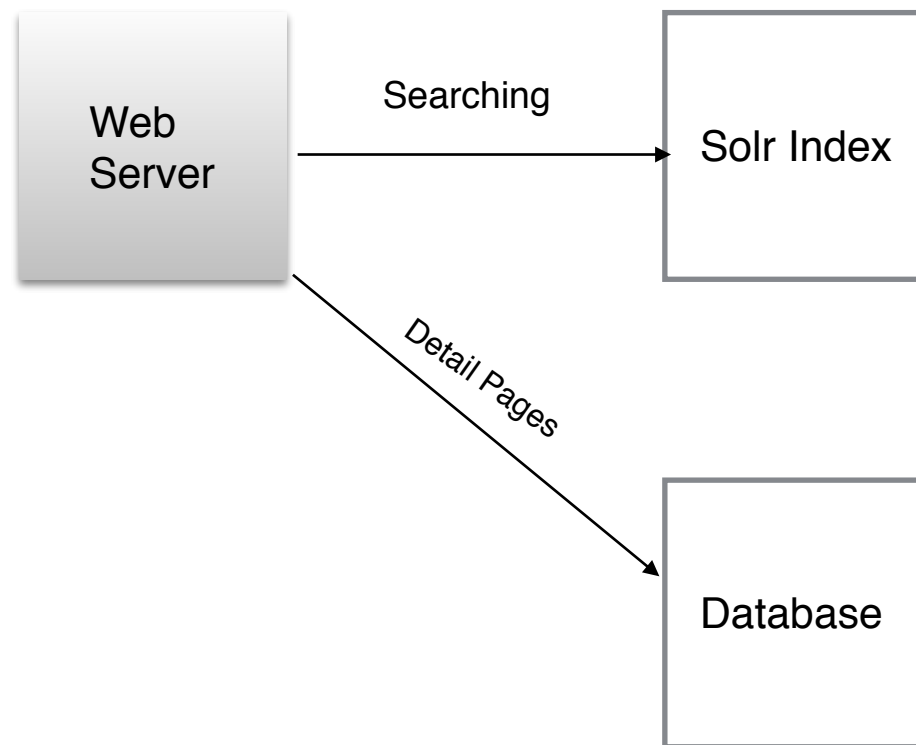
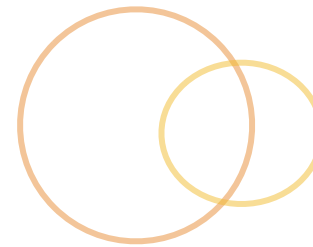
- Built on Lucene
- Versions up with Lucene



Solr: Indexing



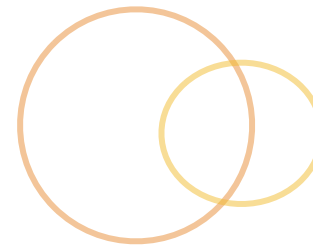
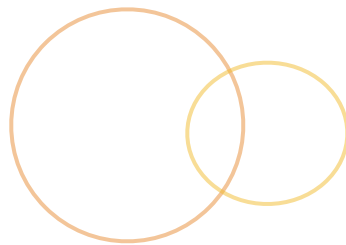
Solr: Searching



Solr: Master/Slave

- ☉ Main server responsible for indexing
- ☉ Slave servers
- ☉ Replication handled within Solr Software

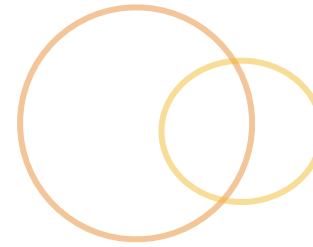
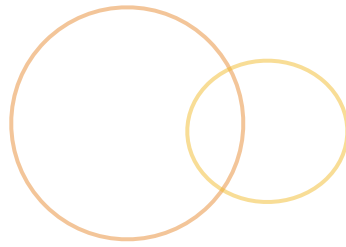
SolrCloud



- ☉ Built to simplify scalability
- ☉ Allows adding nodes to a cluster

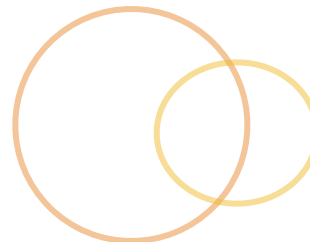
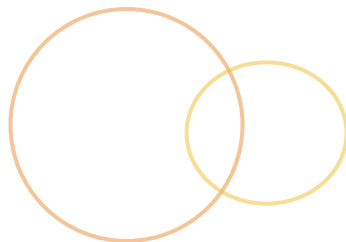
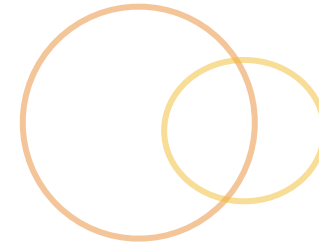


Solr Admin

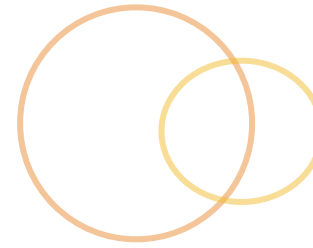


Walkthrough

Part 1: Searching

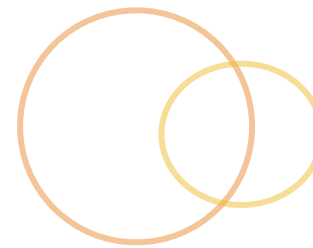


Searching: The Basics



- ② q=<field>:<value>
- ② Wildcards: * and ?
- ② Fuzzy searches: toy~ or toy~1
- ② Proximity: “toy story”~10
- ② Required: +toy
- ② Ranges:
 - ② Year:[1995 TO 1999] - Inclusive
 - ② Title:{Aladdin TO Boy} - Exclusive
- ② Boosting: ^10

Lab: Searching

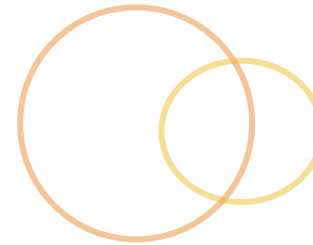
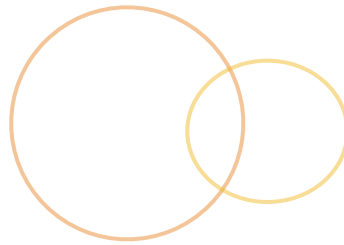


- ② Try to find the movies with Woody Harrelson
- ② Find your favorite movie
- ② Get movies that were made in 2010 and 2011
- ② Return movies that were about disasters in 1999
- ② Find movies that have the word “world” and may have “danger”
- ② Find movies with a fuzzy search for fight
- ② Find movies with Deep in the title or the Plot
- ② Same as previous but boost documents with it in the title by 10

Searching: Common Params



- 🕒 sort DESC:ASC
- 🕒 Field List: `fl=title,plot,year`
- 🕒 Filter Query:
 - 🕒 Applies filter without impacting the score
 - 🕒 `fq=Year:1995&fq=actors:"Tom Hanks"`
- 🕒 Write Type: `wt=json`
- 🕒 Indent to pretty up the return: `indent=true`
- 🕒 Start (offset): `start=10`
- 🕒 Debug/ExplainOther



