# Search/Discovery "Under the Hood"
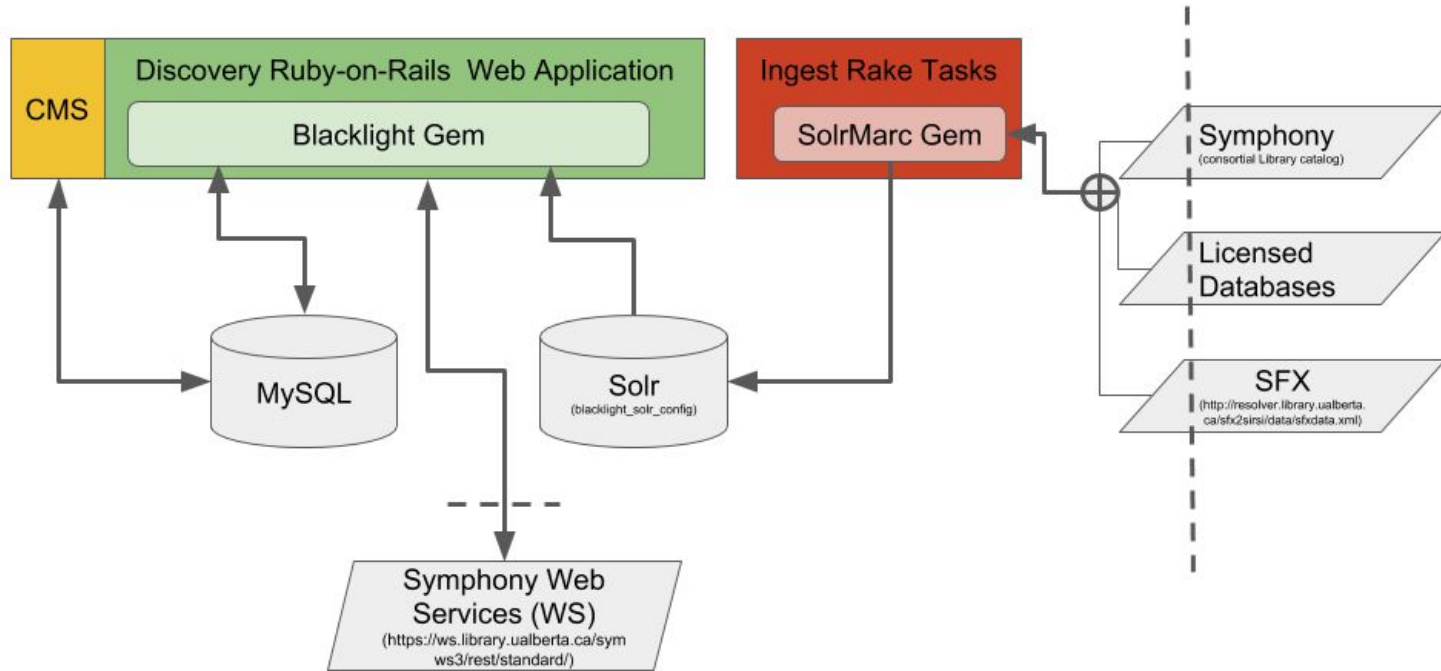
Tricia Jenkins and Sean Luyk | Spring Training 2019

# Outline

- Search in libraries
- Search trends
- Search "under the hood"

The Discovery Technology Stack

- Open Source Apache Project since 2007
- Webserver providing search capabilities
- Based on Apache Lucene
- Main competitor: Elastic Search
- Powers:

"

---

"Compared with the research tradition developed in information science and subsequently diffused to computer science, **the historical antecedents for understanding information retrieval in librarianship and indexing are far longer but less widely influential today**"

Warner, Julian. *Human Information Retrieval*. MIT Press: 2010

# Search in Libraries

# Search Goal #1

Retrieve all relevant documents for a user query, while retrieving as few non-relevant documents as possible
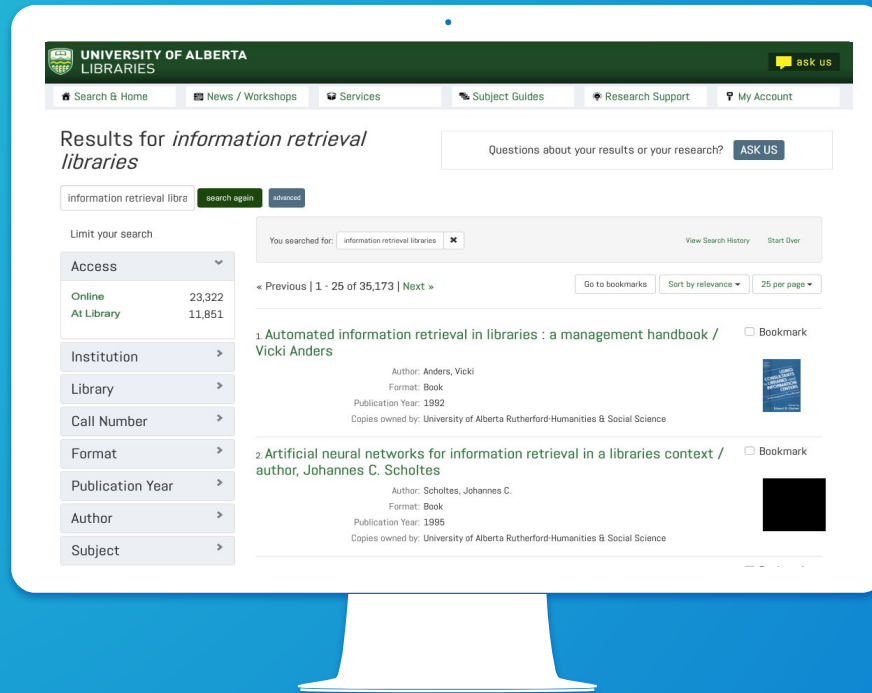
# What makes search results "relevant"?

