# CS 5604 Information Storage and Retrieval

# Solr Team Final Presentation

## Presenters:

Liuqing Li,  Ye Wang,  Anusha Pillai,  Ke Tian

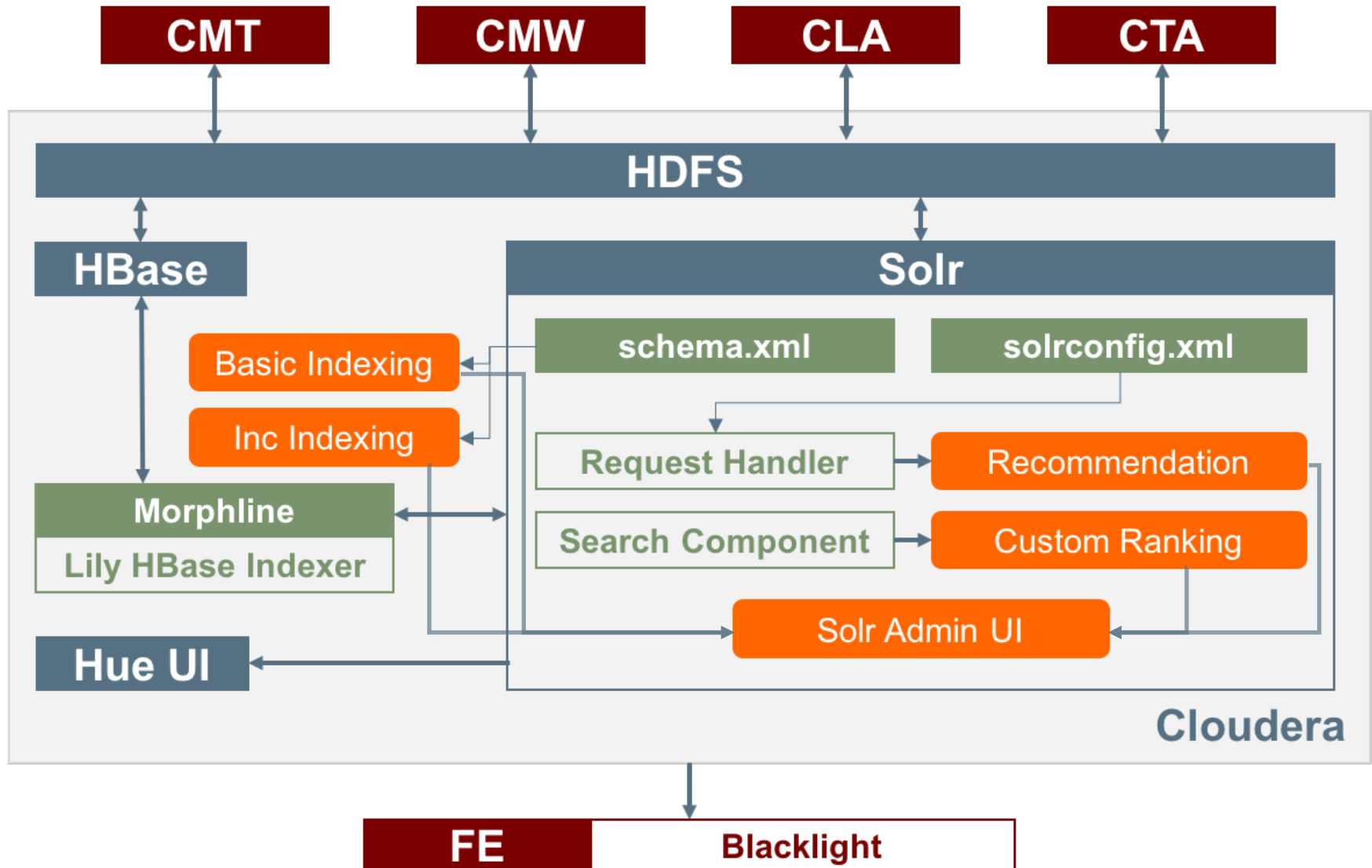{liuqing, yewang16, anusha89, ketian} @vt.edu

