
Solr Team

CS5604: Cloudera Search in IDEAL

Nikhil Komawar, Ananya Choudhury, Rich Gruss

Tuesday May 5, 2015

Department of Computer Science
Virginia Tech, Blacksburg

Outline

1. Schema design
 2. Indexing
 3. Custom Search
-

System Overview

1. CPU
 - a. Intel i5 Haswell Quad core 3.3 Ghz Xeon
 2. RAM
 - a. 660 GB in total
 - b. 32 GB in each of the 19 Hadoop nodes
 - c. 4 GB in the manager node
 - d. 16 GB in the tweet DB nodes
 - e. 16 GB in the HDFS backup node
 3. Storage
 - a. 60 TB across Hadoop, manager, and tweet DB nodes
 - b. 11.3 TB for backup
 4. Number of nodes
 - a. 19 Hadoop nodes
 - b. 1 Manager node
 - c. 2 Tweet DB nodes
 - d. 1 HDFS backup node
-

Solr Schema: Design

```
<!-- twitter |fields -->
<field name="id" type="string" indexed="true" multiValued="false"/>
<field name="tweet_id" type="long" indexed="true" multiValued="false"/>
<field name="collection" type="text_general" indexed="true" multiValued="false"/>
<field name="text" type="text_general" indexed="true" multiValued="false"/>
<field name="created_at" type="tdate" indexed="true" multiValued="false"/>
<field name="source" type="string" indexed="true" multiValued="false"/>
<field name="user_screen_name" type="string" indexed="true" multiValued="false"/>
<field name="user_id" type="string" indexed="true" multiValued="false"/>
<field name="lang" type="string" indexed="true" multiValued="false"/>
<field name="retweet_count" type="tint" indexed="true" multiValued="false"/>
<field name="favorite_count" type="tint" indexed="true" multiValued="false"/>
<field name="_version_" type="tlong" indexed="true" multiValued="false"/>
<field name="contributors_id" type="string" indexed="true" multiValued="false"/>
<field name="coordinates" type="string" indexed="true" multiValued="false"/>
<field name="urls" type="string" indexed="true" multiValued="false"/>
<field name="hashtags" type="string" indexed="true" multiValued="false"/>
<field name="user_mentions_id" type="string" indexed="true" multiValued="false"/>
<field name="in_reply_to_user_id" type="string" indexed="true" multiValued="false"/>
<field name="in_reply_to_status_id" type="string" indexed="true" multiValued="false"/>
```

```
<fields>
  <!-- webpages |fields -->
  <field name="id" type="string" indexed="true" multiValued="false"/>
  <field name="collection" type="string" indexed="true" multiValued="false"/>
  <field name="title" type="text_general" indexed="true" multiValued="false"/>
  <field name="domain" type="string" indexed="true" multiValued="false"/>
  <field name="url" type="string" indexed="true" multiValued="false"/>
  <field name="text" type="text_general" indexed="true" multiValued="false"/>
```

```
<!-- class analysis fields -->
<field name="ner_people" type="string" indexed="true" multiValued="false"/>
<field name="ner_locations" type="string" indexed="true" multiValued="false"/>
<field name="ner_dates" type="string" indexed="true" multiValued="false"/>
<field name="ner_organizations" type="string" indexed="true" multiValued="false"/>
<field name="cluster_id" type="string" indexed="true" multiValued="false"/>
<field name="cluster_label" type="string" indexed="true" multiValued="false"/>
<field name="classification_vector_json" type="string" indexed="true" multiValued="false"/>
<field name="classification_labels" type="string" indexed="true" multiValued="true"/>
<field name="social_vector_json" type="string" indexed="true" multiValued="false"/>
<field name="social_importance" type="double" indexed="true" multiValued="false"/>
<field name="lda_dict_json" type="string" indexed="true" multiValued="false"/>
<field name="lda_topics" type="string" indexed="true" multiValued="false"/>
<field name="urls_multiple" type="string" indexed="true" multiValued="false"/>
<field name="hashtags_multiple" type="string" indexed="true" multiValued="false"/>
<field name="ner_people_multiple" type="string" indexed="true" multiValued="false"/>
<field name="ner_locations_multiple" type="string" indexed="true" multiValued="false"/>
<field name="ner_dates_multiple" type="string" indexed="true" multiValued="false"/>
<field name="ner_organizations_multiple" type="string" indexed="true" multiValued="false"/>
<field name="lda_topics_multiple" type="string" indexed="true" multiValued="false"/>
<field name="classification_labels_multiple" type="string" indexed="true" multiValued="false"/>
```