🕒 Term Frequency

- 🕒 Search: The Brown Cow
- 🕒 Documents more relevant where “the”, “brown” and “cow” occur the most

🕒 Inverse Document Frequency

- 🕒 Adds weight to terms that are infrequent across all documents
- 🕒 “The” occurs many times across all documents
- 🕒 “brown” and “cow” get more weighting

Lab: Common Parameters



- Return json and indent it
- Get a list of movies with Title, Year, score and imdbRating for Adventure and show 50 records
- Sort the previous example on the highest rated and most recent
- Search for “batman” in plot and title but use fq to filter to the year 2013
- Use debug=true to parse through the matching of the previous example
- Play around with searching and see what else you can find

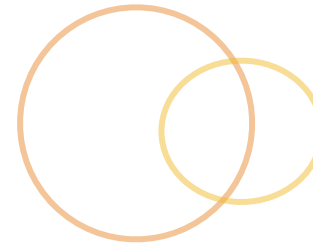
Solr: Facets

- ⦿ Similar to group by
- ⦿ Allows for faceted navigation
- ⦿ Field must be indexed
- ⦿ Facets on tokens not stored data
- ⦿ `facet=true` Turns on the facet system
- ⦿ `facet.field` Picks a field to facet on
- ⦿ `facet.limit` Sets a max amount of facet values
- ⦿ `facet.mincount` Sets minimum doc count for a facet

Solr: Facets

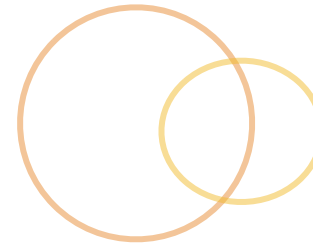
- ② `facet.prefix` Simply filters facets to prefix
- ② `facet.contains` (`contains.ignorecase`): Filters facets to ones that contain the string
- ② `facet.offset`: Starts the facets after the offset

Solr: Facet Ranges



- ② Useful for know ranges of values like ratings
- ② `facet.range`: Which field to do a range facet
- ② `facet.range.start (end)`: Start and end of range
- ② `facet.range.gap`: The amount of each range
- ② All properties can be set on a per field basis with `f.<fieldname>.facet.<parameter>`

Solr: Facet Intervals



- ② Useful for arbitrary groupings
- ② `facet.interval`: set the field
- ② `facet.interval.set`: set an interval
- ② * can be used as any for upper and lower bound
- ② Don't forget you can use
`f.<fieldname>.facet.interval.set`

Lab: Facets

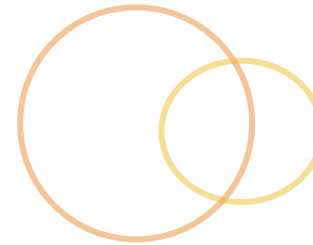
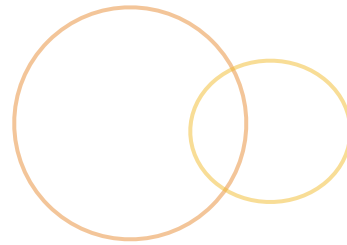
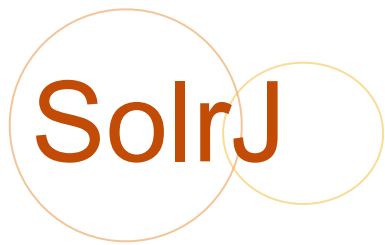


- ② Facet on Year and genres
- ② Facet on actors but limit to only 5 facets
- ② Facet on Year. Limit them to years with more than 1000 movies in the data set
- ② Facet on all Baldwin brothers using contains
- ② Facet on Year in ranges of 5 years mincount of 1
- ② Interval Facet on imdbRating. 0 to 3, 3 to 5, 5 to 7, 8 to 9, 9 to 10

Discussion: Searching

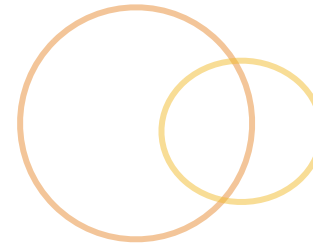


- ② How do you currently use solr?
- ② What are some other ways you could see using solr?



- ⦿ Java client for solr
- ⦿ Extremely limited documentation (1 page)
- ⦿ Versions up with the version of solr
- ⦿ Abstracts many of the communication details
- ⦿ Adding documents easier with java based class
- ⦿ Uses a fast java serialization

SolrJ: Connecting



//Stand alone server example

```
String urlString = "http://localhost:8983/solr/movies";  
SolrClient solr = new HttpSolrClient  
    .Builder(urlString).build();
```

//SolrCloud Connection

```
String zkHostString = "localhost:9983";  
SolrClient solr = new CloudSolrClient.Builder()  
    .withZkHost(zkHostString).build();
```

SolrJ: Searching By Id



```
SolrDocument doc;  
try {  
    //Simple Search by Id  
    doc = solr.getById(id: "tt0113497");  
} catch (SolrServerException e) {  
    e.printStackTrace();  
} catch (IOException e) {  
    e.printStackTrace();  
}
```

SolrJ: Search



```
// Simple search for toy with a few fields  
SolrQuery query = new SolrQuery();  
query.setQuery("Title:future");  
query.setFields("id", "Title", "Plot", "Year");  
query.setRows(10);
```


SolrJ: Sorting



```
// Sorting done with a SortClause  
SolrQuery.SortClause sort;  
sort = new SolrQuery.SortClause( item: "Title",  
    SolrQuery.ORDER.asc);  
query.setSort(sort);
```

SolrJ: Search Response



```
//SolrDocument basically wrapped Map  
String title = (String)doc.get("Title");  
Integer metascore = (Integer)doc.getFieldValue(name: "Metascore");  
Collection<Object> genres = doc.getFieldValues(name: "genres");  
Collection<String> names = doc.getFieldNames();
```

SolrJ: Facets

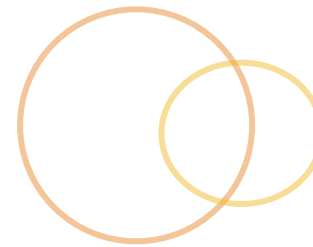
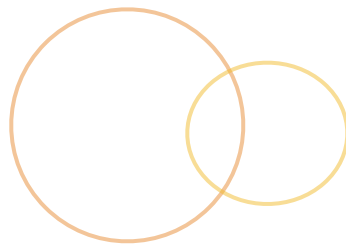
// Facets

```
query.setFacet(true);  
query.addFacetField("Year", "Rated", "genres");  
query.setFacetMinCount(1);  
query.addNumericRangeFacet("Metascore", 1, 100, 10);  
query.setFacetLimit(100);
```

SolrJ: Facet Response



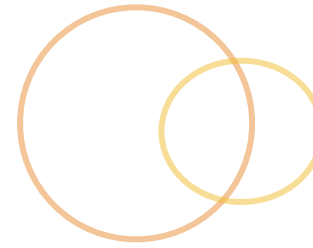
```
// Call solr with search
QueryResponse response = solr.query(query);
// Get Facet Fields
List<FacetField> fields = response.getFacetFields();
for(FacetField ff : fields){
    String name = ff.getName();
    int valueCount = ff.getValueCount();
    // Get individual values/counts
    for(FacetField.Count fieldCount : ff.getValues()){
        String value = fieldCount.getName();
        long count = fieldCount.getCount();
    }
}
```