It's all about expectations...

# Search Relevance is Hard

Users: relevant to <u>me</u>
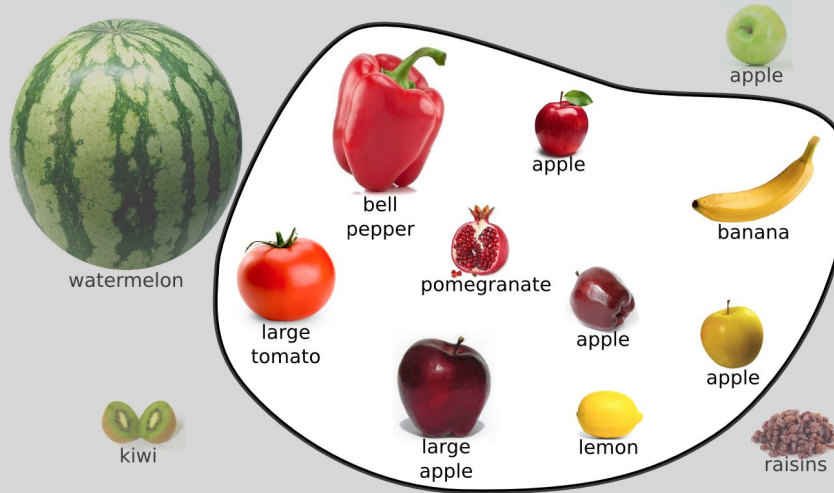


Technologists: relevant as defined by the model
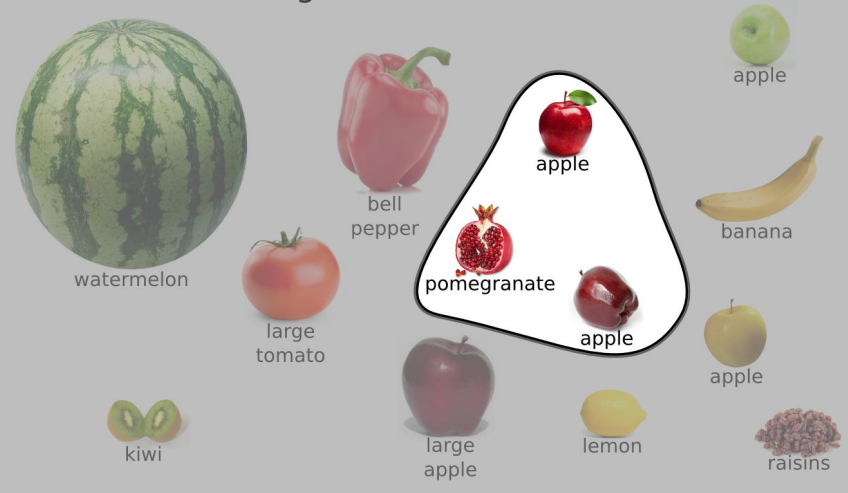
# Expectations for Precision Vary

# Relevance and Precision are Always at Odds

## Search query: "apples"

Berryman, John. "Search Precision and Recall by Example" <https://opensourceconnections.com/blog/2016/03/30/search-precision-and-recall-by-example>.

# Search Goal #2

Provide users with a good search experience

# What makes for a "good" user experience?

How do we know if we're providing users with a good search experience?

"

"To design the best UX, pay attention to what users **do**, not what they **say**. Self-reported claims are unreliable, as are user speculations about future behavior. Users do not know what they want."

Nielsen, Jakob. "First Rule of Usability? Don't Listen to Users"

<https://www.nngroup.com/articles/first-rule-of-usability-dont-listen-to-users/>

How do our users search?

———

What are their priorities?

———

How do different user groups search?

# Search Trends in Libraries

# Focus on Delivery, Ditch Discovery (Utrecht)

- Improve delivery at point of need (e.g. Google Scholar)
- Don't invest in discovery. Let users use the systems they already do
- Provide good information on the best search engines for different kinds of materials

# Coordinated Discovery (UW-Madison)

- Show users information categories
- Connect searches across the categories, and recommend relevant resources from other categories
- Promote serendipitous discovery
- Present different metadata for different categories
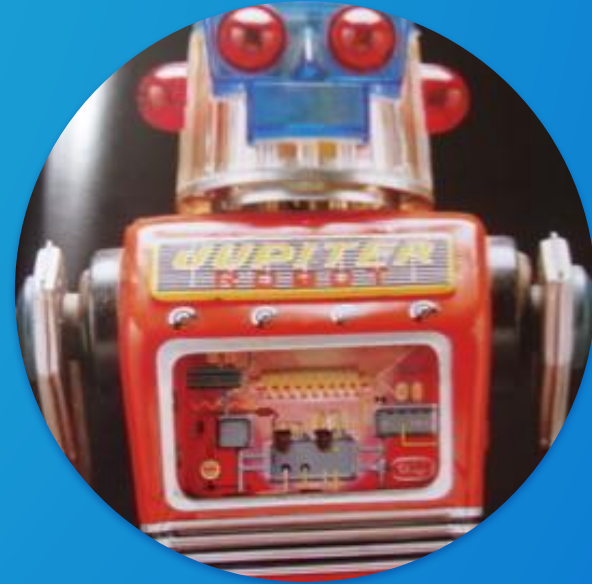- UI = not bento, but also not jambalaya