## Instructor:

Dr. Edward A. Fox

Virginia Polytechnic Institute and State University

Blacksburg, VA, 24061

December 6, 2016

**VirginiaTech**

*Invent the Future*

# Outline

- Background

- Implementation

- Problems Faced

- Lessons Learned

- Future Work

- Acknowledgement

**Solr Team Final Presentation**

VirginiaTech
*Invent the Future*

# Background — Overview

**Solr Team Final Presentation**

VirginiaTech
*Invent the Future*

# Background — Updates

| | **Spring 2016** | **Fall 2016** |
|---|---|---|
| **schema.xml** | Coarse grained | Fine grained |
| | No copyfields | Copyfields for all fields search |
| | Create stopwords.txt & profanity.txt | Update the two files |
| **morphlines.conf** | Two field types: string and text | Multiple field types |
| | Field "time" => string | Field "time" => datetime |
| | No multiple-valued fields | Multiple-valued field parser |
| **Basic Indexing** | Small collection | 1.2 billion tweets dataset |
| **Incremental Indexing** | Virtual Cloudera (VC) | VC & Hadoop Cluster (HC) |
| **Recommendation** | Brief description | Implemented in VC & HC |
| **Custom Ranking** | Brief description | Implemented in VC & HC |
| **Solr Admin UI** | Brief description | Detailed description |
| | Limited faceted search | Detailed faceted search |

**Solr Team Final Presentation**

VirginiaTech
1872
Invent the Future

# Implementation — Basic Indexing

- Live Mode
  - Continuous stream of HBase cell updates into live search indexers
  - Simple and efficient
  - Cannot handle big data

- Batch Mode
  - Batch index tables in HBase by using MapReduce jobs
  - Write index files into HDFS (/user/cs5604f16_solr/...)
  - Can handle big data

VirginiaTech
*Invent the Future*

- schema.xml: fields configuration
  - field (e.g., ideal-cs5604f16-fake)
    - # of fields: 30
    - Types: string (22), text_general (2), int (2), float (2), long (1), date (1)
    - Stored: True (17), False (13)

```xml
<field name="t_month_i" type="int" indexed="true" stored="true"/>
<field name="hashtags_s" type="string" indexed="true" stored="false" multiValued="true"/>
```

  - dynamicField: matching multiple fields, using wildcard

```xml
<dynamicField name="*_s"  type="string"  indexed="true"  stored="true" />
<dynamicField name="*_ss" type="string"  indexed="true"  stored="true" multiValued="true"/>
```

  - copyField

```xml
<copyField source="*_ss" dest="text" maxChars="3000"/>
```

5

VirginiaTech
1872
*Invent the Future*

- stopword.txt and profanity.txt
  - stopword.txt: tf-idf value will not be calculated
  - profanity.txt: quick response for such search queries
  - Solr loads the two files while reading schema.xml

```
<!-- Case insensitive stop word removal.
    -->
<filter class="solr.StopFilterFactory"
        ignoreCase="true"
        words="lang/stopwords_en.txt"
        />
<filter class="solr.LowerCaseFilterFactory"/>
<filter class="solr.EnglishPossessiveFilterFactory"/>
<filter class="solr.KeywordMarkerFilterFactory" protected="protwords.txt"/>
```

**Source:**
https://pypi.python.org/pypi/many-stop-words
http://www.freewebheaders.com/full-list-of-bad-words-banned-by-google/

6

**Solr Team Final Presentation**

**VirginiaTech**
*Invent the Future*

- morphlines.conf: mapping and parsing

**Mapping data from HBase to Solr**

```
mappings: [
# tweet : cleantext
{
        inputColumn: "tweet:cleantext"
        outputField: "raw_cleantext_s"
        type: string
        source: value

}
```

**Split multiple values into list**

```
split {
        inputField : "topic_label_s"
        outputField : "topic_label_ss"
        separator : ";"
        isRegex : false
        addEmptyStrings : false
        trim : true
}
```

```
"topic_label_s":
"twitter;social;media;text"
```

```
"topic_label_ss": [
  "twitter",
  "social",
  "media",
  "text"
],
```

7

**Solr Team Final Presentation**

# Implementation — Basic Indexing

- Index the big dataset

| | | ideal-cs5604f16 | ideal-cs5604f16-1204 |
|---|---|---|---|
| **Dataset** | | All collections (raw tweets) | All collections (raw tweets + processed data) |
| **Indexing** | **# of DataNode** | 18 | 17 |
| | **Space Cost** | 392.33 GB | 399.21 GB |
| | **Time Cost** | | |
| | **Mapping** | 1h21m | 1h45m |
| | **Reducing** | 5h11m | 5h13m |
| | **Merging** | 3h18m | 3h10m |
| | **Total** | 9h50m | 10h8m |

