

# INDEXING AND SEARCHING RDF DATASETS

Improving Performance of Semantic Web Applications with  
Lucene, SIREn and RDF

Mike Hugo  
Entagen, LLC

slides and sample code can be found at

<https://github.com/mjhugo/rdf-lucene-siren-presentation>

# AGENDA

- The problem
- RDF / NTriples overview
- Lucene
- SIREn
- Searching RDF with SIREn
- Questions

# THE PROBLEM

- Build a web application utilizing RDF data that can:
  - Provide responsive search results and/or “auto-complete” or “type-ahead” functionality
  - Query literals in a case in-sensitive manner
  - Find relationships between entities (and counts)
  - Scale well over millions of triples
  - Provide a responsive user interface

# WHAT'S A TRIPLE?

subject

predicate

object

# WHAT'S A TRIPLE?

subject

“Imatinib”

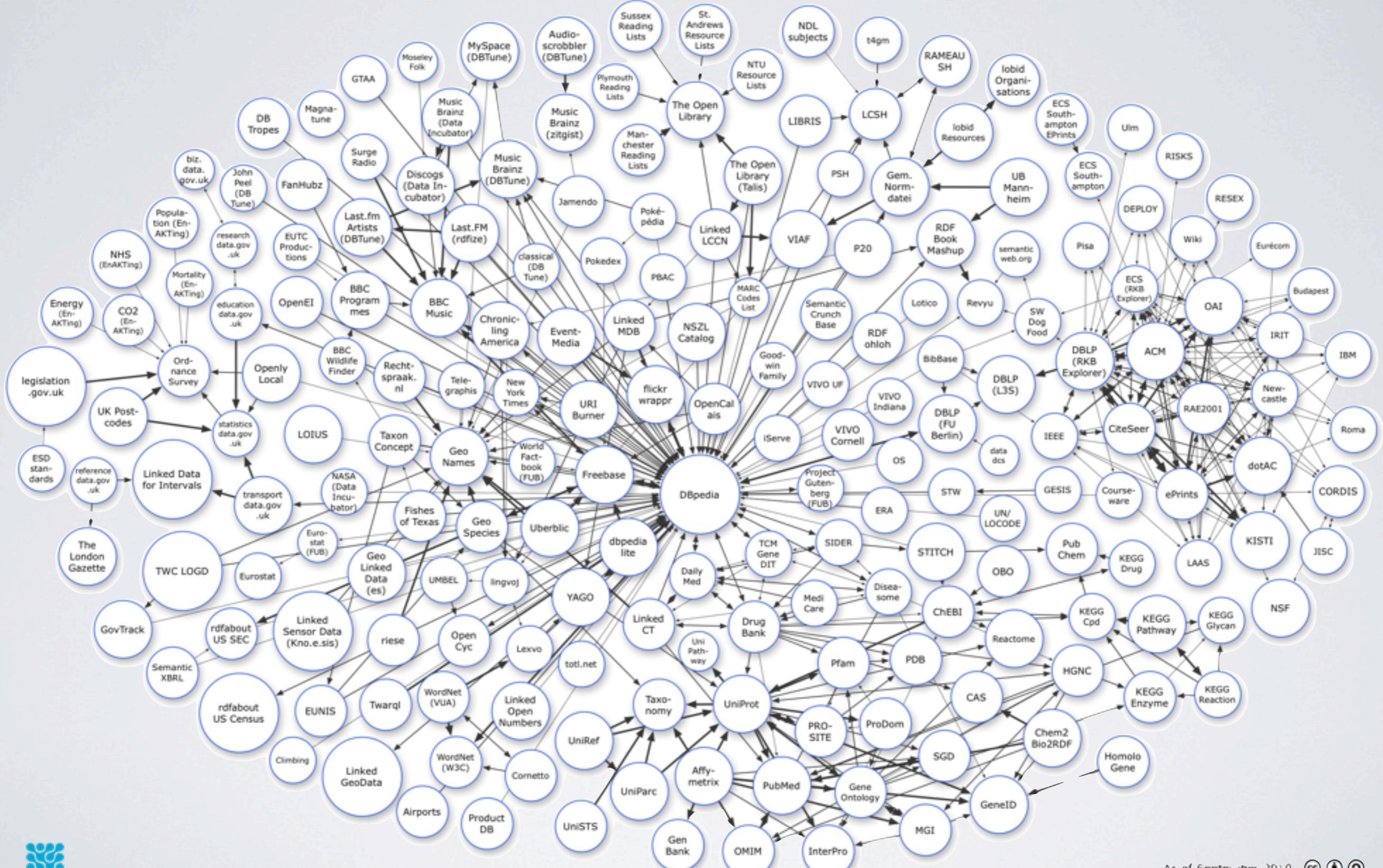
predicate

type

object

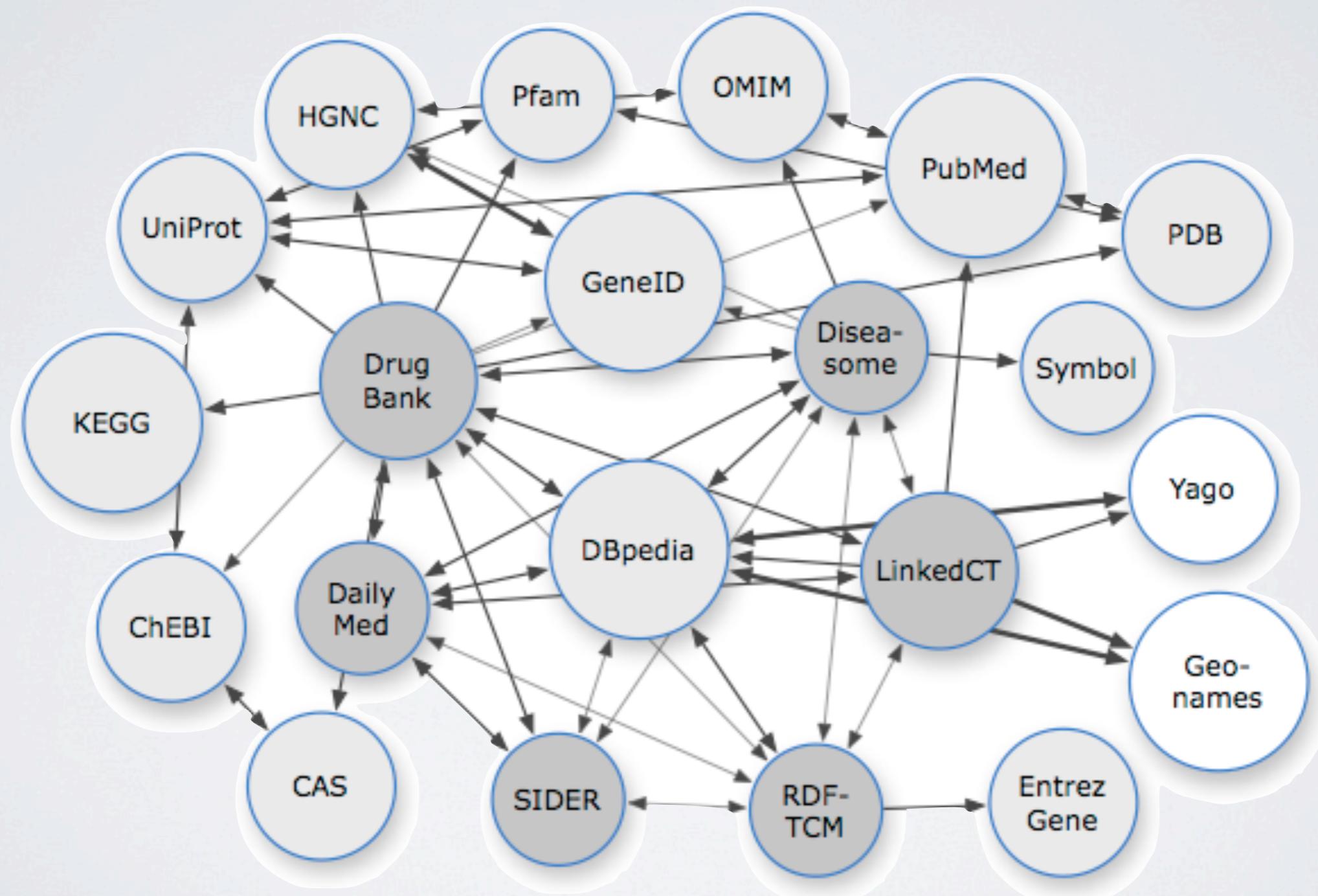