https://www.library.wisc.edu/experiments/coordinated-discovery/

# New Developments

# Machine Learning/AI Assisted Search

- Use supervised/unsupervised machine learning to improve search relevance
- Use real user feedback (result clicks) and/or document features (e.g. quality) to train a learning to rank (LTR) model
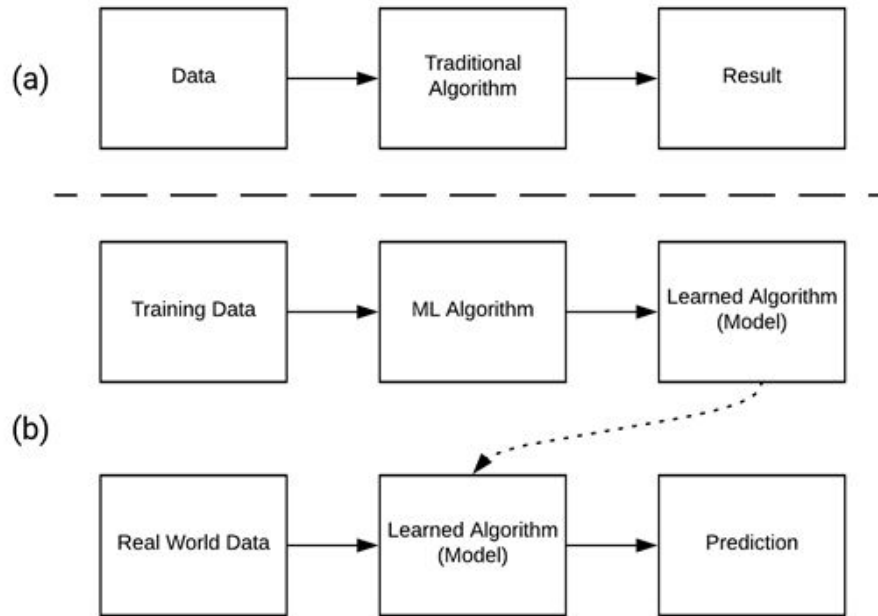
# Machine Learning (in a nutshell)



Figure 1. A traditional algorithm (a) versus a machine learning algorithm (b).

Harper, Charlie. "Machine Learning and the Library or: How I Learned to Stop Worrying and Love My Robot Overlords." *Code4Lib Journal* 41 <https://journal.code4lib.org/articles/13671>

# Machine Learning-Powered Discovery

Some examples...

- Carnegie Museum of Art Teenie Harris Archives
  - Automated metadata improvement, facial recognition: https://github.com/cmoa/teenie-week-of-play


- Capacity building: Fantastic Futures, Stanford Library AI Initiative/Studio

# Clustering/Visualization

- Use cluster analysis methods to group similar objects
- Example: Carrot2 (open source clustering engine)
- Example: Stanford's use of Yewno

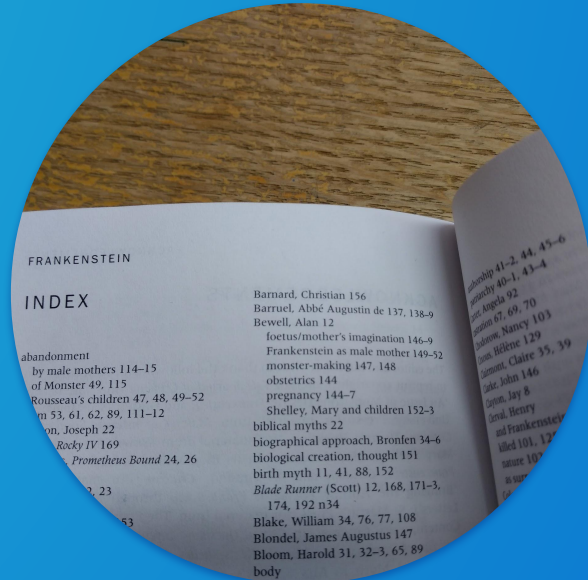# Search Under the Hood

# Index

If you are trying to find a subject in a book, where do you look first?

# Indexing Concepts

## Inverted Index

A searchable index that lists every word and the documents that contain those words, similar to an index in the back of a book which lists words and the pages on which they can be found. Finding the term before the document saves processing resources and time.