VirginiaTech
*Invent the Future*

# Implementation — Incremental Indexing

- Purpose
  - Process a continuous stream of HBase cell updates into live search indexes (Near Real-Time, NRT Indexing)
  - Solve the problem of frequent inserts, deletes and updates

- How does it work?
  - Enabling HBase replication (columnfamily)
  - Pointing an NRT Indexer Service at an HBase table
  - Starting an NRT Indexer Service

- Our work

**Source:**
http://www.cloudera.com/documentation/enterprise/5-6-x/topics/
search_config_hbase_indexer_for_search.html

8

VirginiaTech
1872
*Invent the Future*

# Implementation — Incremental Indexing

**Create and check the NRT indexer**

```
liuqing — cs5604f16_solr@node1:~ — ssh cs5604f16_solr@hadoop.dlib.vt.edu — 78×21

[cs5604f16_solr@node1 ~]$ hbase-indexer add-indexer --name NRTindexer --indexe
r-conf ~/ideal-cs5604f16-fake-morphline/morphline-hbase-mapper.xml --connectio
n-param solr.zk=node1.dlrl:2181,node2.dlrl:2181,node3.dlrl:2181,node4.dlrl:218
1,solr2.dlrl:2181/solr --connection-param solr.collection=ideal-cs5604f16-fake
 --zookeeper node1.dlrl:2181,node2.dlrl:2181,node3.dlrl:2181,node4.dlrl:2181,s
olr2.dlrl:2181
```

```
liuqing — cs5604f16_solr@node1:~ — ssh cs5604f16_solr@hadoop.dlib.vt.edu — 78×21

[cs5604f16_solr@node1 ~]$ hbase-indexer list-indexers
ZooKeeper connection string not specified, using default: localhost:2181

Number of indexes: 1

NRTindexer
  + Lifecycle state: ACTIVE
  + Incremental indexing state: SUBSCRIBE_AND_CONSUME
  + Batch indexing state: INACTIVE
  + SEP subscription ID: Indexer_NRTindexer
  + SEP subscription timestamp: 2016-11-24T19:26:45.331-05:00
  + Connection type: solr
  + Connection params:
    + solr.collection = ideal-cs5604f16-fake
    + solr.zk = node1.dlrl:2181,node2.dlrl:2181,node3.dlrl:2181,node4.dlrl:218
1,solr2.dlrl:2181/solr
```

**Solr Team Final Presentation**

VirginiaTech
1872
*Invent the Future*

# Implementation — Incremental Indexing

**Restart the HBase Solr Indexer service**

**Restart the service in HC**

**Restart the service in VC**

**Solr Team Final Presentation**

VirginiaTech
*Invent the Future*

# Implementation — Incremental Indexing

**Check the results in HBase and Solr Admin UI**

**Solr Team Final Presentation**

Virginia Tech
1872
*Invent the Future*

- Types
  - **Textual similarity based**
  - Collaborative filtering

- More Like This Component
  - Identifies similar documents to search result documents.
  - Can be configured as a **request handler** or search component
  - Uses term vectors to compute similarity.
  - Term vector can be calculated during query runtime or precomputed during indexing
  - Extracts highest matching terms based on tf-idf similarity

12

- schema.xml
  - Set stored = true
  - Set termVectors = true (for calcalating tf-idf)
    - After making changes, reindexing is mandatory

- solrconfig.xml
  - Enable mlt

```xml
<requestHandler name="/mlt" class="solr.MoreLikeThisHandler">
    <lst name="defaults">
        <str name="rows">5</str>
        <str name="mlt.fl">text_txt</str>
        <str name="mlt.mintf">1</str>
    </lst>
</requestHandler>
```

  - Define other configuration parameters
  - e.g., mlt.fl, mlt.mintf, mlt.mindf, mlt.maxdf, mlt.qf

