# Hacking Lucene for Custom Search Results

Doug Turnbull

OpenSource Connections

# Hello

## Me

@softwaredoug

dturnbull@o19s.com

## Us

http://o19s.com

- Trusted Advisors in Search, Discovery & Analytics

# Tough Search Problems

- We have demanding users!

Switch these two!

| | Patent ▼ | ID | Published |
|---|---|---|---|
| 1 | Cat cage door safety switch | CN202222241U 201120333888.2 | May 23, 2012 Sep 7, 2011 |
| 2 | Active cat eye | CN2383947Y 99233726.7 | Jun 21, 2000 Jul 19, 1999 |
| 3 | Two a door of cat eye | CN202202765U 201120287835.1 | Apr 25, 2012 Aug 8, 2011 |
| 4 | A cat sand groove | CN202014496U 201120029786.1 | Oct 26, 2011 Jan 28, 2011 |
| 5 | Cat special charging basket | CN101990845A 200910167190.5 | Mar 30, 2011 Aug 31, 2009 |
| 6 | The cat handheld cartridge | CN101965804A 201010501705.3 | Feb 9, 2011 Sep 30, 2010 |
| 7 | Cat special charging basket | CN201718305U 200920177806.2 | Jan 26, 2011 Aug 31, 2009 |

OpenSource Connections

# Tough Search Problems

- Demanding users!

| Patent ▼ | ID | Published |
|---|---|---|
| 1 A pedal shopping cart | CN202186472U 201120303228.X | Apr 11, 2012 Aug 19, 2011 |
| 2 A portable cart | CN201494478U 200920225545.7 | Jun 2, 2010 Aug 26, 2009 |
| 3 Supermarket using a shopping cart handle with billboard | CN201011609Y 200620153178.0 | Jan 23, 2008 Dec 1, 2006 |
| 4 A shopping cart | CN101648572A 200810048858.X | Feb 17, 2010 Aug 15, 2008 |
| 5 With the trolley of the steering device | CN102390417A 201110299049.8 | Mar 28, 2012 Sep 28, 2011 |
| 6 Displaying shopping cart for commodity position | CN202243588U 201120342156.X | May 30, 2012 Sep 9, 2011 |
| 7 A workpiece cart for drying room box | CN202057194U 201120133588.X | Nov 30, 2011 Apr 29, 2011 |
| 8 The vehicle cart device | CN102452405A 201110037837.X | May 16, 2012 Feb 16, 2011 |

WRONG!

Make search do
what is in my head!

OpenSource Connections ●

# Tough Search Problems

- Our Eternal Problem:

In one ear    Out the other

This is how a search engine works!

/dev/null

○ Customers don't care about the technology field of Information Retrieval: **they just want results**

○ BUT we are constrained by the tech!

# Satisfying User Expectations

- Easy: The Search Relevancy Game:
  - Solr/Elasticsearch **query operations** (boosts, etc)
  - **Analysis** of query/index to enhance matching



- Medium: Forget this, lets write some Java
  - Solr/Elasticsearch **query parsers**. Reuse existing Lucene Queries to get closer to user needs

# That Still Didn't Work

- Look at him, he's angrier than ever!

- For the toughest problems, we've made search complex and brittle

- **WHACK-A-MOLE:**
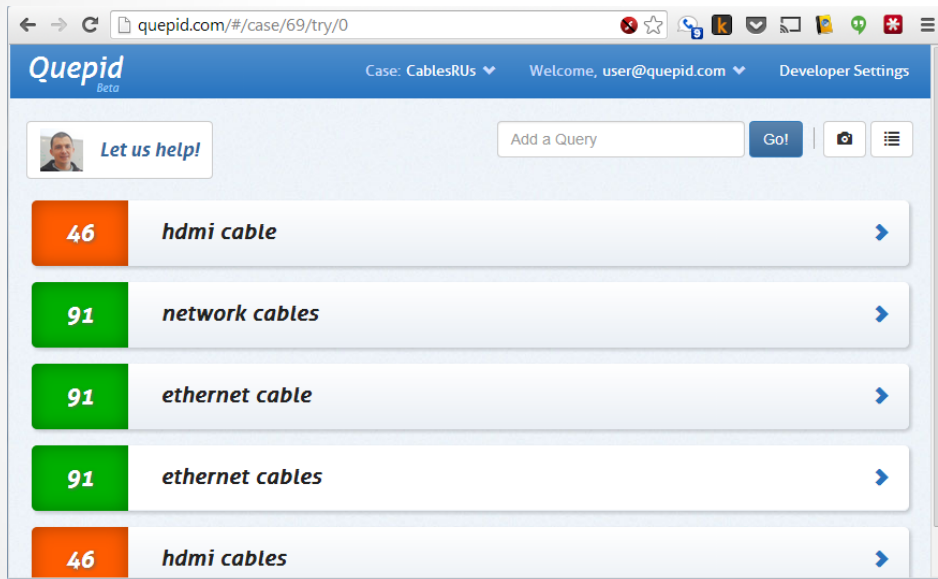  - o Fix one problem, cause another
  - o We give up,

# Next Level

- Hard: **Custom Lucene Scoring** – implement a query and scorer to explicitly control matching and scoring



This is the Nuclear Option!

# Shameless Plug

- How do we know if we're making progress?