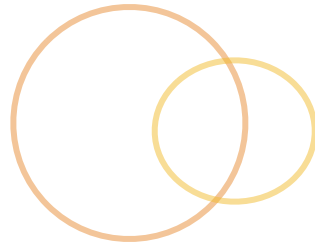
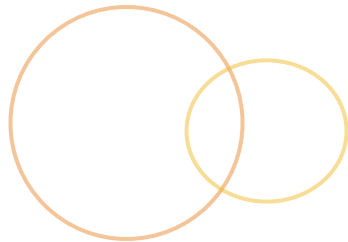
- ② Search for Toy Story in the Title field
 - ② Return: id, Title, Plot, Year
- ② Order the previous search by Year descending
- ② Print the data via the previous query
 - ② Make the output “Title (Year)”
- ② Search for horror movies with “nightmare” in the Title or Plot field
- ② Facet on Year
- ② Print out the facet info as “Year (Count)”

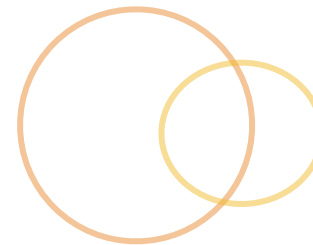
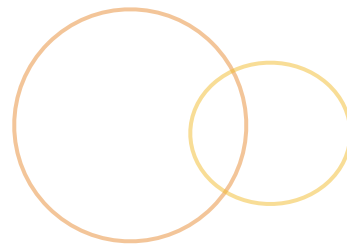
End of day survey

🕒 <http://bit.ly/solr5-8-17>



Part 2: Schemas and Indexes





- Fields store data
- Fields have a type
- Parameters
 - stored: values can be retrieved
 - multiValued: Can store multiple values
 - indexed: Can be searched/faceted
 - required: Document will error if value missing
 - type: The type of field
 - docValues: Orders values to be more efficient

Basic FieldTypes

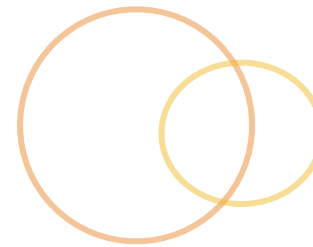
- ⦿ Data types
- ⦿ No tokenization
- ⦿ string: most common
- ⦿ int: `solr.TrieIntField`
- ⦿ long: `solr.TrieLongField`
- ⦿ float: `solr.TrieFloatField`
- ⦿ double: `solr.TrieDoubleField`
- ⦿ date: `solr.TrieDateField`

Lab: Basic Fields



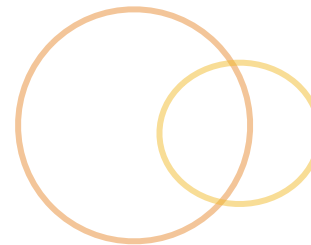
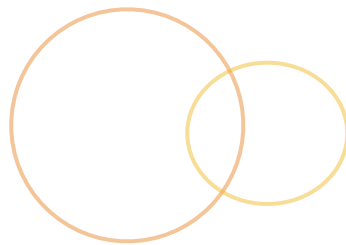
- ① `bin/solr create_core -c labmovies -d <path>/movies-start`
- ① Open file: `server/solr/labmovies/conf/managed_schema`
- ① Find the id field in the file and add fields under. Add remaining fields with built in types from the following: Country, Rated, Language, imdbVotes, Type, Poster, Metascore, Year, actors, genres, directors, writers, Runtime, imdbID, Released, imdbRating

Text FieldTypes



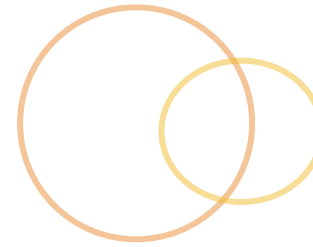
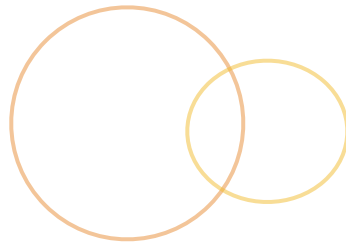
- ⦿ The magic is here
- ⦿ Many pre configured including text_general
- ⦿ Have several steps in processing data
- ⦿ Are the core power of solr

Analizers



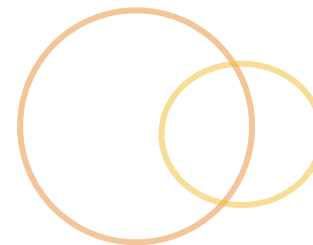
- ② List of instructions to run
- ② 2 Times Run
 - ② Index: Documents going in
 - ② Query: Searching existing
- ② Contain steps to process data
- ② First Step Tokenizer

Tokenizers

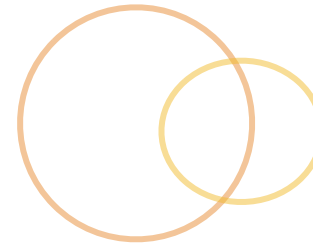
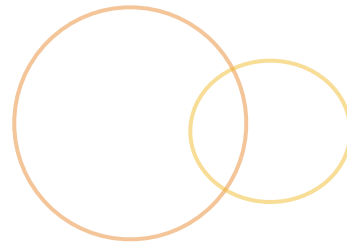


- Tokenizers take input and break it into “tokens”
 - The Quick Brown Fox Jumps Over the Lazy Dog
 - “The”, “Quick”, “Brown”, “Fox”, “Jumps”, “Over”, “the”, “Lazy”, “Dog”
- Can tokenize in a variety of ways including not breaking it up at all

Tokenizer Overview

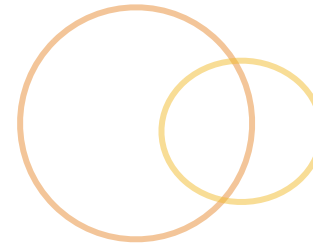


- Standard: Splits based on whitespace and punctuation
- Classic: Standard but doesn't support Unicode Annex
- Keyword: Keeps entire input as a single token
- Letter: Tokens are strings of contiguous letters
- LowerCase: Delimits on non-letters and lowercases
- N-Gram: generates n-gram tokens
- Path Hierarchy: Replaces a delimiter with another value building up a path
- Regular expression: Tokenizes based on regular expression as a delimiter
- Url Email: Splits on white space and tokens preserving email and urls
- Whitespace: Splits on whitespace only



- ② Alter tokens
- ② Produce more tokens
- ② Suppress tokens
- ② Dozens exist

Filters: Heavily Used



- Lowercase Filter: Lowercases all tokens
- Stop Filter: Takes a file and removes all tokens that match words in the file
- Stemmers (porter/snowball): Remove word conjunctions and plurization
- Synonym Filter: Token matched to file and all synonyms added as tokens
- Pattern Replace Filter: Replaces a regular expression match with a string
- WordDelimiter Filter: splits tokens up further into words useful for breaking words apart or concatenating hyphens