Pharmaceutical Drug

# LINKING OPEN DATA



# LINKING OPEN DRUG DATA

<http://esw.w3.org/HCLSIG/LODD>



# WHAT'S A TRIPLE?

<http://www4.wiwiiss.fu-berlin.de/drugbank/resource/drugs/DB00619>

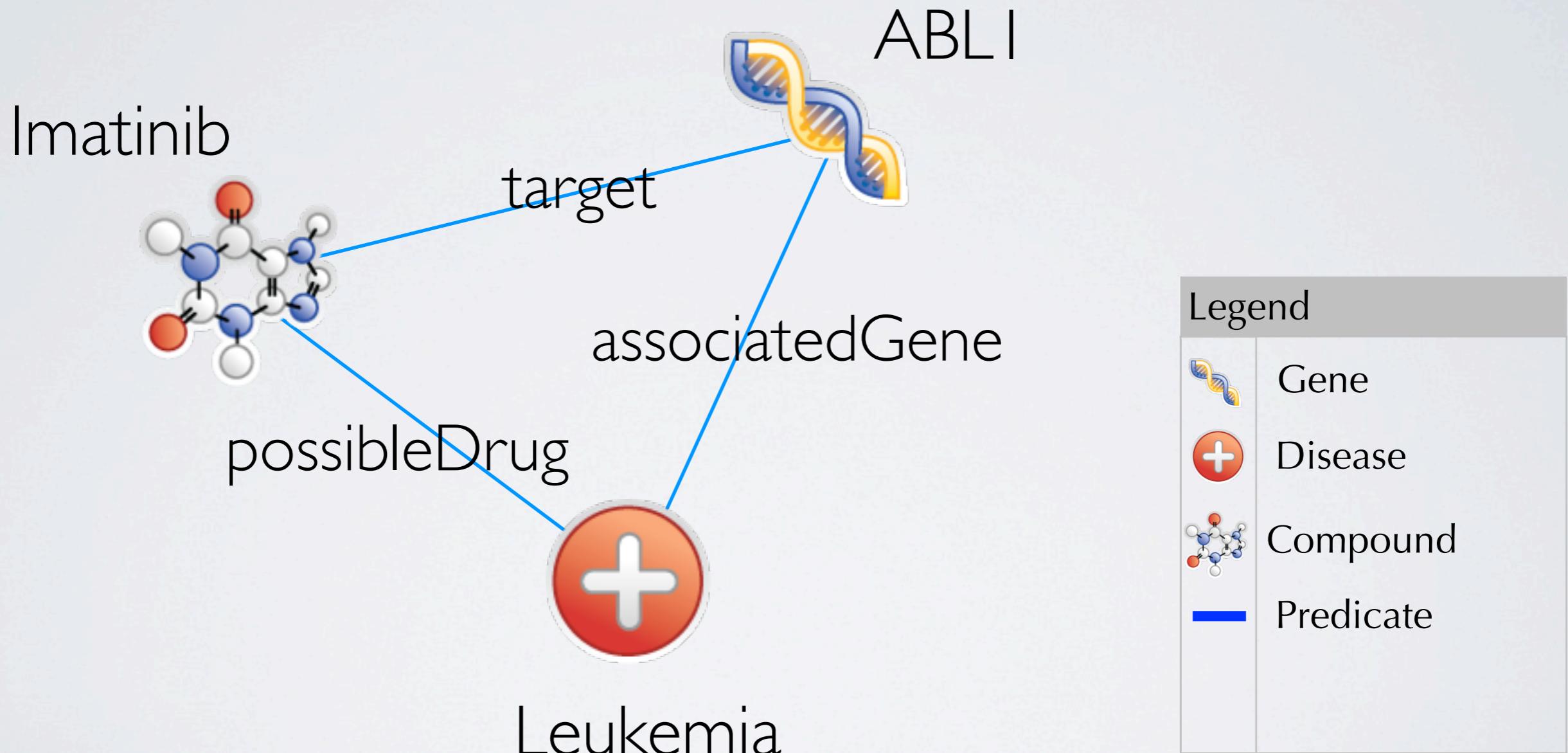
<b>subject</b>	<b>predicate</b>	<b>object</b>
<u>DB00619</u>	label	<b>“Imatinib”</b>
<u>DB00619</u>	type	drugbank:drugs
<u>DB00619</u>	brandName	<b>“Gleevec”</b>
<u>DB00619</u>	drugbank:target	<u>targets:17</u>