- Quepid! – our search test driven workbench

# Lucene Lets Review

- At some point we wrote a Lucene index to a directory

- Boilerplate (open up the index):

```
Directory d = new RAMDirectory();
IndexReader ir = DirectoryReader.open(d);
IndexSearcher is = new IndexSearcher(ir);
```

Boilerplate setup of:
- <u>Directory</u> Lucene's handle to the FS
- <u>IndexReader</u> – Access to Lucene's data structures
- <u>IndexSearcher</u> – use index searcher to perform search

# Lucene Lets Review

- ## <u>Queries</u>:

Make a Query and Search!
- TermQuery: basic term
  search for a field

```
Term termToFind = new Term("tag", "space");
TermQuery spaceQ = new TermQuery(termToFind);
termToFind = new Term("tag", "star-trek");
TermQuery starTrekQ = new TermQuery(termToFind);
```
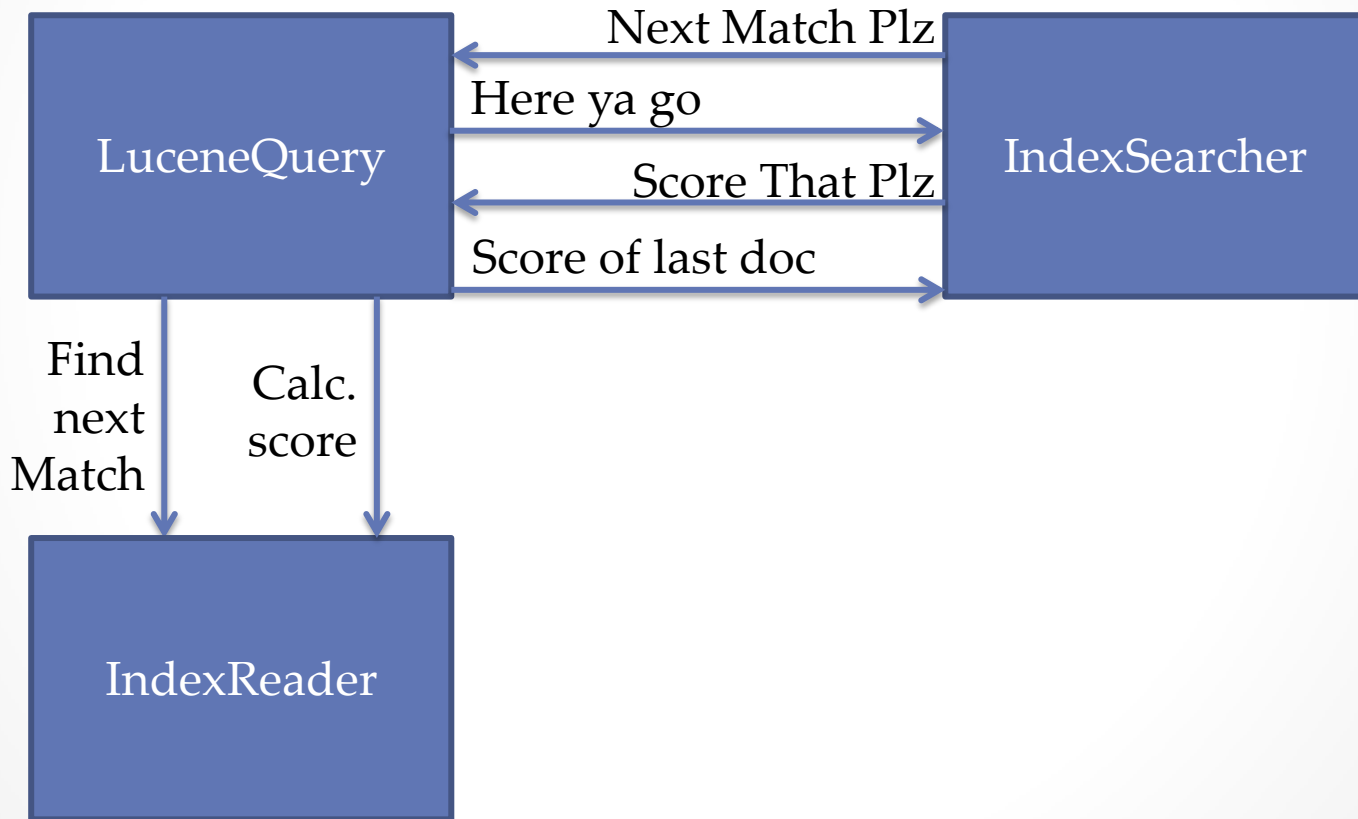
- ## <u>Queries That Combine Queries</u>

```
BooleanQuery bq = new BooleanQuery();
BooleanClause bClause = new BooleanClause(spaceQ, Occur.MUST);
BooleanClause bClause2 = new BooleanClause(starTrekQ, Occur.SHOULD);
bq.add(bClause);
bq.add(bClause2);
```

# Lucene Lets Review

- Query responsible for specifying search behavior
- Both:
  - **Matching** – what documents to include in the results