Solr Schema: Effect

Solr idx size depends on:

- Number of fields
 - Stored vs. not stored
 - Type of field
 - example: string vs. text
 - Index issues:
 - Recommended to add H/W
 - Alternatively, design schema and tune Java/Solr configurations
-

Solr Schema: Future

- Fewer stored fields
 - NRT indexer
 - Consider using facet.method=fcs, helps during first request
 - Index size lowering
 - $D \times S + U$
 - D is the document count
 - S is the the size of the data type
 - ints - 4bytes, 8 bytes for doubles, U cumulative size of the unique field values
-

Document ID: Design

```
hbase(main):002:0> scan 'tweets', {LIMIT => 1}
COLUMN+CELL
column=analysis:class, timestamp=1430618259668, value=POSITIVE
column=analysis:cluster_id, timestamp=1430016647887, value=Jan.2
column=analysis:cluster_label, timestamp=1430016647887, value=
column=original:collection, timestamp=1428850721220, value=Jan.2
column=original:coordinates, timestamp=1428850721220, value=0.0,
column=original:created_at, timestamp=1428850721220, value=13490
column=original:hashtags, timestamp=1428850721220, value=Egypti
column=original:lang, timestamp=1428850721220, value=English
column=original:source, timestamp=1428850721220, value=twitter-t
column=original:text_clean, timestamp=1430618259668, value= E
[1](status)
s/skip-link-focus-fix.js?ver=20130115'></script>
er=4.1.1'></script>\x0A\x0A</body>\x0A</html>\x0
column=original:title, timestamp=1430091885753,
column=original:url, timestamp=1430091885753, va
1 row(s) in 0.1650 seconds
100 row(s) in 2.9180 seconds
hbase(main):003:0>
hbase(main):009:0>
```

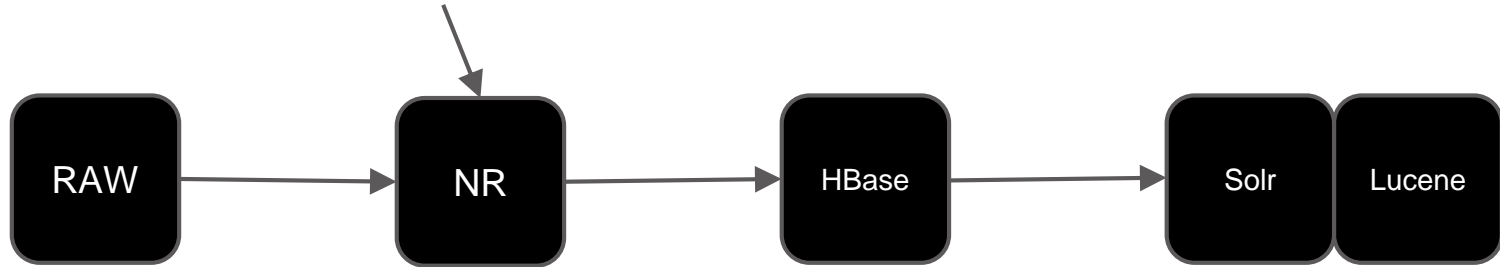
Syntax: Collection_name--Counter_value
Use: Noise, Reduction, HBase, Solr/Lucene

Document ID: Effect

- Affects Lucene indexing
 - Does not affect Solr index size
 - Fastest: Zero padded sequential
 - Slowest: Random UUID generated using some languages (UUID v4)
 - Our current design
 - Performance tradeoffs
 - Somewhere in the middle
-