13

VirginiaTech
1872
*Invent the Future*

- Request Handler



**Link:**
https://drive.google.com/open?id=0B2iasHDgHqGyYUk0R3RkVktkM2M

# Implementation — Recommendation

- Search Component



**Link:**
https://drive.google.com/open?id=0B2iasHDgHqGyU0doVEpidlh3c2c

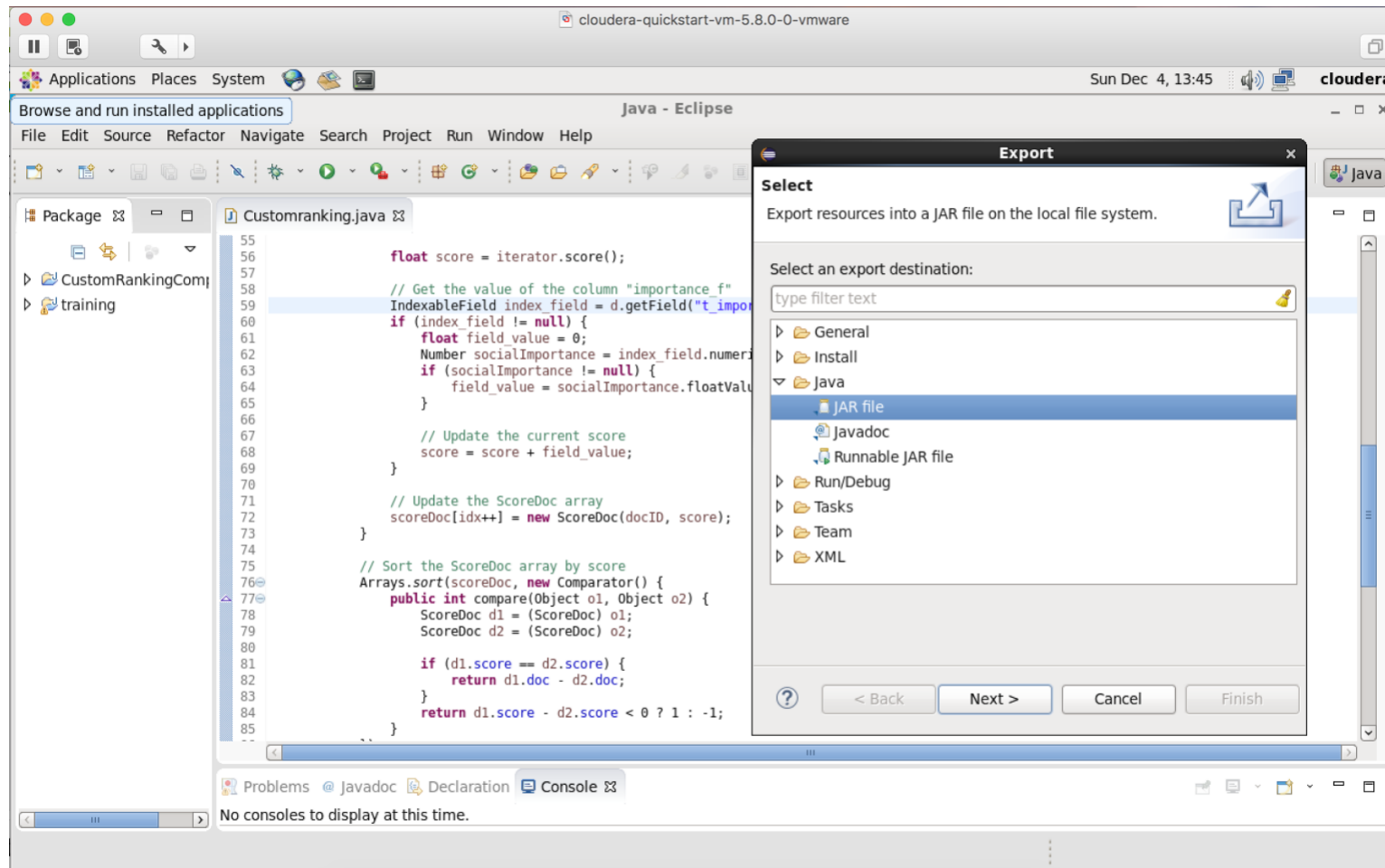VirginiaTech
1872
*Invent the Future*

- ## Purpose
  - ### Customize and optimize the ranked results

- ## How does it work?
  - ### Search Component
    - prepare(): pre-processing, invoked before query is executed
    - processing(): post-processing, invoked after all the results are fetched
  - ### Custom Scoring

$$Score = Doc_{score,Solr} + Doc_{importance}$$
$$+ W_{topic} \times Doc_{score,topic} + W_{cluster} \times Doc_{score,cluster}$$