  - **Scoring** – how relevant is a result to the query by assigning a score
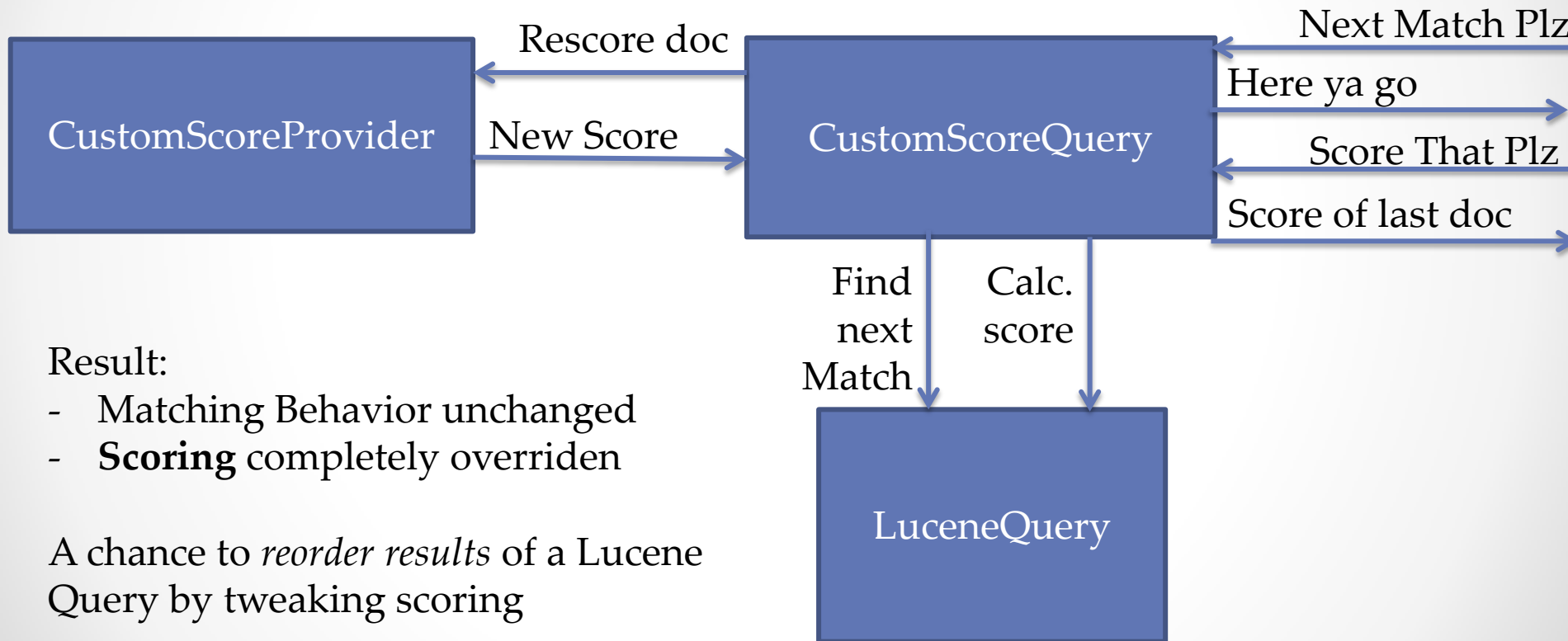
# Lucene Queries, 30,000 ft view

Aka, "not really accurate, but what to tell your boss to not confuse them"

LuceneQuery

IndexSearcher

Next Match Plz

Here ya go

Score That Plz

Score of last doc

Find next Match

Calc. score

IndexReader

# First Stop CustomScoreQuery

- Wrap a query but override its score

```
Rescore doc
CustomScoreProvider  ←────────  CustomScoreQuery
                     New Score →
```

Next Match Plz →
Here ya go →
Score That Plz →
Score of last doc →

CustomScoreQuery:
- Find next Match
- Calc. score
→ LuceneQuery

Result:
- Matching Behavior unchanged
- **Scoring** completely overriden

A chance to *reorder results* of a Lucene Query by tweaking scoring

# How to use?

- Use a normal Lucene query for **matching**

```
Term t = new Term("tag", "star-trek");
TermQuery tq = new TermQuery(t);
```

- Create & Use a CustomQueryScorer for **scoring** that wraps the Lucene query

```
CountingQuery ct = new CountingQuery(tq);
```

# Implementation

- Extend CustomScoreQuery, provide a CustomScoreProvider

```java
protected CustomScoreProvider getCustomScoreProvider(
            AtomicReaderContext context) throws IOException {
    return new CountingQueryScoreProvider("tag", context);
}
```

(boilerplate omitted)

OpenSource Connections

# Implementation

- CustomScoreProvider rescores each doc with IndexReader & docId

```
// Give all docs a score of 1.0
public float customScore(int doc, float subQueryScore, float
valSrcScores[]) throws IOException {
    return (float)(1.0f); // New Score
}
```

# Implementation

- Example: Sort by number of terms in a field

```java
// Rescores by counting the number of terms in the field
public float customScore(int doc, float subQueryScore, float
valSrcScores[]) throws IOException {
    IndexReader r = context.reader();
    Terms tv = r.getTermVector(doc, _field);
    TermsEnum termsEnum = null;
    termsEnum = tv.iterator(termsEnum);
    int numTerms = 0;
    while((termsEnum.next()) != null) {
        numTerms++;
    }
    return (float)(numTerms); // New Score

}
```

# CustomScoreQuery, Takeaway

- **SIMPLE!**
  - Relatively few gotchas or bells & whistles (we will see lots of gotchas)
- Limited
  - No tight control on what matches