ID generation in pipeline

ID addition



Document ID: Future

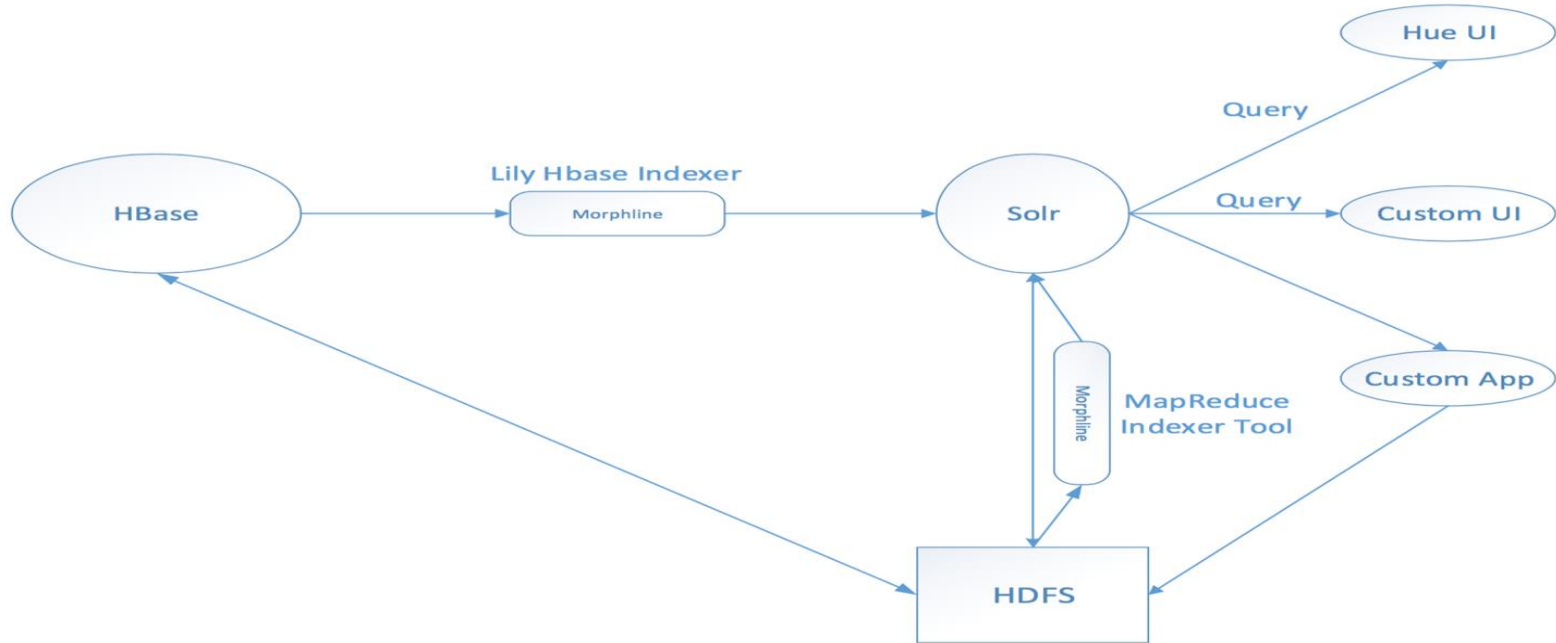
Preprocessing

- Current:
 - Concurrency
 - `test_and_set`
 - Batch processing
- Recommendation:
 - Binary encoded UUID
 - Parallel
 - CPU hit

Querying

- Current:
 - Faster fetching
 - Lower disk I/O
- Recommendation:
 - Sequentially assigned value or unhashed timestamp
 - Batch processing vs. asynchronous processing

Indexing



Hbase Indexer - Morphline

```
morphlines: [  
  {  
    commands: [  
      {  
        extractHBaseCells {  
          mappings: [{  
            inputColumn: "original:text_clean"  
            outputField: "text"  
            type: string  
            source: value  
          }  
          {  
            inputColumn: "analysis:ner_people"  
            outputField: "ner_people_multiple"  
            type: string  
            source: value  
          }  
        ]  
      }  
      {  
        split {  
          inputField: "ner_people_multiple"  
          outputField: "ner_people"  
          separator: "|"  
        }  
      }  
    }  
  }  
]
```

maps to HBase
column family and
column qualifier

defined in
schema.xml

Indexing - Jobs

hue Home Query Editors v Data Browsers v Search v File Browser Job Browser cs5604s15_solr ?

Job Browser

Username Text

Succeeded Running Failed Killed

Logs	ID	Name	Status	User	Maps	Reduces	Queue	Priority	Duration	Submi
	1427130979177_12390	org.apache.solr.hadoop.ForkedMapReduceIndexerTool/ForkedTreeMergeMapper	SUCCEEDED	cs5604s15_solr	100%	100%	root.cs5604s15_solr	N/A	2h:26m:46s	05/01/15 23:30:28
	1427130979177_12346	HBaseMapReduceIndexerTool/HBaseIndexerMapper	SUCCEEDED	cs5604s15_solr	100%	100%	root.cs5604s15_solr	N/A	26m:41s	05/01/15 23:03:41
	1427130979177_5798	org.apache.solr.hadoop.ForkedMapReduceIndexerTool/ForkedTreeMergeMapper	SUCCEEDED	cs5604s15_solr	100%	100%	root.cs5604s15_solr	N/A	9m:9s	04/30/15 13:52:20
	1427130979177_5775	HBaseMapReduceIndexerTool/HBaseIndexerMapper	SUCCEEDED	cs5604s15_solr	100%	100%	root.cs5604s15_solr	N/A	10m:2s	04/30/15 13:42:12

Results



Dashboard

Logging

Cloud

Core Admin

Java Properties

Thread Dump

webpages_shar...