## Stemming

A stemmer is basically a set of mapping rules that maps the various forms of a word back to the base, or stem, word from which they derive.

```
02831cga a2200481 a 4500
001 2117026
007 vf cbaho
008 930913s1993    abc020          vleng d
035    $a ocm30704841
040    $b eng
055  3 $a Z 710 $b F494 1993
090    $a Z 710 F494 1993 $b AEU
090  0 $a Z 710 F559 1993 $b ARDC
090    $a Z 710 F494 1993 $b AEC
245 00 Finding Frankenstein $h [videorecording] : $b an introduction to the University of Alberta Library system / 
260    $a Edmonton, Alta. : $b Vicom, $c c1993.
300    $a 1 videocassette (20 min.) : $b sd., col. ; $c 1/2 in.
336    $a two-dimensional moving image $b tdi $2 rdacontent
337    $a video $b v $2 rdamedia
338    $a videocassette $b vf $2 rdacarrier
500    $a VHS.
500    $a Known as: University of Alberta Library instruction video.
500    $a Available in French with title: A la recherche de Frankenstein : une initiation au système de bibliothèque de l'Albert
596    $a 38 42 43 48
610 20 $a University of Alberta. $b Library.
650  0 $a Library orientation $x Aids and devices.
650  0 $a Library orientation for college students.
650  0 $a Information services $x User education.
710 2  $a Vicom Ltd.
740 0  $a Introduction to the University of Alberta Library system.
740 0  $a University of Alberta Library instruction video.
740 2  $a A la recherche de Frankenstein.
740 4  $a Une initiation au système de bibliothèque de l'Alberta.
926    $a Z 710 F494 1993 $w LC $c 9 $i 0162000388809 $d 4/12/2001 $l ON_SHELF $m UASCITECH $n 8 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 5 $i 0162000388742 $d 10/12/2016 $e 9/21/2016 $l ON_SHELF $m UAHSS $n 10 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 6 $i 0162000388759 $d 10/12/2016 $e 9/29/2016 $l ON_SHELF $m UAHSS $n 3 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 7 $i 0162000388767 $d 10/12/2016 $e 9/29/2016 $l ON_SHELF $m UAHSS $n 4 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 8 $i 0162000388775 $d 10/12/2016 $e 9/29/2016 $l ON_SHELF $m UAHSS $n 9 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 2 $i 0162000388718 $d 2/1/1996 $l ON_SHELF $m UAHLTHSC $n 2 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993
926    $a Z 710 F494 1993 $w LC $c 3 $i 0162000388726 $d 4/10/2013 $e 11/21/2006 $l ON_SHELF $m UARCRF $n 1 $p $150.00 $r Y $s Y $t MAG_MEDIA $u 9/13/1993 $o .STAFF. A0074008
926    $a Z 710 F494 1993 $w LC $c 1 $i 0162000388700 $l READONSITE $m UARCRF $p $150.00 $r Y $s Y $t NO_LOAN $u 9/13/1993 $o .STAFF.  A0018532
926    $a Z 710 F494 1993 $w LC $c 10 $i 0162009685411 $l ON_SHELF $m UARCRF $p $150.00 $r Y $s Y $t MAG_MEDIA $u 8/17/1999 $o .STAFF.  A0022869
926    $a Z 710 F494 1993 $w LC $c 11 $i 0162009685429 $l ON_SHELF $m UARCRF $p $150.00 $r Y $s Y $t MAG_MEDIA $u 8/17/1999 $o .STAFF.  A0022869
```

Finding Frankenstein [videorecording] : an introduction to the University of Alberta Library system / produced for University of Alberta Library

| | |
|---|---|
| Additional authors/performers: | Vicom Ltd. |
| Format: | Video or Projection |
| Published: | Edmonton, Alta: Vicom |
| Year: | 1993 |
| Physical Details: | 1 videocassette (20 min.) : sd., col. ; 1/2 in |
| General Note: | VHS. -- Known as: University of Alberta Library instruction video. -- Available in French with title: A la recherche de Frankenstein : une initiation au système de bibliothèque de l'Alberta. |
| Object type: | videorecording |

# Another example

Frankenstein : or, The modern Prometheus.(The 1818 text) Edited, with variant readings, an introd., and notes, by James Rieger

Author: Shelley, Mary Wollstonecraft, 1797-1851

Format: Book

Published: Indianapolis: Bobbs-Merrill

Year: 1974

Physical Details: xiv, 287 p. illus. 21 cm

ISBN: 0672514575

Series: The Library of literature

```
00961cam a2200265 | 4500
001 38596
008 073074s1974    inu              0 eng l
010    $a     72080409
020    $a 0672514575
035    $a ocm00415598
040    $a DLC $b eng
049    $a aeu $b eng
050 0  $a PR 5397 $b F82 1974
090 00 $a PR 5397 F82 1974 $b AEU
100 1  $a Shelley, Mary Wollstonecraft, $d 1797-1851.
245 10 $a Frankenstein : $b or, The modern Prometheus.(The 1818 text) $c Edited, with variant readings, an introd., and n
260    $a Indianapolis, $b Bobbs-Merrill $c [1974.]
300    $a xiv, 287 p. $b illus. $c 21 cm.
336    $a text $b txt $2 rdacontent
337    $a unmediated $b n $2 rdamedia
338    $a volume $b nc $2 rdacarrier
490 0  $a The Library of literature
596    $a 43
740 4  $a The modern Prometheus.
926    $a PR 5397 F82 1974 $w LC $c 1 $i 000000895789 $d 9/4/2018 $e 9/4/2018 $k CHECKEDOUT $l ON_SHELF $m UAHSS $n 94 $p $150.00 $r M $s Y $t BOOK $u 10/25/1988
```

https://search.library.ualberta.ca/catalog/38596

# Marc Mapping

```
27   # Title fields
28   #    primary title
29   title_t = custom, getLinkedFieldCombined(245[a-z])
30   title_display = 245[a-bd-z]
31   title_vern_display = custom, getLinkedField(245a)
32
33   #    additional title fields
34   title_addl_t = custom, getLinkedFieldCombined(130[a-z]:240[a-z]:210ab:222ab:242abnp:243[a-gk-s]:246[a-gnp]:247[a-gnp])
35   title_added_entry_t = custom, getLinkedFieldCombined(700[gk-pr-t]:710[fgk-t]:711fgklnpst:730[a-gk-t]:740anp)
36   title_series_t = custom, getLinkedFieldCombined(440anpv:490av)
37   title_sort = custom, getSortableTitle
38   edition_tesim = 250a
39   alternate_display_tesim = 880a
40   responsibility_display = 245c
```

https://github.com/ualbertalib/discovery/blob/master/config/SolrMarc/symphony_index.properties