- If this satisfies your requirements: You should get off the train here

# Lucene Circle Back

- I care about overriding **scoring**
  - CustomScoreQuery

- I need to control custom scoring and **matching**
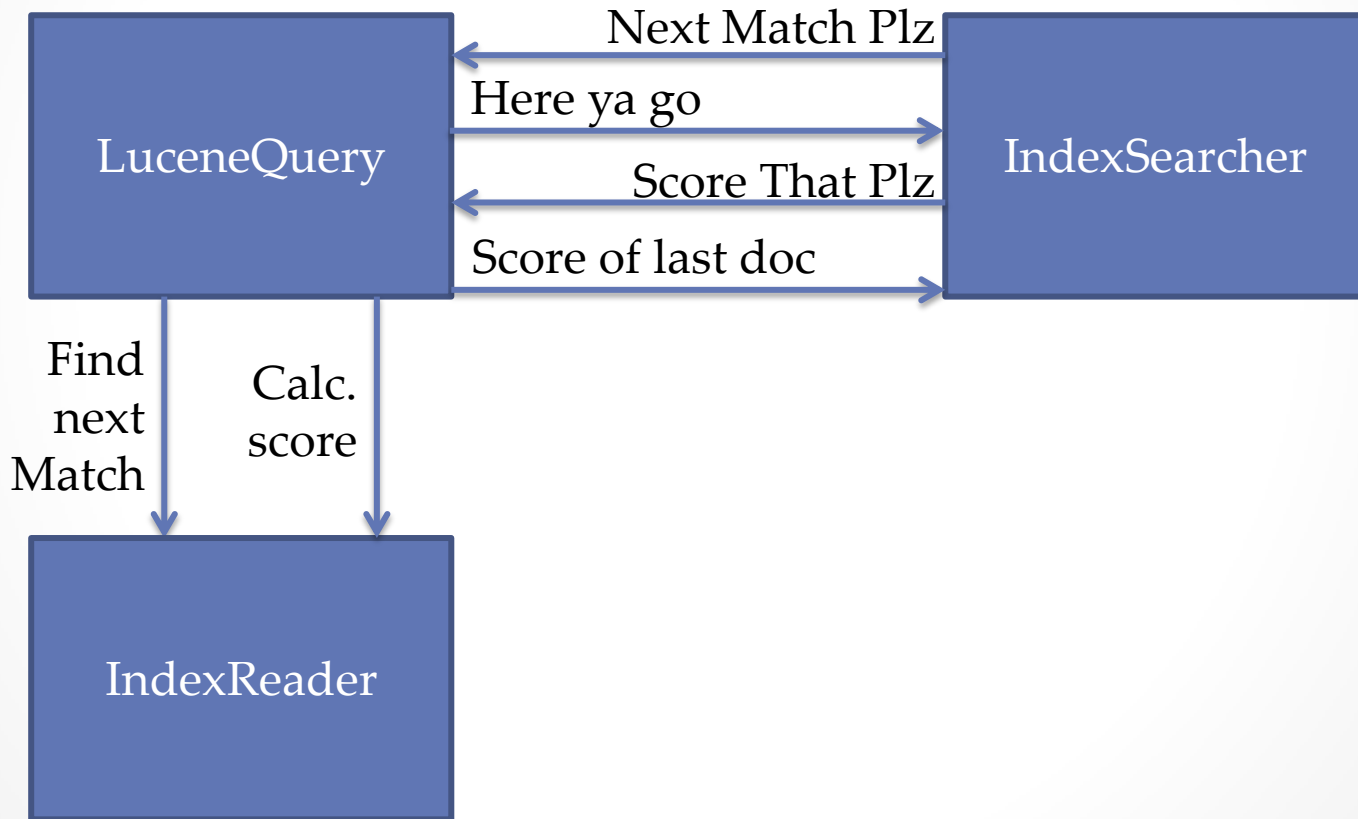  - Custom Lucene Queries!

# Example – Backwards Query

- Search for terms backwards!
  - Instead of banana, lets create a query that finds **ananab** matches and scores the document (5.0)
  - But lets also match forward terms **(banana),** but with a lower score (1.0)

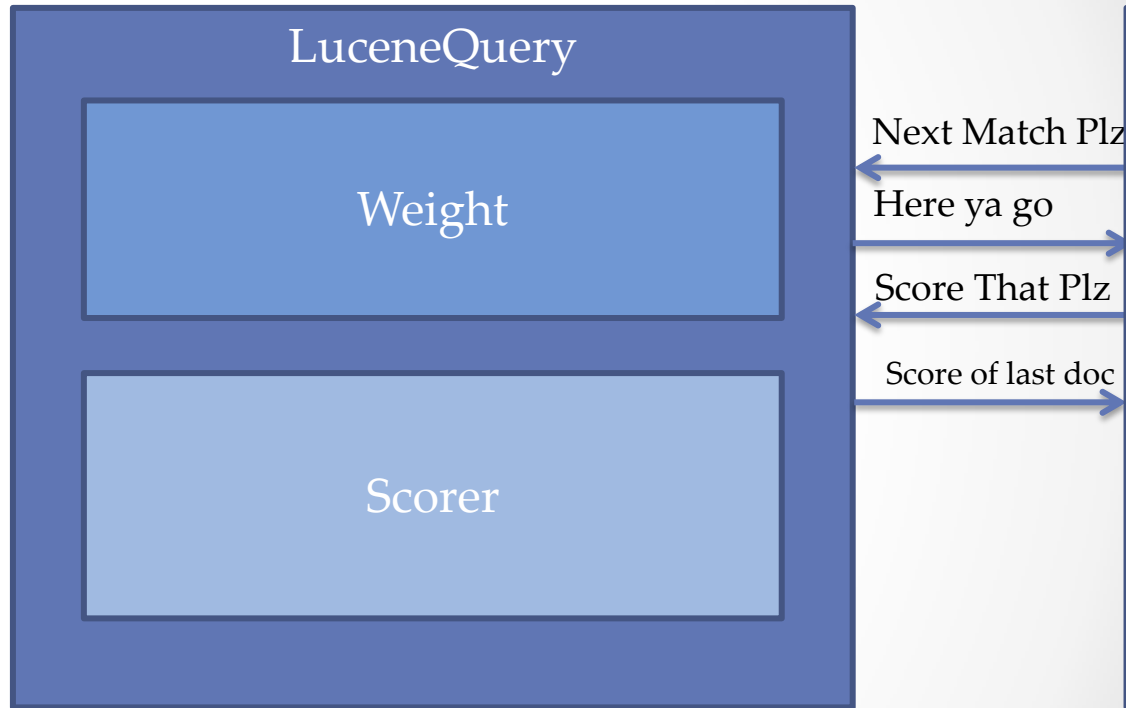- Disclaimer: its probably possible to do this with easier means!