  - ### Re-ranking

VirginiaTech
1872
*Invent the Future*

# Implementation — Custom Ranking

**Build and copy jar file into Hadoop Cluster**

**Solr Team Final Presentation**

**Modify the solrconfig.xml**

```
liuqing — cs5604f16_solr@node1:~/ideal-cs5604f16-fake/conf — ssh cs5604
<requestHandler name="/custom" class="solr.SearchHandler">
<!-- default values for query parameters can be specified, these
                will be overridden by parameters in the request
    -->
 <lst name="defaults">
    <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>Customranking</str>
    </arr>

</requestHandler>

<searchComponent name="Customranking" class="cs5604f16.solr.Customranking">
</searchComponent>
```

```
liuqing —
<lib dir="../../..
<lib dir="../../..

<lib dir="../../..
<lib dir="../../..

<lib dir="../../..
<lib dir="../../..

<lib dir="../../../contrib/velocity/lib" regex=".*\.jar" />
<lib dir="../../../dist/" regex="solr-velocity-\d.*\.jar" />

<lib dir="/home/cs5604f16_solr/bin/" regex=".*\.jar" />
```

**Solr Team Final Presentation**

Virginia Tech
1872
*Invent the Future*

# Implementation — Custom Ranking

**Update the instanceDir**

**Reload the collection**

**Check the results in Solr Admin UI**

**Solr Team Final Presentation**

VirginiaTech
*Invent the Future*

# Implementation — Solr Admin UI



DashBoard: provide basic functions for users to choose. (Logging to check Solr logs for debugging)

Core Selector: select the core (dataset) for queries

Solr instance Information: current versions, JVM information

Choose ideal-cs5604f16-fake for querying    19

# Implementation — Solr Admin UI



**Request–Handler (qt)**
`/select`

**1** → The request-handler: /select

**— common**

**q**
`text_txt:happy`

**2** → The query event: q

→ Parameters for query: fq (filter queries) sort (descending or ascending)

**fq**
[ ] ➖ ➕

**sort**

→ Execute query

**start, rows**
`0`  `10`

→ Results outputs: json format

**fl**

**df**

**Raw Query Parameters**
`key1=val1&key2=val2`

**wt**
`json` ⬍

☑ indent
☐ debugQuery

**3**

☐ dismax
☐ edismax
☐ hl
☐ facet
☐ spatial
☐ spellcheck

**4**  **Execute Query**

**5**

```
"response": {
    "numFound": 1      Result statistics
    "start": 0,
    "docs": [
        {
            "title_s": "Police Shadow Journalists Charlie Hebdo Gathering in Beijing",
            "sub_urls_ss": [
                "http://online.barrons.com/home-page|http://bigcharts.marketwatch.com|http:
            ],
            "organization_s": "",
            "url_s": "http://blogs.wsj.com/chinarealtime/2015/01/09/police-shadow-journal
            "location_ss": [
                "com"
            ],
            "text_txt": [    Field name
                "Journalists hold signs saying. am Charlie in French and Chinese on Thursda
            ],
```