Filters: Less Common



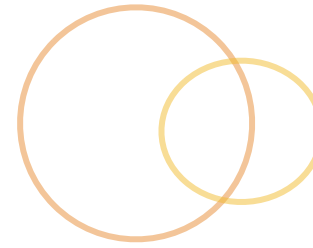
- ⦿ Managed Stopwords/Synonyms: Allows management via rest api
- ⦿ N-Gram Filter: Breaks up token into small chunks
- ⦿ Phonetic Filter: translates tokens into phonetic equivalent. Can replace or add
- ⦿ Length Filter: Filters tokens by min/max length

Lab: Text Fields

- ④ Create a field for Plot and Title
- ④ Create a field type for Awards
- ④ Create a field for awards using your own field type
- ④ Use the analyzer admin to ensure these fields work

End of day survey

📍 <http://bit.ly/solr5-8-17>



Solr Primer

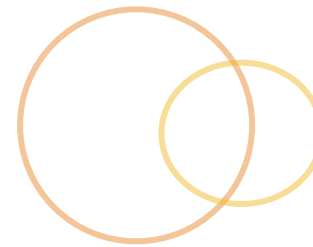
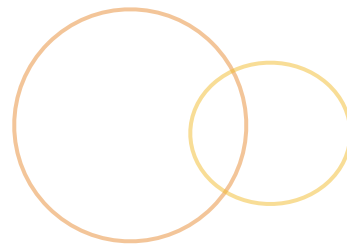
Day 2

Macy's

jeffnb@gmail.com



CopyField

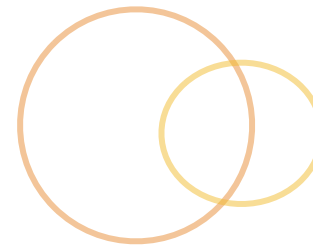


- 🕒 Shortcut to copy one field into another

```
<copyField source="<sourcefield>" dest="<destfield>"/>
```

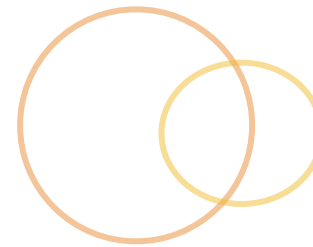
- 🕒 Useful for when facets and search are needed

Lab: CopyField



- Use copyField to copy the following to text:
 - actors, genres, directors, writers, Plot, Year, Title, Awards, language, Rated
- Create fields for (use a name like actors_text):
 - actors
 - genres
 - directors
 - writers
- Use copy field to copy the source fields above to the new _text versions

Dynamic Fields

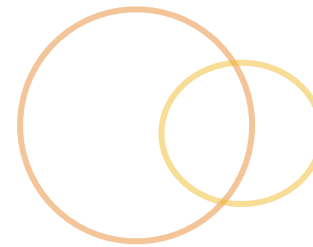


- DynamicField allows wild carding field names

```
<dynamicField name="*_en" type="text_en"
indexed="true" stored="true" multiValued="true"/>
```

- Any field with _en at the end automatically populates into that field
- Can be referenced directly

Lab: DynamicFields

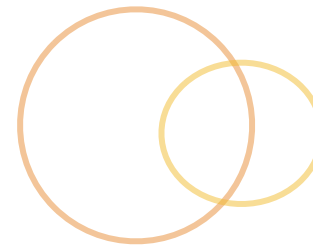


- 🕒 Change all `_text` fields to simply copy into `_en` fields
- 🕒 Explore a few of the other dynamic fields

Discussion

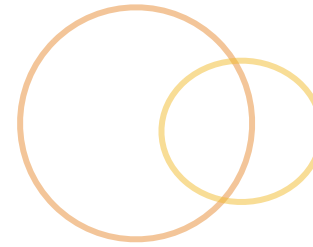
- ④ Pick a feature Content, Product, Reviews
- ④ What are major fields needed?
- ④ What types are the fields?
- ④ What analysis could be used?

Schema API



- ⦿ Allows for schema changes dynamically
- ⦿ Allows for adding, changing, deleting fields and field types
- ⦿ Caveat: Changing a field does not change data
- ⦿ Main entry point:
 - ⦿ `http://<host:port>/solr/<core>/schema`
- ⦿ Reading fields or a specific field:
 - ⦿ `/solr/<core>/schema/fields/`
 - ⦿ `/solr/<core>/schema/fields/<fieldname>`

Lab: Schema API

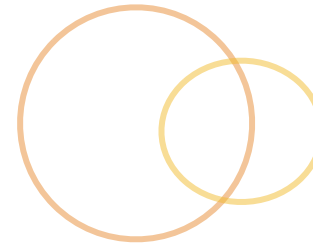


- ⦿ Retrieve the list of the defined fields for the index
- ⦿ Retrieve the Plot field specifically
- ⦿ Add a new field RTReview (int) for storing rotten tomatoes
- ⦿ Add a string field to store tags
- ⦿ Add a field type to search tags should trim, lowercase, and stem them
- ⦿ Create a copy directive to copy the first tags field into the new field type

SolrJ: Schema

- ⦿ Package Exists
- ⦿ No documentation page
- ⦿ Uses solr.request.schema
- ⦿ Use a series of request objects to change schema
- ⦿ Very similar to REST API from last section
- ⦿ Link in the git repo

SolrJ: List Fields



```
// Get all fields
```

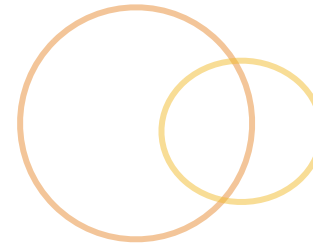
```
SchemaRequest.Fields fieldsRequest =
```

```
    new SchemaRequest.Fields();
```

```
NamedList fieldsResponse = solr.request(fieldsRequest);
```

```
List fieldsList = (List)fieldsResponse.get("fields");
```

SolrJ: Field Schema



// Getting a field

```
SchemaRequest.Field request =  
    new SchemaRequest.Field( fieldName: "Title");  
NamedList fieldResponse = solr.request(request);
```

// Don't be fooled this is not extending map

```
SimpleOrderedMap fieldMap =  
    (SimpleOrderedMap)fieldResponse.get("field");  
String name = (String)fieldMap.get("name");  
String type = (String)fieldMap.get("type");  
Boolean indexed = (Boolean)fieldMap.get("indexed");
```

SolrJ: Add Field



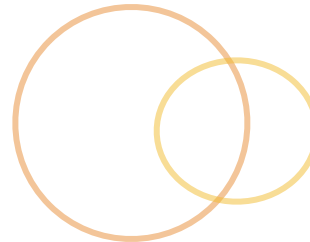
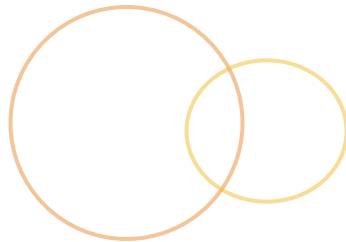
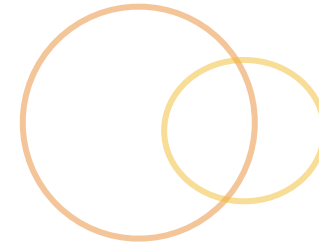
```
// Create the map for the values
```

```
Map<String, Object> valueMap = new HashMap<>();  
valueMap.put("name", "amazon");  
valueMap.put("type", "string");  
valueMap.put("stored", true);  
valueMap.put("indexed", false);  
SchemaRequest.AddField addFieldRequest =  
    new SchemaRequest.AddField(valueMap);  
//Returns the newly created field  
NamedList addResponse = solr.request(request);
```