https://github.com/o19s/lucene-query-example/

# Lucene Queries, 30,000 ft view

Aka, "not really accurate, but what to tell your boss to not confuse them"



LuceneQuery

IndexSearcher

Next Match Plz

Here ya go

Score That Plz

Score of last doc

Find next Match

Calc. score

IndexReader

# Anatomy of Lucene Query

A Tale Of Three Classes:

- *Queries* Create *Weights:*
  - Query-level stats for this search
  - Think "IDF" when you hear weights

- *Weights* Create *Scorers:*
  - Heavy Lifting, reports **matches** and returns a **score**

**LuceneQuery**

Weight

Scorer

Next Match Plz

Here ya go

Score That Plz

Score of last doc

Weight & Scorer are inner classes of Query

Find next Match

Calc. score

# Backwards Query Outline

```
class BacwkardsQuery {


    class BackwardsScorer {
        // matching & scoring functionality
    }


    class BackwardsWeight {
        // query normalization and other "global" stats

        public Scorer scorer(AtomicReaderContext context, …)
    }


    public Weight createWeight(IndexSearcher)

}
```

# How are these used?

When you do:

```
Query q = new BackwardsQuery();
idxSearcher.search(q);
```

This Setup Happens:

```
Weight w = q.createWeight(idxSearcher);
normalize(w);
foreach IndexReader idxReader:
    Scorer s = w.scorer(idxReader);
```

Important to know how Lucene is calling your code

# Weight

What should we do with our weight?

```
Weight w = q.createWeight(idxSearcher);
normalize(w);
```

**IndexSearcher Level Stats**
- Notice we pass the IndexSearcher when we create the weight
    - Weight tracks IndexSearcher level statistics used for scoring

**Query Normalization**
- Weight also participates in query normalization

Remember – its your Weight! **Weight can be a no-op and just create searchers**

# Weight & Query Normalization

Query Normalization – an optional little ritual to take your Weight instance through:

What I think my weight is

```
float v = weight.getValueForNormalization();
float norm = getSimilarity().queryNorm(v);
weight.normalize(norm, 1.0f);
```

Normalize that weight against global statistics

Pass back the normalized stats

# Weight & Query Normalization

```
float v = weight.getValueForNormalization();
float norm = getSimilarity().queryNorm(v);
weight.normalize(norm, 1.0f);
```

- For TermQuery:
  - The result of all this ceremony is the *IDF* (inverse document frequency of the term).

- This code is fairly abstract
  - All three steps are pluggable, and can be totally ignored

# BackwardsWeight

- Custom Weight that completely ignores query normalization:

```java
@Override
public float getValueForNormalization() throws IOException {
    return 0.0f;
}

@Override
public void normalize(float norm, float topLevelBoost) {
    // no-op
}
```
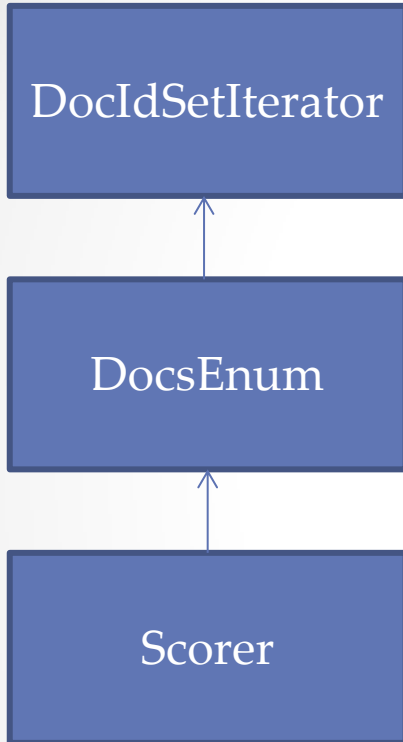
# Weights make Scorers!

```
@Override
public Scorer scorer(AtomicReaderContext context, boolean
scoreDocsInOrder, boolean topScorer, Bits acceptDocs) throws
IOException {

    return new BackwardsScorer(...);
}
```