22

**Solr Team Final Presentation**

Virginia Tech
1872
*Invent the Future*

# Implementation — Solr Admin UI



**1** The faceted search query: range

**2** Faceted search field: t_month_i

**3** Parameters, true when enabled

**4** Search Results: counts

**5** Search Results: details

**Solr Team Final Presentation**

Virginia Tech
*Invent the Future*

# Problem Faced

## Cloudera and OS

Virtual Cloudera seems slow and often crashes due to the memory

Not familiar with the whole architecture at the beginning

Versions of Cloudera and Solr
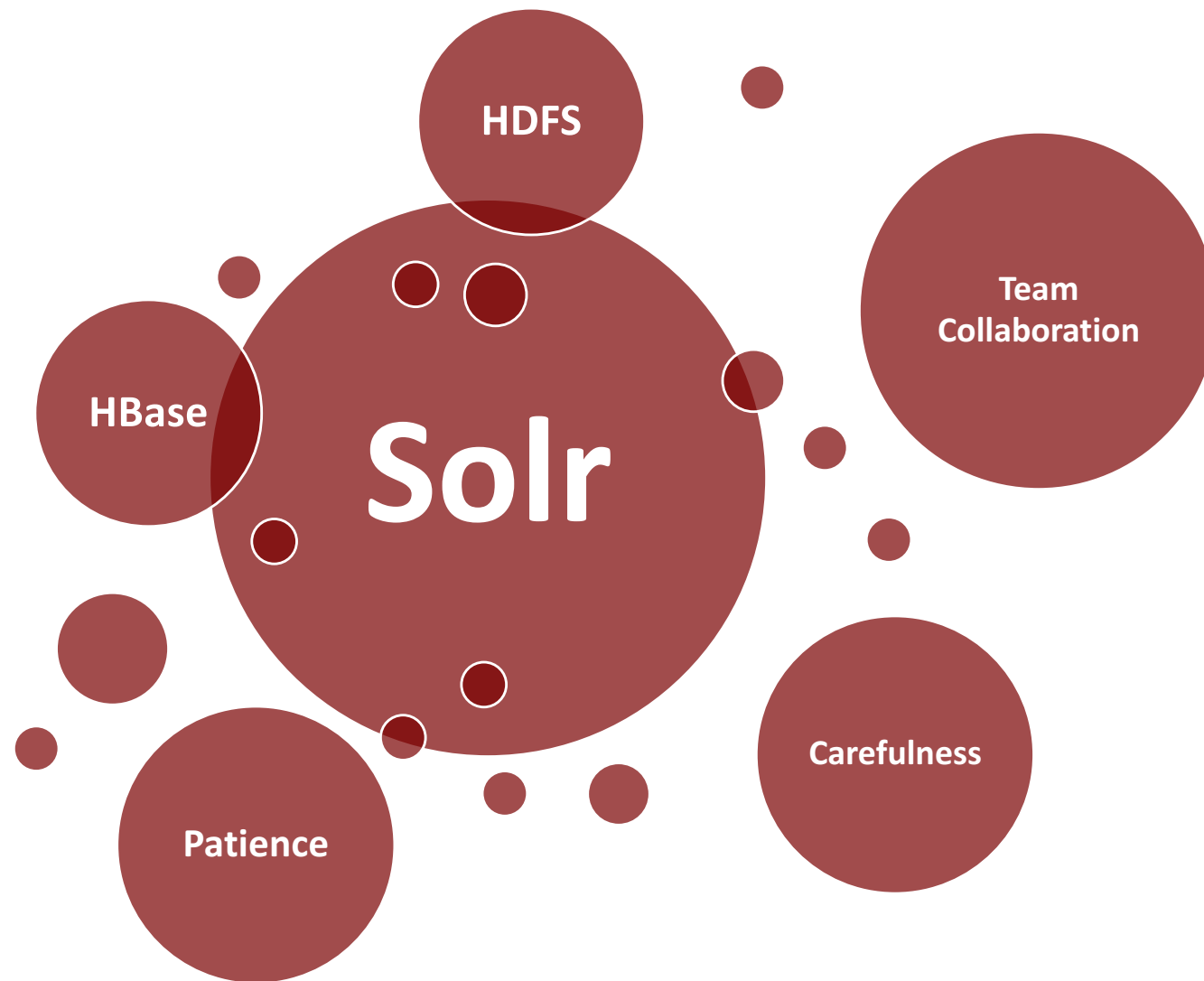
## Data

Consistency check

Not enough real data available to perform tests

Not much information available regarding logs to perform collaborative filtering

## Collaboration

Communication and modification

**Solr Team Final Presentation**

**VirginiaTech**
*Invent the Future*

# Lessons Learned

**Solr Team Final Presentation**

# Future Work

## Search

Customize more request handlers

Deal with the profanity issue

## Custom Ranking

Customize more search components

## Recommendation

Create a custom recommendation component (Probabilities – CTA team)

Implement the collaborative filtering (Log files – FE team)

## Solr

Figure out SolrCloud, multiple Solr nodes in Cloudera Search

**Solr Team Final Presentation**

VirginiaTech
1872
*Invent the Future*

# Acknowledgement

## Projects

| | |
|---|---|
| NSF IIS - 1319578 | III: Small: Integrated Digital Event Archiving and Library (IDEAL) |
| NSF IIS - 1619028 | III: Small: Collaborative Research: Global Event and Trend Archive Research (GETAR) |

## Teams

CMT, CMW, CLA, CTA, FE teams

## Persons

| | |
|---|---|
| Instructor | Dr. Edward A. Fox |
| GRA | Sunshin Lee |

**Solr Team Final Presentation**

VirginiaTech
1872
*Invent the Future*

# Thank you !

## Questions?

Virginia Tech

*Invent the Future*