Overview

Analysis

Config

Dataimport

Documents

Ping

Statistics

Last Modified: 4 days ago

Num Docs: 12326

Max Doc: 12326

Deleted Docs: 0

Version: 9

Segment Count: 1

Optimized:

Current:

Replication (Master)

	Version	Gen	Size
Master (Searching)	1430412557297	5	27.04 MB
Master (Replicable)	-	-	-

Admin Extra

Webpages



Dashboard

Logging

Cloud

Core Admin

Java Properties

Thread Dump

tweets_shard1_...

Overview

Analysis

Config

Dataimport

Documents

Ping

Statistics

Last Modified: about 16 hours ago

Num Docs: 3444305

Max Doc: 3444305

Deleted Docs: 0

Version: 3

Segment Count: 1

Optimized:

Current:

Replication (Master)

	Version	Gen	Size
Master (Searching)	1430730434200	2	764.78 MB
Master (Replicable)	-	-	-

Admin Extra

Tweets

Search

1: Adjust boost levels in Query Parser (solrconfig.xml)

```
<requestHandler name="/select" class="solr.SearchHandler">
  <lst name="defaults">
    <str name="defType">edismax</str>
    <str name="qf">
      text
      collection^3
      hashtags^3
      cluster_label^2.5
      lda_topics^2.0
      ner_people^2.0
      ner_locations^2.0
      ner_organizations^2.0
    </str>
    <str name="echoParams">explicit</str>
    <int name="rows">10</int>
    <str name="df">text</str>
    <str name="fl">*,score</str>
  </lst>
  <arr name="last-components">
    <str>idealSocialBoostComponent</str>
    <str>idealTopicSupplementComponent</str>
  </arr>
</requestHandler>
```

```
<searchComponent name="idealSocialBoostComponent" class="edu.vt.dlib.ideal.solr.IDEALSocialBoostComponent"/>
<searchComponent name="idealTopicSupplementComponent" class="edu.vt.dlib.ideal.solr.IDEALTopicSupplementComponent"/>
```

These numbers are **conjectural**. A more disciplined approach:

$$\text{Relevance} = f(q,d)$$

For each of 100 queries, induce a logistic regression model. Find average contribution of each field to relevance.

Search



Dashboard

Logging

Cloud

Core Admin

Java Properties

Thread Dump

tweets_shard1_...

Overview

Analysis

Config

DataImport

Documents

Raw Query Parameters

key1=val1&key2=val2

wt

json

indent

debugQuery

dismax

edismax

q.alt

qf

mm

pf

ps



Stephanie Sparkles @stripesnsparkle · Feb 1

Winter Storm Coming! #Chicago #snow #winter #winterstorm #cold #outside
#white #babyitscoldoutside #lights... fb.me/3Uwz74a9p



```
{
  "text": "Winter Storm Coming storm",
  "user_screen_name": "stripesnsparkle",
  "urls": [
    "http://t.co/zuzA8x1bhe"
  ],
  "tweet_id": 3834306245955107000,
  "cluster_id": "winter_storm_S--c5",
  "collection": "winter_storm_S",
  "lang": "English",
  "id": "winter_storm_S--299751",
  "source": "twitter-search",
  "created_at": "2015-02-01T06:07:32Z",
  "cluster_label": "storm",
  "user_id": "2914255473",
  "hashtags": [
    "Chicago",
    "snow",
    "winter",
    "winterstorm",
    "cold",
    "outside",
    "white",
    "babyitscoldoutside"
  ],
  "coordinates": [
    "0.0,0.0"
  ],
  "_version_": 150022883354922000
},
```


Search

2: Reorder documents on the basis of Social Importance Score.

solr.SearchHandler

```
for( SearchComponent c : components ) {
    c.prepare(rb);
}

...

for( SearchComponent c : components ) {
    c.process(rb);
}
```

edu.vt.dlib.ideal.solr.SocialBoostComponent

```
DocIterator iterator = docList.iterator();
while (iterator.hasNext()) {

    int docId = iterator.nextDoc();

    float score = iterator.score();

    Document d = rb.req.getSearcher().doc(docId);
    float socialBoost = 0;
    IndexableField socialImportanceField = d.getField("social_importance");
    if(socialImportanceField != null) {
        Number socialImportance = socialImportanceField.numericValue();
        if (socialImportance != null) {
            socialBoost = socialImportance.floatValue();
        }
        score = score + socialBoost;
    }
    scoreDocs[idx++] = new ScoreDoc(docId, score);
}
```