- Scorers Have Two Jobs:
  - **Match!** – iterator interface over matching results
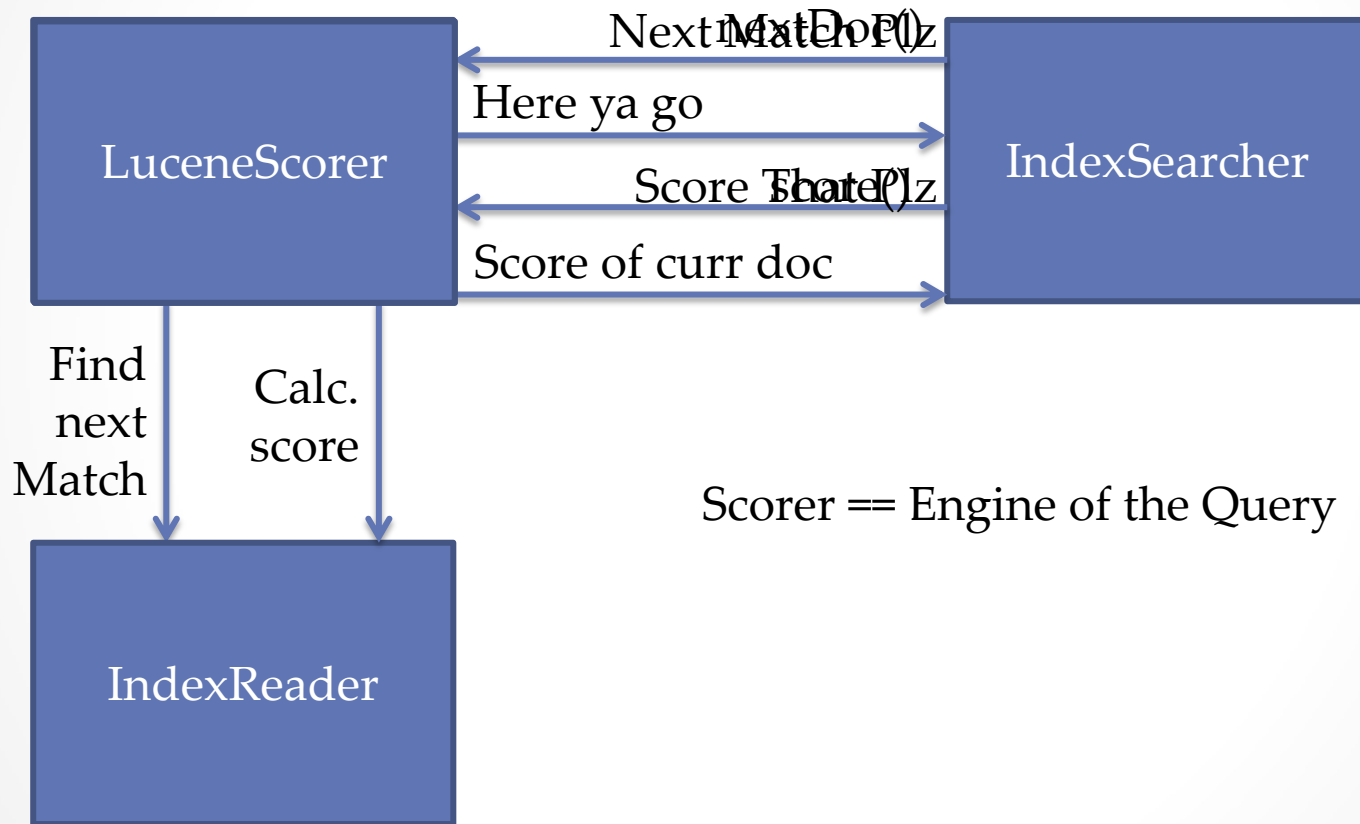  - **Score!** – score the current match

# Scorer as an iterator

DocIdSetIterator

DocsEnum

Scorer

- Inherits the following from DocsEnum:

- **nextDoc()**
  - Next match

- **advance(int docId)** –
  - Seek to the specified docId

- **docID()**
  - Id of the current document we're on

OpenSource Connections ●

# In other words…

- Remember THIS?          …Actually…

LuceneScorer

IndexSearcher

NextMatchPlz
Next Doc()

Here ya go

Score That()
Score Plz

Score of curr doc

Find next Match

Calc. score

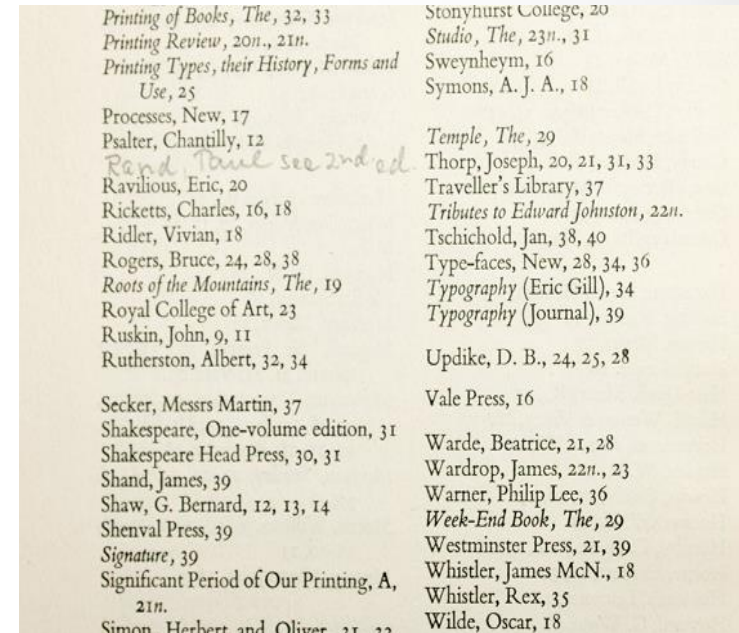IndexReader

Scorer == Engine of the Query

# What would nextDoc look like

- Remember search is an **inverted index**
  - Much like a book index
  - Fields -> Terms -> Documents!