# Analysis Chain

Index Analyzer: org.apache.solr.analysis.TokenizerChain

Tokenizer: **org.apache.lucene.analysis.standard.StandardTokenizerFactory**
- class: solr.StandardTokenizerFactory
- luceneMatchVersion: 6.6.0

Token Filters: **org.apache.lucene.analysis.icu.ICUFoldingFilterFactory**
- class: solr.ICUFoldingFilterFactory
- luceneMatchVersion: 6.6.0

**org.apache.lucene.analysis.core.StopFilterFactory**
- words: stopwords.txt
- class: solr.StopFilterFactory
- ✔ ignoreCase
- luceneMatchVersion: 6.6.0

**org.apache.lucene.analysis.snowball.SnowballPorterFilterFactory**
- language: English
- class: solr.SnowballPorterFilterFactory
- luceneMatchVersion: 6.6.0

# Finding Frankenstein [videorecording] : an introduction to the University of Alberta Library system

| ST | Finding | Frankenstein | videorecording | an | introduction | to | the | University | of | Alberta | Library | system |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ICUFF | finding | frankenstein | videorecording | an | introduction | to | the | university | of | alberta | library | system |
| SF | finding | frankenstein | videorecording | | introduction | | | university | | alberta | library | system |
| SF | find | frankenstein | videorecord | | introduct | | | univers | | alberta | librari | system |

# Frankenstein : or, The modern Prometheus.(The 1818 text)

| ST | Frankenstein | or | The | modern | Prometheus | The | 1818 | text |
|---|---|---|---|---|---|---|---|---|
| ICUFF | frankenstein | or | the | modern | prometheus | the | 1818 | text |
| SF | frankenstein | | | modern | prometheus | | 1818 | text |
| SF | frankenstein | | | modern | prometheus | | 1818 | text |

# Inverted Index

| word | documents |
|---|---|
| frankenstein | doc1, doc2 |
| edit | doc2 |
| system | doc1 |
| modern | doc2 |
| introd | doc2 |
| introduct | doc1 |
| jame | doc2 |
| librari | doc1 |
| videorecord | doc1 |
| note | doc2 |

| word | documents |
|---|---|
| produc | doc2 |
| prometheus | doc2 |
| read | doc2 |
| rieger | doc2 |
| find | doc1 |
| text | doc2 |
| univers | doc1 |
| variant | doc2 |
| alberta | doc1 |
| 1818 | doc2 |

# Document Term Frequency

# Now repeat for many different attributes

We use a dynamic schema which defines many common types that can be used for searching, display and faceting. We apply these to title, author, subject, etc.

| Use Case | indexed | stored | multiValued | omitNorms | termVectors | termPositions | docValues |
|---|---|---|---|---|---|---|---|
| search within field | true | | | | | | |
| retrieve contents | | true[8] | | | | | true[8] |
| use as unique key | true | | false | | | | |
| sort on field | true[7] | | false[9] | true [1] | | | true[7] |
| highlighting | true[4] | true | | | true[2] | true [3] | |
| faceting [5] | true[7] | | | | | | true[7] |
| add multiple values, maintaining order | | | true | | | | |
| field length affects doc score | | | | false | | | |
| MoreLikeThis [5] | | | | | true [6] | | |

# Search Concepts

## DisMax

DisMax stands for Maximum Disjunction. The DisMax query parser takes responsibility for building a good query from the user's input using Boolean clauses containing multiple queries across fields and any configured boosts.

## Boosting

Applying different weights based on the significance of each field.

# DisMax

```xml
<str name="qf">
    id^100000
    isbn_t^100000
    issn_t^100000
    lc_callnum_display^100000
    title_unstem_search^100000
    title_tesim^...
    subtitle_unstem_search^100000
    title_t^25000
    subtitle_t^25000
    title_addl_unstem_search^25000
    title_addl_t^25000
    earlier_title_tesim^25000
    later_title_tesim^25000
    title_added_entry_unstem_search^1500
    title_addl_entry_t^1250
    ...
    subject_topic_unstem_search^1000
    subject_t^500
    author_addl_unstem_search^250
    author_t^100
    author_addl_t^50
    contents_tesim^50
    databasedescription_tesim^50
    subject_addl_unstem_search^250
    subject_addl_t^50
    gmd_tesim^50
    summary_holdings_tesim^50
    title_series_unstem_search^25
    local_note_tesim^25
    general_note_tesim^25
    awards_note_tesim^25
    title_series_t^10
    section_number_tesim^10
    section_name_tesim^10
<!-- text -->
<!-- source -->
</str>
```

## mm

Minimum "Should" Match: specifies a minimum number of clauses that must match in a query.

`<str name="mm">6&lt;90%</str>`

## qf

Query Fields: specifies the fields in the index on which to perform the query.