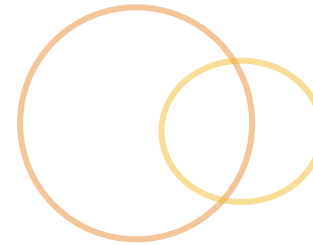
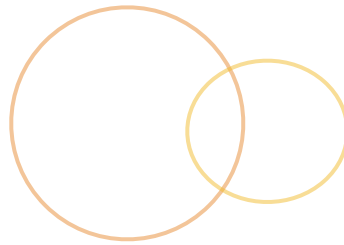
Schemaless

- ⦿ Solr can run in schemaless (guessing) mode
- ⦿ Data will determine the schema
- ⦿ New fields are added automatically

Part 3: Indexing

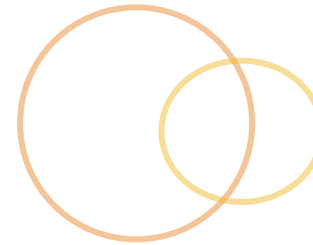
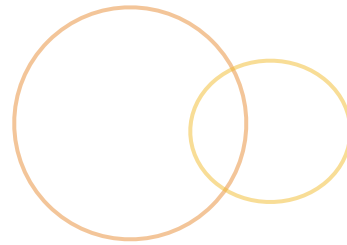


Indexing



- ② Process to add to the index
- ② Common ways:
 - ② `bin/post` file
 - ② Post request to update handler
 - ② Data Import Handler

Deleting

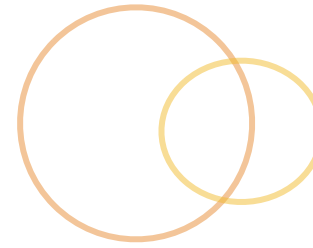


- ⦿ Leaves document in index hidden
- ⦿ Posting with body
 - ⦿ {"delete": { "id": "12345" }}
 - ⦿ {"delete": { "query": "Title:Toy Story" }}

Committing

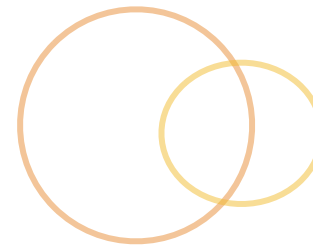
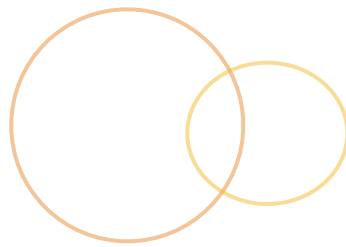
- ⦿ Adding records keeps them in staging
- ⦿ Commit makes the records update the index
- ⦿ Similar to git add vs git commit
- ⦿ bin/post commits automatically
- ⦿ Sending a commit:
 - ⦿ <http://localhost:8983/solr/labmovies/update?commit=true>

Commit: Soft vs Hard



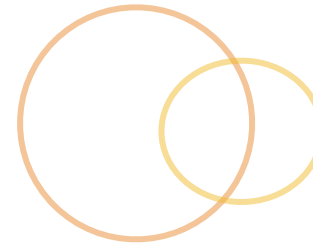
- ⦿ Hard commits write to disk
- ⦿ Soft commits:
 - ⦿ Makes changes visible quicker
 - ⦿ Changes not written
 - ⦿ Crashes can lose data

Optimizing



- ⦿ Optimization cleans up the index
- ⦿ Reduces amount of index files
- ⦿ Commit does not remove deleted docs
- ⦿ Optimize makes the index nice and clean
- ⦿ Can speed up search times

Lab: Load the Data



⦿ Moment of Truth

- ⦿ `bin/post -c labmovies <path>/movie_data_cleaned.json`
- ⦿ Windows: `jar` in `example/exampledocs`

⦿ Fix errors or ask questions

⦿ Delete your least favorite movie

⦿ Run an optimization if num docs and max docs are different

SolrJ: Adding Documents



```
/* ****  
 * Indexing  
 **** */  
SolrInputDocument document = new SolrInputDocument();  
document.addField(name: "Type", value: "movie");  
document.addField(name: "id", value: "tt1234566");  
document.addField(name: "genres",  
    Arrays.asList("horror", "thriller"));  
UpdateResponse addResponse = solr.add(document);
```


SolrJ: Commit/Optimize



//By default will block until hard commit is done
`UpdateResponse commit = solr.commit();`

//By default will block until complete
`UpdateResponse optimize = solr.optimize();`



Solr Primer Day 3



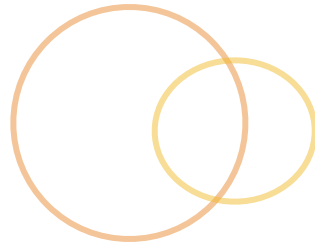
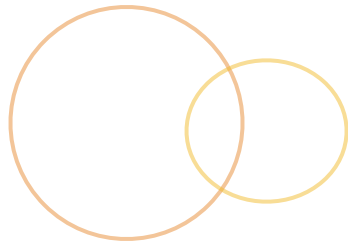
Macy's



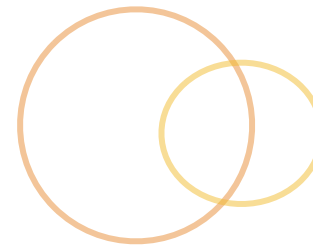
jeffnb@gmail.com



Part 4: Configurations



SolrConfig.xml



☉ Defines:

- ☉ All Handlers
- ☉ Commit Behavior
- ☉ Caching
- ☉ Default behaviors
- ☉ Search Components installed
- ☉ Schemaless configuration
- ☉ Default and environmental variables
- ☉ Listeners

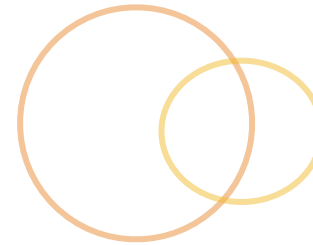
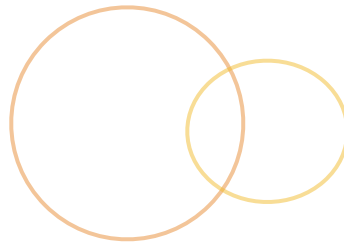
AutoCommit



- ⦿ Allows commits to be run automatically
- ⦿ Useful when trickle changes come in
- ⦿ Helps avoid many small commits

AutoSoftCommit

- ⦿ Commits will be visible
- ⦿ Commits not written to disk
- ⦿ Useful with AutoCommit
 - ⦿ AutoCommit at longer intervals
 - ⦿ AutoSoftCommit at short intervals



🕒 Query

- 🕒 Stores document ids for common queries
- 🕒 Cuts down on multiple index calls to find docs

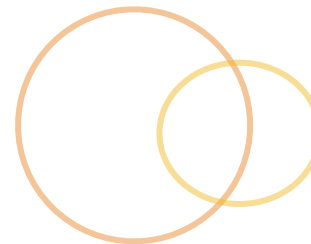
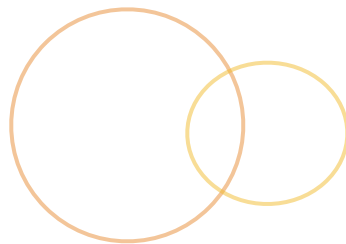
🕒 Filter