IndexReader == our handle to inverted index:

- Much like an index. Given **term**, return **list of doc ids**
- **TermsEnum:**
  - Enumeration of terms (actual logical index of terms)
- **DocsEnum**
  - Enum. of corresponding docIDs (like list of pages next to term)
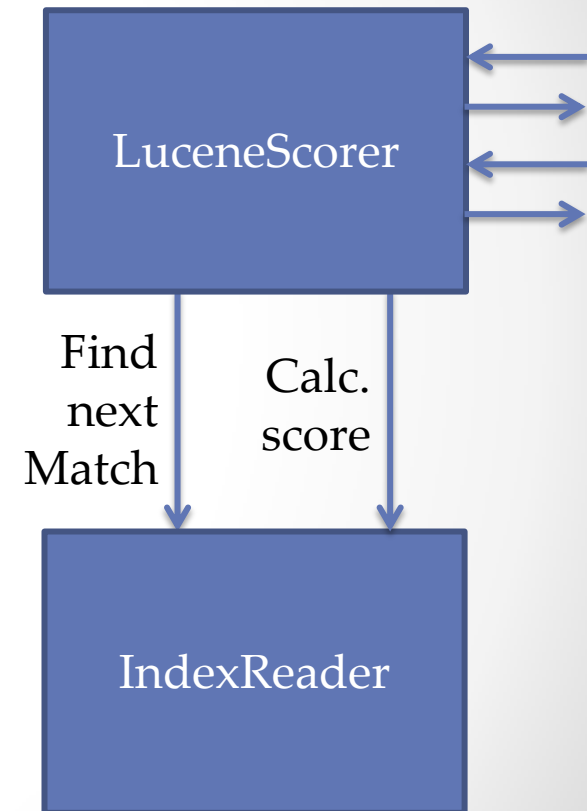
# What would nextDoc look like?

- TermsEnum to lookup info for a Term:

```
final TermsEnum termsEnum =
    reader.terms(term.field()).iterator(null);
termsEnum.seekExact(term.bytes(), state);
```

- Each term has a DocsEnum that lists the docs that contain this term:

```
DocsEnum docs = termsEnum.docs(acceptDocs, null);
```

LuceneScorer

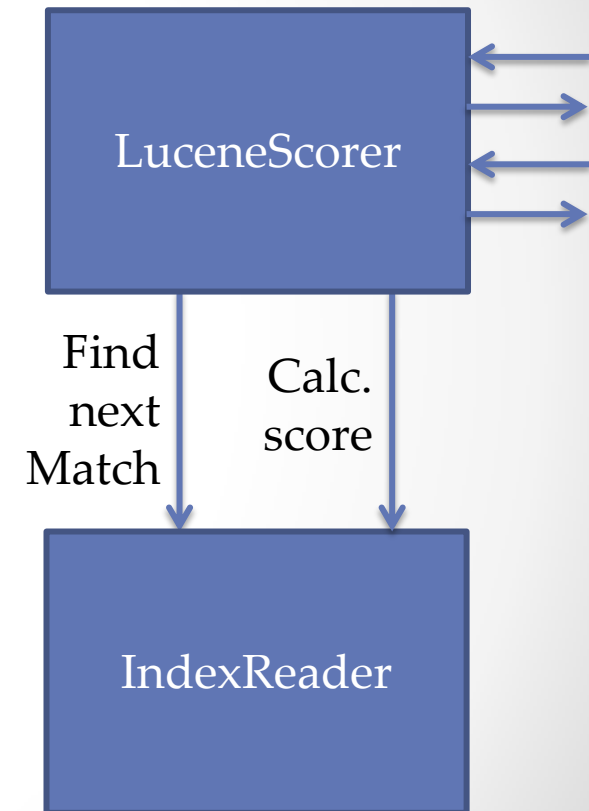Find next Match

Calc. score

IndexReader

# What would nextDoc look like?

- Wrapping this enum, now I can return matches for this term!

```
@Override
 public int nextDoc() throws IOException {
   return docs.nextDoc();
 }
```

- You've just implemented TermQuery!

LuceneScorer

Find next Match

Calc. score

IndexReader

# BackwardsScorer nextDoc

- Recall our Query has a Backwards Term (ananab):

```
public BackwardsQuery(String field, String term) {
    backwardsTerm = new Term(field, new StringBuilder(term).reverse().toString());
    ...
}
```

- Later, when creating a Scorer. Get a handle to **DocsEnum** for our backwards term:

```
public Scorer scorer(AtomicReaderContext context, boolean scoreDocsInOrder,
                     boolean topScorer, Bits acceptDocs) throws IOException {

    Term bwdsTerm = BackwardsQuery.this.backwardsTerm;
    TermsEnum bwdsTerms = context.reader().terms(bwdsTerm.field()).iterator(null);
    bwdsTerms.seekExact(bwdsTerm.bytes());
    DocsEnum bwdsDocs = bwdsTerms.docs(acceptDocs, null);
```