## q

Defines the raw input strings for the query.

i.e. frankenstein

# Simplified Dismax

frankenstein

○

title^100000

subject^1000

author^250

=

title:frankenstein^100000 OR
subject:frankenstein^1000 OR
author:frankenstein^250

# frankenstein

```
"(+DisjunctionMaxQuery(((subtitle_t:frankenstein)^25000.0 | (databasedescription_tesim:frankenstein)^50.0 | (gmd_tesim:frankenstein)^50.0 |
(isbn_t:frankenstein)^100000.0 | (lc_callnum_display:frankenstein)^100000.0 | (subject_addl_t:frankenstein)^50.0 | (general_note_tesim:frankenstein)^25.0 |
(title_addl_t:frankenstein)^25000.0 | (subject_t:frankenstein)^500.0 | (later_title_tesim:frankenstein)^25000.0 |
(subject_addl_unstem_search:frankenstein)^250.0 | (title_series_unstem_search:frankenstein)^25.0 | (issn_t:frankenstein)^100000.0 |
(subject_topic_unstem_search:frankenstein)^1000.0 | (title_unstem_search:frankenstein)^100000.0 | (awards_note_tesim:frankenstein)^25.0 |
(section_name_tesim:frankenstein)^10.0 | (earlier_title_tesim:frankenstein)^25000.0 | (title_addl_unstem_search:frankenstein)^25000.0 |
(title_t:frankenstein)^25000.0 | (title_tesim:frankenstein)^100000.0 | (publisher_tesim:frankenstein)^1000.0 | (id:frankenstein)^100000.0 |
(subtitle_unstem_search:frankenstein)^100000.0 | (title_series_t:frankenstein)^10.0 | (local_note_tesim:frankenstein)^25.0 |
(author_unstem_search:frankenstein)^250.0 | (subject_unstem_search:frankenstein)^750.0 | (author_t:frankenstein)^100.0 |
(author_addl_unstem_search:frankenstein)^250.0 | (contents_tesim:frankenstein)^50.0 | (author_addl_t:frankenstein)^50.0 |
(title_added_entry_t:frankenstein)^1250.0 | (summary_holdings_tesim:frankenstein)^50.0 | (title_added_entry_unstem_search:frankenstein)^1500.0 |
(section_number_tesim:frankenstein)^10.0)~0.01) DisjunctionMaxQuery(((subtitle_t:frankenstein)^250000.0 | (databasedescription_tesim:frankenstein)^500.0 |
(gmd_tesim:frankenstein)^500.0 | (isbn_t:frankenstein)^1000000.0 | (lc_callnum_display:frankenstein)^1000000.0 | (subject_addl_t:frankenstein)^500.0 |
(general_note_tesim:frankenstein)^250.0 | (title_addl_t:frankenstein)^25000.0 | (subject_t:frankenstein)^5000.0 | (later_title_tesim:frankenstein)^25000.0 |
(subject_addl_unstem_search:frankenstein)^2500.0 | (title_series_unstem_search:frankenstein)^250.0 | (source:frankenstein)^100000.0 |
(issn_t:frankenstein)^1000000.0 | (subject_topic_unstem_search:frankenstein)^10000.0 | (title_unstem_search:frankenstein)^1000000.0 |
(awards_note_tesim:frankenstein)^250.0 | (section_name_tesim:frankenstein)^100.0 | (earlier_title_tesim:frankenstein)^25000.0 |
(title_addl_unstem_search:frankenstein)^250000.0 | (title_t:frankenstein)^250000.0 | (title_tesim:frankenstein)^1000000.0 |
(publisher_tesim:frankenstein)^10000.0 | (id:frankenstein)^1000000.0 | (text:frankenstein)^10.0 | (subtitle_unstem_search:frankenstein)^1000000.0 |
(title_series_t:frankenstein)^100.0 | (local_note_tesim:frankenstein)^250.0 | (author_unstem_search:frankenstein)^2500.0 |
(subject_unstem_search:frankenstein)^7500.0 | (author_t:frankenstein)^1000.0 | (author_addl_unstem_search:frankenstein)^2500.0 |
(subject_topic_facet:frankenstein)^6250.0 | (contents_tesim:frankenstein)^500.0 | (author_addl_t:frankenstein)^500.0 |
(title_added_entry_t:frankenstein)^12500.0 | (summary_holdings_tesim:frankenstein)^500.0 | (title_added_entry_unstem_search:frankenstein)^15000.0 |
(section_number_tesim:frankenstein)^100.0)~0.01))/no_coord"
```

# Show Your Work

# Boolean Model + Vector Space Model

### Boolean query

A document either matches or does not match a query.
AND, OR, NOT

### IDF

Inverse document frequency deals with the problem of terms that occur too often in the collection to be meaningful for relevance determination.

### TF

Term frequency is the number of times a term occurs in a document. A document that mentions a query term more often has more to do with that query and therefore should receive a higher score.