- 🕒 Stores document ids for common filter queries
- 🕒 Speeds up common filters considerably

🕒 Document

- 🕒 Stores commonly requested fields
- 🕒 Helps speed up document retrieval

Searchers



- ⦿ Similar to a thread
- ⦿ Searchers process the request
- ⦿ Searchers have their own caches
- ⦿ Have events:
 - ⦿ Trigger when the first searcher starts
 - ⦿ Trigger when new searchers start
 - ⦿ Useful for warming with common queries

Lab: Commits and Caches



- ② Set up auto commit in the movies solrconfig file
 - ② Commit every 5 minutes
 - ② Soft commit every minute
 - ② Ensure new searcher is opened on commits
- ② Experiment tuning the caches
 - ② Enable a filter cache which auto warms with 200 entries
 - ② Enable a query cache with 20 entries auto warmed
 - ② Look at the admin Plugins->Cache and checkout the cache levels. Try queries and watch them change
- ② Run 3 queries on startup when searcher is warmed

Search Handlers

- ⦿ End points for querying
- ⦿ Several come with solrconfig
- ⦿ Built in ones can be overridden
- ⦿ New ones can be created (edismax from earlier)
- ⦿ Specifies all behaviors
- ⦿ `class="solr.SearchHandler"`

Search Handlers

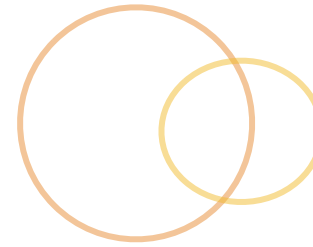
Defaults for queries

- ⦿ `<lst name="defaults">`
- ⦿ Specifies any query defaults
- ⦿ Useful for fl, rows and wt
- ⦿ Within a search handler

Appends

- ⦿ `<lst name="appends">`
- ⦿ Allows values to queries
- ⦿ User cannot override
- ⦿ Example: `<str name="fq">inStock:true</str>`

Search Handlers [cont]



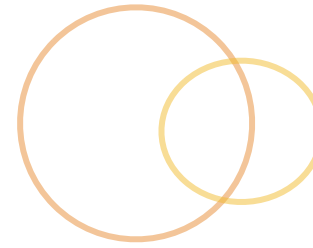
⦿ InitParams

- ⦿ Specifies a set of parameters across specified handlers
- ⦿ Overridden by individual handlers

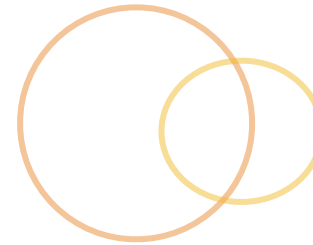
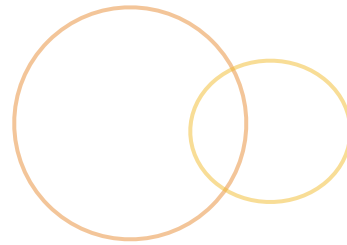
⦿ Invariants

- ⦿ `<lst name="invariants">`
- ⦿ Used to limit options on queries
- ⦿ Example: `<str name="facet.field">genres</str>`
- ⦿ Once specified no changes are allowed query time for that name

Lab: Search Handler



- 🕒 Create a new search handler for Action movies
- 🕒 Specify these defaults:
 - 🕒 fl with Title, Plot, Year, genres, imdbRating
 - 🕒 rows of 50
 - 🕒 wt of json
 - 🕒 turn on indent
- 🕒 Append genres_en:Action to all queries
- 🕒 Allow faceting on Year and actors



⦿ Advantages

- ⦿ Intended to not throw errors
 - ⦿ Syntax more forgiving
 - ⦿ Can search over many fields
 - ⦿ Useful for user searches
- ## ⦿ Extremely powerful
- ## ⦿ Flexible for different needs
- ## ⦿ Run with: `defType=edismax`

eDismax parameters

- q - Term to search for

- qf

- Fields to use to search with boosts

```
<str name="qf">  
  Title^10 Plot^5  
</str>
```

- mm - Minimum amount of terms to match

- pf

- Boosts based on fields whole phrases appear in.

- Same syntax as qf

- bq

- Boosts based on a query

- bq=genres:action^3.0

Lab: eDismax Handler

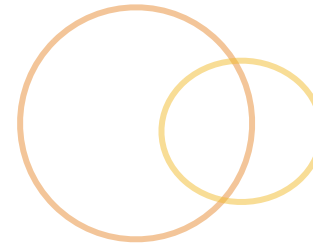
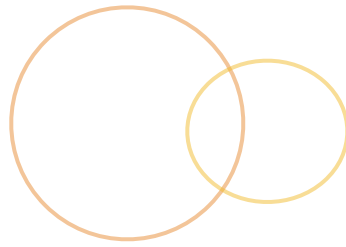


- ② Create a new search handler “edismax”
- ② For search fields:
 - ② Title and Plot are the highest boosted
 - ② actors, directors and writers have a lower boost
 - ② Year and Language are searched but not boosted
- ② All terms should match
- ② Boost (pf) Title and plot again
- ② Rows should be 30
- ② Fields returned should be Title, Plot and Year

Discussion

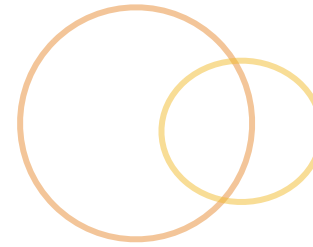
- ⦿ Now that you have a grasp on query options
- ⦿ What fields do you think are most important?
- ⦿ What would you want to return by default?

Config API



- ⦿ Rest API exists
- ⦿ Doesn't appear to be supported in SolrJ
- ⦿ Can administer the config through the api
- ⦿ Stand alone mode would only replicate with a master server

Config API: Reading



- Entire Config

- <http://localhost:8983/solr/movies/config/>

- Specific Section

- <http://localhost:8983/solr/movies/config/requestHandler>

- Specific Component

- <http://localhost:8983/solr/movies/config/requestHandler?componentName=/select>

Config API: Common Commands



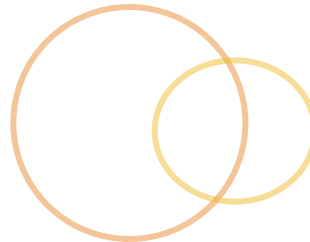
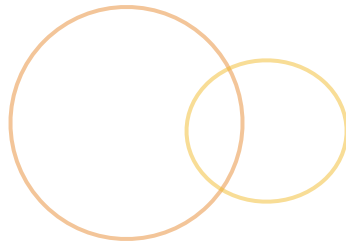
- ⦿ Updating caches
- ⦿ Changing auto commit values
- ⦿ JMX changes

Config API: General Purpose

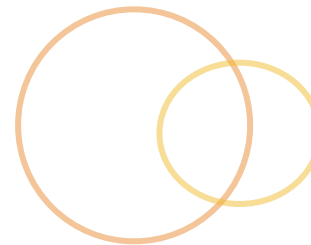
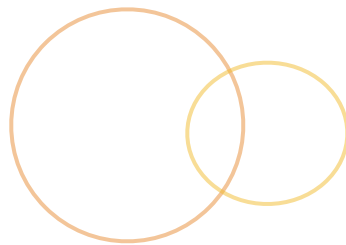