Terrifying and verbose Lucene speak for:
1. Seek to term in field via TermsEnum
2. Give me a DocsEnum of matching docs
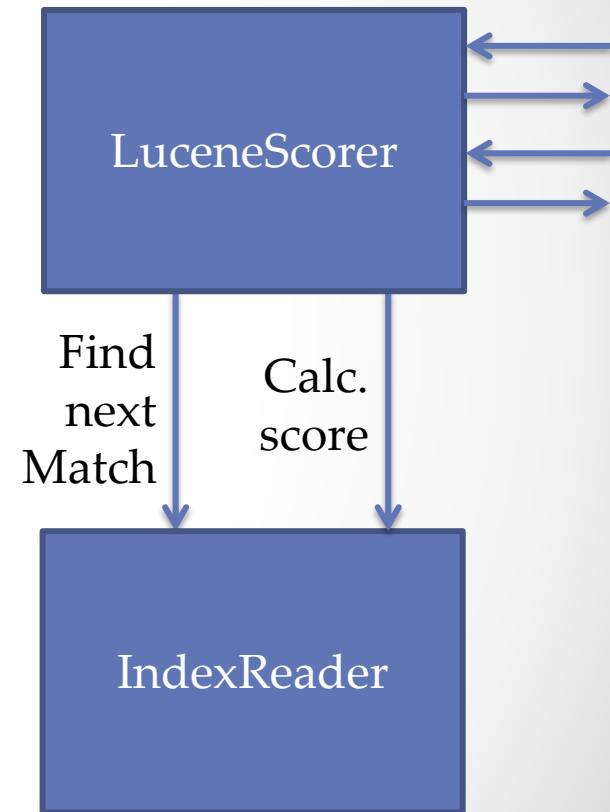
# BackwardsScorer nextDoc

- Our scorer has **bwdDocs** and **fwdDocs**, our nextDoc just walks both:

```
@Override
public int nextDoc() throws IOException {
    int currDocId = docID();
    // increment one or both
    if (currDocId == backwardsScorer.docID()) {
        backwardsScorer.nextDoc();
    }
    if (currDocId == forwardsScorer.docID()) {
        forwardsScorer.nextDoc();
    }
    return docID();
}
```

# Scorer for scores!

- Score is easy! Implement **score**, do whatever you want!

```
@Override
public float score() throws
IOException {
    return 1.0f;
}
```

LuceneScorer

Find next Match

Calc. score

IndexReader

# BackwardsScorer Score

- Recall, match a backwards term (ananab)score = 5.0, fwd term (banana) score = 1.0
- We hook into docID, update score based on current posn

We call docID() in nextDoc()

```
@Override
public int docID() {
    int backwordsDocId = backwardsScorer.docID();
    int forwardsDocId = forwardsScorer.docID();
    if (backwordsDocId <= forwardsDocId && backwordsDocId != NO_MORE_DOCS) {
        currScore = BACKWARDS_SCORE;
        return backwordsDocId;
    } else if (forwardsDocId != NO_MORE_DOCS) {
        currScore = FORWARDS_SCORE;
        return forwardsDocId;
    }
    return NO_MORE_DOCS;
}
```
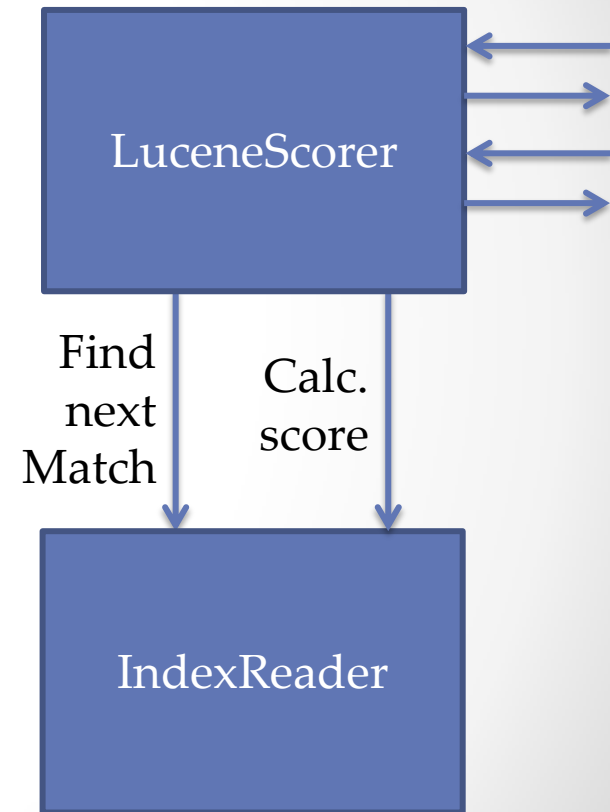
Currently positioned on a bwds doc, set currScore to 5.0

Currently positioned on a fwd doc, set currScore to 1.0

# BackwardsScorer Score

- For completeness sake, here's our **score**:

```
@Override
public float score() throws
IOException {
    return currScore;
}
```

```
┌─────────────────┐  ←──
│                 │  ──→
│  LuceneScorer   │  ←──
│                 │  ──→
└─────────────────┘
    │         │
 Find      Calc.
 next      score
 Match       │
    ↓         ↓
┌─────────────────┐
│                 │
│   IndexReader   │
│                 │
└─────────────────┘
```

# So many gotchas!

- Ultimate POWER! But You will have weird bugs:

  - Do all of your searches return the results of your first query?
    - In **Query** Implement **hashCode** and **equals**
  - Weird/Random Test Failures
    - Test using LuceneTestCase to ferret out common Lucene bugs
      - Randomized testing w/ different codecs etc
      - IndexReader methods have a certain ritual and very specific rules, (enums must be primed, etc)

# Extras

- Query **rewrite** method
  - Optional, recognize you are a complex query, turn yourself into a simpler one
    - BooleanQuery with 1 clause -> return just one clause
- Weight has optional **explain**
  - Useful for debugging in Solr
  - Pretty straight-forward API

# Conclusions!

- These are nuclear options!
  - You can achieve SO MUCH before you get here (at much less complexity)
  - There's certainly a way to do what you've seen without this level of control

- Fun way to learn about Lucene!

# QUESTIONS?