Search

3: Supplement result list, if necessary, by retrieving documents from collections that include topics related to the search. Get topic list straight from HBase.



```
@Override
public void process(ResponseBuilder rb) throws IOException {

    DocListAndSet results = rb.getResults();

    SortSpec sortSpec = rb.getSortSpec();
    int len = sortSpec.getCount();
    int offset = sortSpec.getOffset();

    DocList docList = rb.getResults().docList;
    if (docList.size() < MIN_DOCS) {
        String collectionName = getDominantCollection(rb);
        TopicModel topicModel = getCollectionTopicModel(collectionName);
        String topTopics = topicModel.getTopTopics();

        QueryParser queryParser = new QueryParser("name", new StandardAnalyzer());
        Query query = null;
        try {
            query = queryParser.parse("q=" + topTopics);
        } catch (ParseException e) {
            e.printStackTrace();
        }

        TopDocs additionalDocs = rb.req.getSearcher().search(query, 100);
        ScoreDoc[] scoreDocs = additionalDocs.scoreDocs;
        int totalHits = additionalDocs.totalHits + docList.matches();

        int[] docs = new int[scoreDocs.length];
        float[] scores = new float[scoreDocs.length];
    }
}
```

```
private TopicModel getCollectionTopicModel(String collectionName) throws IOException {

    TopicModel topicModel = new TopicModel(collectionName);
    HTable table = new HTable(conf, "collection_metadata");
    Get get = new Get(Bytes.toBytes(collectionName));
    byte[] val = table.get(get).getValue(b("analysis"), b("lda"));
    String jsonTopicModel = new String(val, StandardCharsets.UTF_8);

    Map<String, Double> map = new Gson().fromJson(jsonTopicModel, new TypeToken<HashMap<String, Double>>() {}.getType());
    topicModel.setTopicProbabilities(map);

    return topicModel;
}
```

Search Results

(first 1000 results)

```
{'query': 'election', 'num_results': 637498, 'precision': 0.998, 'time': 0.05295705795288086}
{'query': 'elect', 'num_results': 3247, 'precision': 0.978, 'time': 0.04558682441711426}
{'query': 'revolution', 'num_results': 13048, 'precision': 0.95, 'time': 0.04502081871032715}
{'query': 'uprising', 'num_results': 1769, 'precision': 0.851, 'time': 0.04298877716064453}
{'query': 'storm', 'num_results': 429329, 'precision': 0.999, 'time': 0.04975700378417969}
{'query': 'winter', 'num_results': 409987, 'precision': 0.999, 'time': 0.04920697212219238}
{'query': 'ebola', 'num_results': 306827, 'precision': 1.0, 'time': 0.04514813423156738}
{'query': 'disease', 'num_results': 6802, 'precision': 0.993, 'time': 0.041940927505493164}
{'query': 'bomb', 'num_results': 33924, 'precision': 0.857, 'time': 0.040463924407958984}
{'query': 'explosion', 'num_results': 1224, 'precision': 0.284, 'time': 0.04803609848022461}
{'query': 'crash', 'num_results': 274014, 'precision': 0.995, 'time': 0.04688715934753418}
{'query': 'plane crash', 'num_results': 193046, 'precision': 1.0, 'time': 0.056591033935546875}
{'query': 'shooting', 'num_results': 5366, 'precision': 0.744, 'time': 0.04262495040893555}
{'query': 'paris shooting', 'num_results': 446, 'precision': 0.446, 'time': 0.20793604850769043}
{'query': 'terrorist attack', 'num_results': 1143, 'precision': 0.768, 'time': 0.042675018310546875}
```

Future Work

1. Relevance feedback to derive logistic regression, evaluate with F1 score.
 2. More Solr nodes -> more index shards.
 3. Boost by length
 4. Innovative inputs (social graph)
 5. More sophisticated traversal of hierarchical classification
-

Acknowledgement

We are especially thankful to

- The NSF grant IIS - 1319578, III: Small: Integrated Digital Event Archiving and Library (IDEAL) for the funding that supported the infrastructure and the data used in the project.
 - Dr. Edward A. Fox for being the guide on the side for us and for coordinating efforts of all the teams to help us make it a successful class project.
 - The GTA (Sunshin), the GRA (Mohamed) and other students of the class who supported us with ideas as well as efforts during the semester.
 - The authors and the contributors of open source projects, blogs and wiki pages from where we borrowed some ideas and solutions.
-

Thank you!

Q & A