- 🕒 Request Handlers
- 🕒 Search Components
- 🕒 Init params

Part 4: Optimization

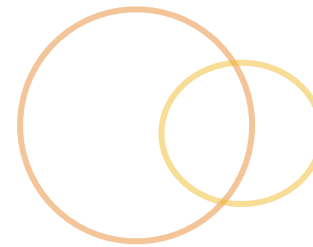


Overview



- ⦿ Solr is big
- ⦿ No silver bullet
- ⦿ Should be based on user behavior
- ⦿ Dependent on your architecture

System Level



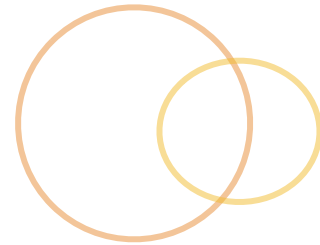
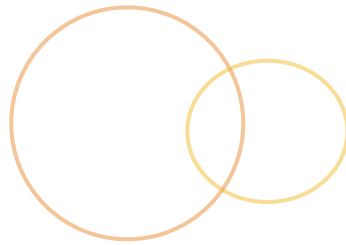
☉ Memory

- ☉ Do not give Solr all the memory
- ☉ At least half should be left for the OS
- ☉ Turn up OS file cache to get index in memory
- ☉ Make sure open file limit is high enough

☉ Hardware

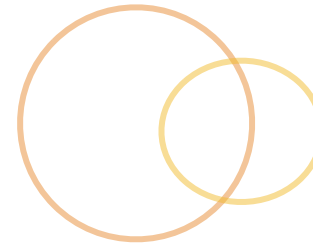
- ☉ Master/Slave servers work much better on bare metal
- ☉ SSDs are great speed improvements if possible
- ☉ Ram can boost performance

Caching



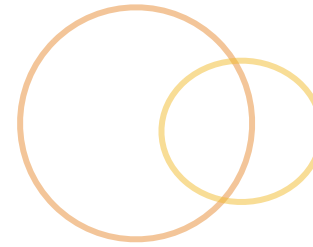
- 🕒 Look at the logs
- 🕒 Warm searchers with common queries
- 🕒 Autowarm when you can
- 🕒 Start with defaults and increase or decrease levels
- 🕒 Balance memory with caches
- 🕒 Use the cache admin to get hit rates

Distributed Indexes



- ⦿ Indexes slow with size
- ⦿ Threshold highly variable (millions of records)
- ⦿ Sharding: breaking up index
- ⦿ Solr HIGHLY recommends not sharding with M/S setup
- ⦿ If index becomes too large move to SolrCloud

Document Adding



☉ Change Update Format

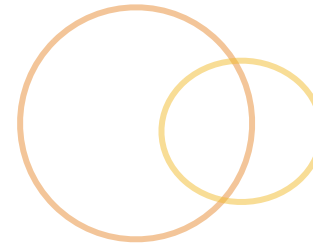
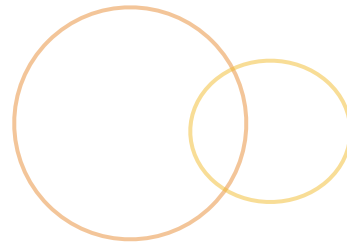
- ☉ Default XML
- ☉ Change to binary
- ☉ Only works after 3.1

☉ Concurrent updates.

- ☉ Uses multiple HTTP connections
- ☉ Speeds up large indexing

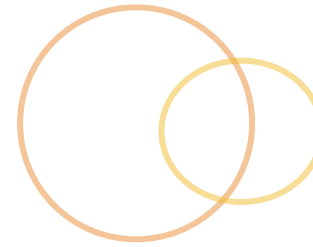
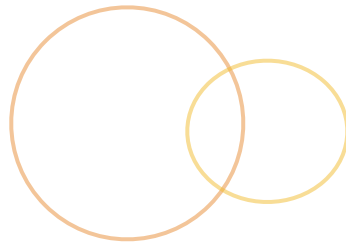
```
ConcurrentUpdateSolrClient conSolr =  
    new ConcurrentUpdateSolrClient.Builder(urlString).build();  
conSolr.setRequestWriter(new BinaryRequestWriter());
```

Commits



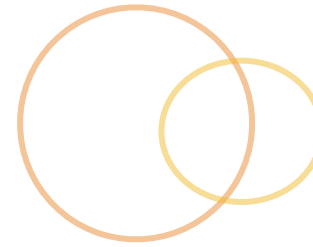
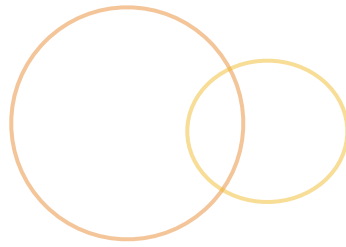
- ⦿ Commits slow throughput
- ⦿ Searches are expensive to open
- ⦿ Be practical about real time data needs
- ⦿ Use soft commits to supplement hard commits
- ⦿ Transaction logs will grow out of control without hard commits

Optimizes



- ⦿ Slow throughput down
- ⦿ Optimizing creates new files
- ⦿ Optimizing merges files making index smaller
- ⦿ Overall better performance when optimized
- ⦿ Better run at low times

Replication



- ⦿ Replication can cause new searchers
- ⦿ Will temporarily slow searches down
- ⦿ Understand real time needs

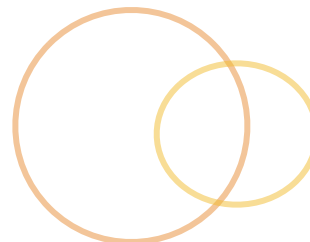
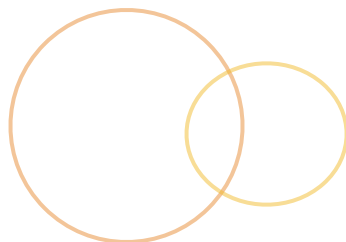
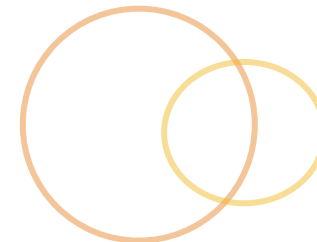
Garbage Collection

- ⦿ Larger Solr instances
- ⦿ Can stop the entire server for seconds
- ⦿ Older versions of GC “stop the world”
- ⦿ CMS better for 1.6
- ⦿ G1 better for 1.7+

Discussion

- What are some of the performance struggles?
- Are you using the optimizations discussed?