# University of Alberta Library

```
"(+(DisjunctionMaxQuery(((subtitle_t:univers)^25000.0 | (databasedescription_tesim:univers)^50.0 | (gmd_tesim:univers)^50.0 | (isbn_t:univers)^100000.0 |
(lc_callnum_display:univers)^100000.0 | (subject_addl_t:univers)^50.0 | (general_note_tesim:univers)^25.0 | (title_addl_t:univers)^25000.0 | (subject_t:univers)^500.0 |
(later_title_tesim:univers)^25000.0 | (subject_addl_unstem_search:university)^250.0 | (title_series_unstem_search:university)^25.0 | (issn_t:univers)^100000.0 |
(subject_topic_unstem_search:university)^1000.0 | (title_unstem_search:university)^100000.0 | (awards_note_tesim:univers)^25.0 | (section_name_tesim:univers)^10.0 |
(earlier_title_tesim:univers)^25000.0 | (title_addl_unstem_search:university)^25000.0 | (title_t:univers)^25000.0 | (title_tesim:univers)^100000.0 | (publisher_tesim:univers)^1000.0 |
(id:univers)^100000.0 | (subtitle_unstem_search:university)^100000.0 | (title_series_t:univers)^10.0 | (local_note_tesim:univers)^25.0 | (author_unstem_search:university)^250.0 |
(subject_unstem_search:university)^750.0 | (author_t:univers)^100.0 | (author_addl_unstem_search:university)^250.0 | (contents_tesim:univers)^50.0 | (author_addl_t:univers)^50.0 |
(title_added_entry_t:univers)^1250.0 | (summary_holdings_tesim:univers)^50.0 | (title_added_entry_unstem_search:university)^1500.0 | (section_number_tesim:univers)^10.0)~0.01)
DisjunctionMaxQuery(((subtitle_t:alberta)^25000.0 | (databasedescription_tesim:alberta)^50.0 | (gmd_tesim:alberta)^50.0 | (isbn_t:alberta)^100000.0 | (lc_callnum_display:alberta)^100000.0 |
(subject_addl_t:alberta)^50.0 | (general_note_tesim:alberta)^25.0 | (title_addl_t:alberta)^25000.0 | (subject_t:alberta)^500.0 | (later_title_tesim:alberta)^25000.0 |
(subject_addl_unstem_search:alberta)^250.0 | (title_series_unstem_search:alberta)^25.0 | (issn_t:alberta)^100000.0 | (subject_topic_unstem_search:alberta)^1000.0 |
(title_unstem_search:alberta)^100000.0 | (awards_note_tesim:alberta)^25.0 | (section_name_tesim:alberta)^10.0 | (earlier_title_tesim:alberta)^25000.0 |
(title_addl_unstem_search:alberta)^25000.0 | (title_t:alberta)^25000.0 | (title_tesim:alberta)^100000.0 | (publisher_tesim:alberta)^1000.0 | (id:alberta)^100000.0 |
(subtitle_unstem_search:alberta)^100000.0 | (title_series_t:alberta)^10.0 | (local_note_tesim:alberta)^25.0 | (author_unstem_search:alberta)^250.0 | (subject_unstem_search:alberta)^750.0 |
(author_t:alberta)^100.0 | (author_addl_unstem_search:alberta)^250.0 | (contents_tesim:alberta)^50.0 | (author_addl_t:alberta)^50.0 | (title_added_entry_t:alberta)^1250.0 |
(summary_holdings_tesim:alberta)^50.0 | (title_added_entry_unstem_search:alberta)^1500.0 | (section_number_tesim:alberta)^10.0)~0.01) DisjunctionMaxQuery(((subtitle_t:librari)^25000.0 |
(databasedescription_tesim:librari)^50.0 | (gmd_tesim:librari)^50.0 | (isbn_t:librari)^100000.0 | (lc_callnum_display:librari)^100000.0 | (subject_addl_t:librari)^50.0 |
(general_note_tesim:librari)^25.0 | (title_addl_t:librari)^25000.0 | (subject_t:librari)^500.0 | (later_title_tesim:librari)^25000.0 | (subject_addl_unstem_search:library)^250.0 |
(title_series_unstem_search:library)^25.0 | (issn_t:librari)^100000.0 | (subject_topic_unstem_search:library)^1000.0 | (title_unstem_search:library)^100000.0 |
(awards_note_tesim:librari)^25.0 | (section_name_tesim:librari)^10.0 | (earlier_title_tesim:librari)^25000.0 | (title_addl_unstem_search:library)^25000.0 | (title_t:librari)^25000.0 |
(title_tesim:librari)^100000.0 | (publisher_tesim:librari)^1000.0 | (id:librari)^100000.0 | (subtitle_unstem_search:library)^100000.0 | (title_series_t:librari)^10.0 |
(local_note_tesim:librari)^25.0 | (author_unstem_search:library)^250.0 | (subject_unstem_search:library)^750.0 | (author_t:librari)^100.0 | (author_addl_unstem_search:library)^250.0 |
(contents_tesim:librari)^50.0 | (author_addl_t:librari)^50.0 | (title_added_entry_t:librari)^1250.0 | (summary_holdings_tesim:librari)^50.0 |
(title_added_entry_unstem_search:library)^1500.0 | (section_number_tesim:librari)^10.0)~0.01)~3 DisjunctionMaxQuery(((subtitle_t:\"univers ? alberta librari\"~3)^250000.0 |
(databasedescription_tesim:\"univers ? alberta librari\"~3)^500.0 | (gmd_tesim:\"univers ? alberta librari\"~3)^500.0 | (isbn_t:\"univers ? alberta librari\"~3)^1000000.0 |
(lc_callnum_display:\"univers ? alberta librari\"~3)^1000000.0 | (subject_addl_t:\"univers ? alberta librari\"~3)^500.0 | (general_note_tesim:\"univers ? alberta librari\"~3)^250.0 |
(title_addl_t:\"univers ? alberta librari\"~3)^25000.0 | (subject_t:\"univers ? alberta librari\"~3)^5000.0 | (later_title_tesim:\"univers ? alberta librari\"~3)^25000.0 |
(subject_addl_unstem_search:\"university ? alberta library\"~3)^2500.0 | (title_series_unstem_search:\"university ? alberta library\"~3)^250.0 | (source:University of Alberta
Library)^100000.0 | (issn_t:\"univers ? alberta librari\"~3)^1000000.0 | (subject_topic_unstem_search:\"university ? alberta library\"~3)^10000.0 | (title_unstem_search:\"university ?
alberta library\"~3)^1000000.0 | (awards_note_tesim:\"univers ? alberta librari\"~3)^250.0 | (section_name_tesim:\"univers ? alberta librari\"~3)^100.0 | (earlier_title_tesim:\"univers ?
alberta librari\"~3)^25000.0 | (title_addl_unstem_search:\"university ? alberta library\"~3)^25000.0 | (title_t:\"univers ? alberta librari\"~3)^250000.0 | (title_tesim:\"univers ? alberta
librari\"~3)^1000000.0 | (publisher_tesim:\"univers ? alberta librari\"~3)^10000.0 | (id:\"univers ? alberta librari\"~3)^1000000.0 | (text:\"univers ? alberta librari\"~3)^10.0 |
(subtitle_unstem_search:\"university ? alberta library\"~3)^1000000.0 | (title_series_t:\"univers ? alberta librari\"~3)^100.0 | (local_note_tesim:\"univers ? alberta librari\"~3)^250.0 |
(author_unstem_search:\"university ? alberta library\"~3)^2500.0 | (subject_unstem_search:\"university ? alberta library\"~3)^7500.0 | (author_t:\"univers ? alberta librari\"~3)^1000.0 |
(author_addl_unstem_search:\"university ? alberta library\"~3)^2500.0 | (subject_topic_facet:University of Alberta Library)^6250.0 | (contents_tesim:\"univers ? alberta librari\"~3)^500.0 |
(author_addl_t:\"univers ? alberta librari\"~3)^500.0 | (title_added_entry_t:\"univers ? alberta librari\"~3)^12500.0 | (summary_holdings_tesim:\"univers ? alberta librari\"~3)^500.0 |
(title_added_entry_unstem_search:\"university ? alberta library\"~3)^15000.0 | (section_number_tesim:\"univers ? alberta librari\"~3)^100.0)~0.01)/no_coord"
```