Part 5: SolrCloud



SolrCloud Overview

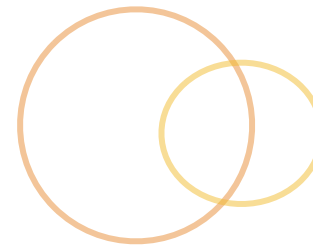
- ☉ Grew out of scalability troubles
- ☉ Closest to ElasticSearch methods
- ☉ Handles replication and sharding
- ☉ Uses ZooKeeper to manage nodes
- ☉ Allows new servers to plug in with ease
- ☉ Has collections not cores

SolrCloud vs. Elasticsearch



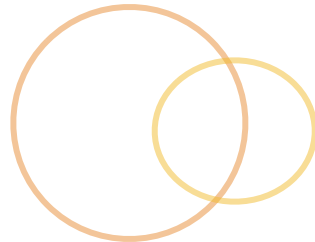
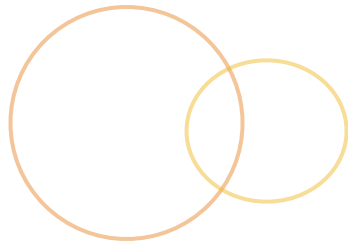
- ⦿ Almost on the same footing
- ⦿ Solr requires an extra server
- ⦿ Elasticsearch is poorly documented
- ⦿ Elasticsearch easier out of the box
- ⦿ Elasticsearch uses JSON which can be frustrating for java
- ⦿ Solr has native Java support

Lab: SolrCloud

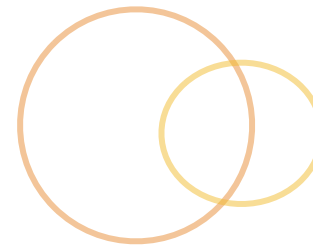


- ② Use bin/solr to create a movies collection
- ② Use 2 shards and 2 replicas
- ② Import movies data
- ② Explore diagrams in admin interface

Part 6: NLP and Solr



What is NLP



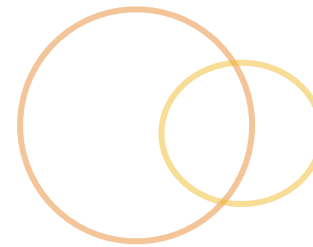
🕒 Natural Language Processing

Natural language processing (NLP) is a field of **computer science**, **artificial intelligence**, and **computational linguistics** concerned with the interactions between **computers** and **human (natural) languages** and, in particular, concerned with programming computers to fruitfully process large **natural language corpora**.

What is NLP

- ④ A way to help computers understand human language
- ④ Instead of looking at each word individually (solr)
- ④ Takes the context of a sentence or document

Human Language



- ⦿ Extremely difficult
- ⦿ New words constantly that don't have any context in existing models
- ⦿ Multiple meanings are a problem: Spoke (speaking) vs Spoke (in a bike wheel)
- ⦿ Modifying words can change the feel of a sentence
“It was a hardly fun” vs “It was incredibly fun”
- ⦿ Only 30% of communication is based on the words said

NLP: Common Tasks

- ② Sentence Breaking
- ② Tokenization
- ② Tagging
- ② Lemmatization/Stemming
- ② Name Entity Extraction

NLP: Sentence Breaking



- ⦿ Splitting input on sentence boundaries
- ⦿ Can be tricky with punctuation ... vs .
- ⦿ Important first step
- ⦿ NLP is all about context
- ⦿ Sentences are a discrete idea.

NLP: Tokenization

- ⦿ Similar to Solr ideas of tokenization
- ⦿ Can take many forms
- ⦿ Usually special versions of white space
- ⦿ Keeps punctuation in the list of tokens

NLP: Tagging



- Processes sentence and tokens
- Tags each token with a type of speech
- Spoke => Spoke/VBD
- Running => Running/VBG

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential 'there'	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VCN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>ilama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>ilamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	"	left quote	<i>' or "</i>
POS	possessive ending	<i>'s</i>	"	right quote	<i>' or "</i>
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	<i>[, (, {, <</i>
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis	<i>],), }, ></i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... - -</i>
RP	particle	<i>up, off</i>			

NLP: Lemmatization

- Similar to stemmers
- Uses tags to determine roots
- Root is always a full word
- Example:
 - Speaking/VBG and Spoke/VB => speak

NLP: Named Entity Extraction

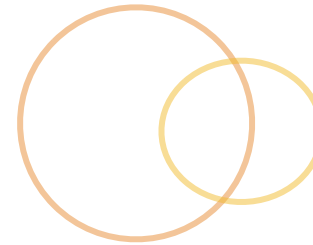
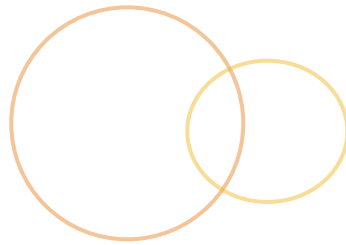


- ⦿ Pulls out names from input
- ⦿ Can be made to identify people vs places
- ⦿ Interesting usage of process

Machine Learning

- ④ How natural language processing is based
- ④ Take training data and run it through model
- ④ Model “learns” from the data to associate certain desired outcomes
- ④ Testing data then run through system to ensure outcome correct

OpenNLP



- ② Apache Project
- ② Handles common Tasks
- ② Has many different models
- ② Models are downloaded separately
- ② Stand alone project
- ② Java available with Maven include

OpenNLP: Pros and Cons



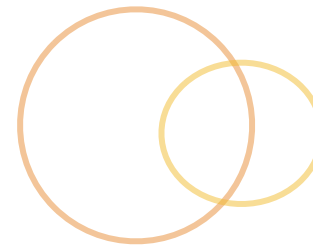
⦿ Pros:

- ⦿ Open source
- ⦿ Rudimentary connection to Solr
- ⦿ Flexible with including models

⦿ Cons:

- ⦿ Small community
- ⦿ Not particularly active
- ⦿ Models aren't included due to licensing

OpenNLP and Solr

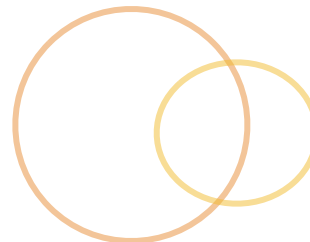
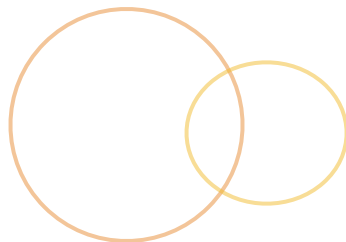
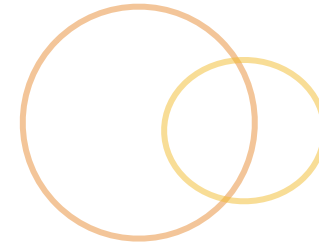


- ❉ No standard integration
- ❉ Famous Lucene-2899 Patch
 - ❉ Written 6 years ago
 - ❉ Updated last year
 - ❉ Missing any way to process named entities
 - ❉ Memory Inefficient
- ❉ Download libs and include in Solr package
- ❉ Allows for common steps mentioned

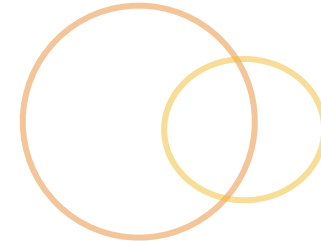
NLP: Considerations

- OpenNLP is not the biggest or best
 - NTK - Python
 - CoreNLP - Stanford
- NLP is good for language not searches
 - Lemmas and Tagging break down in keyword searching
 - “Search Salad” is very difficult to understand
- Lemmatization struggles with new words
- Just a start. Extra work needed to effectively use

Final Discussions



Thank you



Please fill out the surveys