# Show Your Work



43

# Challenges

## Precision vs Recall

Were the documents that were returned supposed to be returned? Were all of the documents returned that were supposed to be returned?

## Phrase searching across fields

"Migrating library data a practical manual"

## Length Norms

matches on a smaller field score higher than matches on a larger field. "Managerial accounting garrison"

## Language

"L'armée furieuse" vs "armée furieuse"

## Minimum "Should" Match

british missions "south pacific"

## Boosting

UAL content or recency.

# Tuning

File  Edit  View  Search  Terminal  Help

1 users, ~0 active

Latest Interval Stats at 19:36:07

Average Times:
Full: 0.090
Connect: 0.000
Latency: 0.042
~Receive: 0.048

Percentiles:
0.0%: 0.001
50.0%: 0.025
90.0%: 0.172
95.0%: 0.410
99.0%: 0.410
99.9%: 0.410
100.0%: 0.410

Response Codes:
200: 81.82% (9)
400: 18.18% (2)
All: 100.00% (11)

Taurus
| v1.13.4 by BlazeMeter.com |

JMeter: Replay_Solr_Logs.jmx
Running...
Elapsed: 95:33:09                                      ETA: N/A

local
mem: 59.800
engine-loop: 7.537
disk-write: 113,89
bytes-recv: 93,501
disk-space: 29.000
conn-all: 46
cpu: 27.400

11 hits  2 fail

Cumulative Stats 95:33:08

Average Times:
Full: 0.145
Connect: 0.000
Latency: 0.137
~Receive: 0.008

Percentiles:
0.0%: 0.000
50.0%: 0.047
90.0%: 0.309
95.0%: 0.592
99.0%: 1.613
99.9%: 1.715
100.0%: 10.816

Response Codes:
200: 77.51% (1096911)
400: 17.33% (245236)
504: 1.28% (18139)
on HTTP response code:
java.lang.Error:
All: 100.00% (1415248)

0.090 avg time ( lat, conn)

| Labels | Hits | Failures | Avg Time |
|---|---|---|---|
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.query= | 244 | 8.61% | 0.267 |
| http://localhost:8983/solr/discovery/suggest?q=%225uccessful+Assessment+fo | 2 | 100.00% | 0.001 |
| http://localhost:8983/solr/discovery/suggest?q=experiencing+the+depths+of+ | 24 | 100.00% | 0.002 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.query= | 1156 | 26.82% | 0.128 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 1157 | 27.14% | 0.354 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=document&id=9549255 | 4172 | 0.00% | 0.003 |
| http://localhost:8983/solr/discovery/suggest?q=%225UCCESSFUL+ASSESSMENT+FO | 8 | 100.00% | 0.001 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.query= | 82 | 7.32% | 0.188 |
| http://localhost:8983/solr/discovery/select?defType=dismax&f.author_displa | 266 | 0.00% | 0.166 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=document&id=8302327 | 3306 | 0.00% | 0.006 |
| http://localhost:8983/solr/discovery/suggest?q=subject+analysis+in+online+ | 16 | 100.00% | 0.006 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 2504 | 1.00% | 0.369 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.query= | 6005 | 0.00% | 0.229 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 2 | 0.00% | 0.892 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 10522 | 0.00% | 0.097 |
| http://localhost:8983/solr/discovery/select?f.author_display.facet.limit=2 | 35928 | 2.01% | 1.315 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 381 | 13.39% | 0.279 |
| http://localhost:8983/solr/discovery/select?wt=ruby&qt=search&facet.field= | 4 | 0.00% | 0.097 |
| http://localhost:8983/solr/discovery/suggest?q=Weiss%2C+G.+ed.+1987.+Hazar | 6 | 100.00% | 0.001 |

09:55:30 WARNING: No details for errors of , dropped
info: [{'cnt': 2, 'tag': None, 'urls': Counter(),
'rc': '400', 'msg': u'Bad Request', 'type': 0}]

# Thanks!

## Any questions?

You can find us at

sean.luyk@ualberta.ca

tricia.jenkins@ualberta.ca