# WHAT'S A TRIPLE?

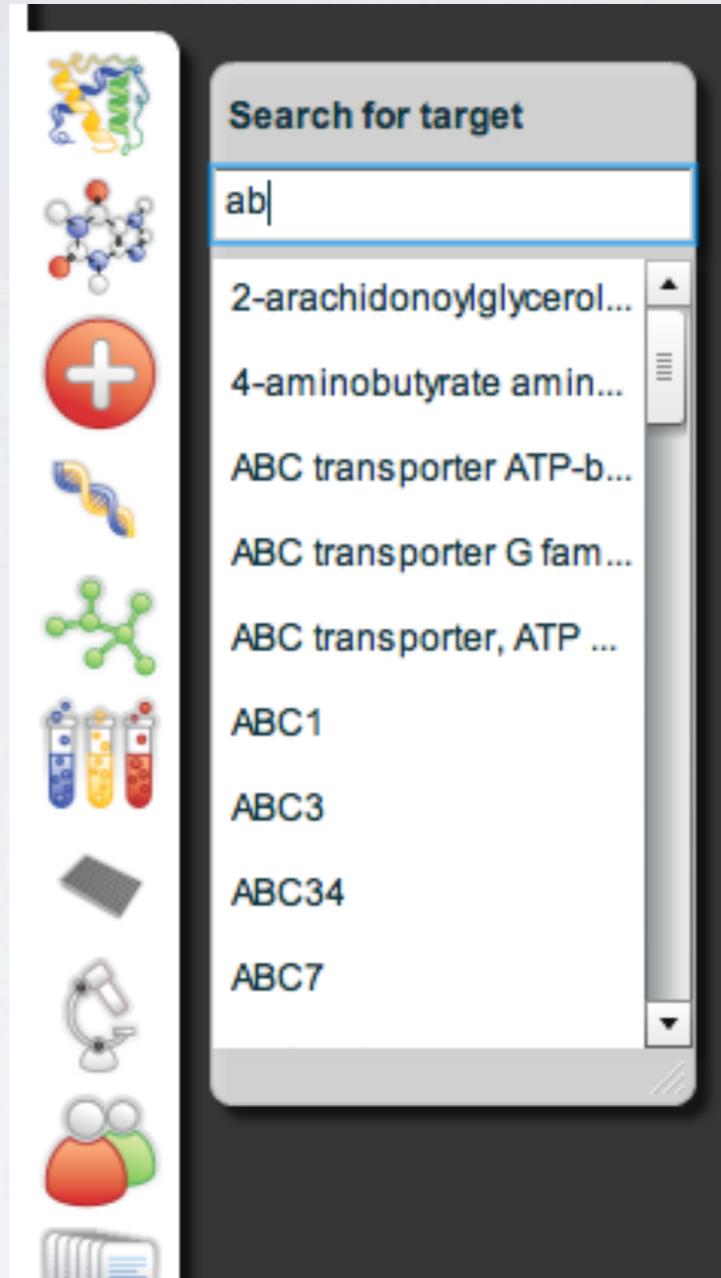
<http://www4.wiwiss.fu-berlin.de/drugbank/resource/targets/17>

<b>subject</b>	<b>predicate</b>	<b>object</b>
<u>DB00619</u>	drugbank:target	<u>targets/17</u>
<u>targets/17</u>	type	drugbank:drugs
<u>targets/17</u>	label	“Proto-oncogene tyrosine-protein kinase ABL1”
<u>targets/17</u>	geneName	“ABL1”

# RELATIONSHIPS



# FAST SEARCHING



```
select id, label  
from targets  
where label ilike '%${queryValue}%'
```

```
select id, label  
from targets  
where label like '%${queryValue}%'
```

```
SELECT ?uri ?label WHERE {  
  ?uri rdfs:label ?label .  
  ?uri rdf:type drugbank:targets .  
  FILTER regex(?label, '\\\\\\Q${queryValue}\\\\E', 'i')  
}
```

query as literal value

case insensitive

```
select id, label  
from targets  
where label like '%${queryValue}%'
```

```
SELECT ?uri ?label WHERE {  
?uri rdfs:label ?label .  
?uri rdf:type drugbank:targets .  
FILTER regex(?label, '\\\\\\Q${queryValue}\\E', 'i')  
}
```

# multiple datasets

```
select id, label  
from targets  
where label like '%${queryValue}%'
```

```
SELECT ?uri ?label WHERE {  
  ?uri rdfs:label ?label .  
  ?uri rdf:type drugbank:targets .  
  FILTER regex(?label, '\\\\\\Q${queryValue}\\\\E', 'i')  
}
```

**multiple datasets  
millions of triples**

# DEMO

Baseline SPARQL Query Performance

**FASTER!**

```
String queryString
if (params.type == 'Exact') {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER (?label = '${params.query}')
        } LIMIT 10
    """
} else {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER regex(?label, '\\\\Q${params.query}\\\\E', 'i')
        } LIMIT 10
    """
}

def s = new Date().time
List results = sparqlQueryService.executeQuery(queryString)
def e = new Date().time
```

```
String queryString
if (params.type == 'Exact') {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER (?label = '${params.query}')
        } LIMIT 10
    """
} else {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER regex(?label, '\\\\Q${params.query}\\\\E', 'i')
        } LIMIT 10
    """
}

def s = new Date().time
List results = sparqlQueryService.executeQuery(queryString)
def e = new Date().time
```

```
String queryString
if (params.type == 'Exact') {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER (?label = '${params.query}')
        } LIMIT 10
    """
} else {
    queryString = """
        SELECT ?uri ?type ?label WHERE {
            ?uri rdfs:label ?label .
            ?uri rdf:type ?type .
            FILTER regex(?label, '\\\\Q${params.query}\\\\E', 'i')
        } LIMIT 10
    """
}

def s = new Date().time
List results = sparqlQueryService.executeQuery(queryString)
def e = new Date().time
```



- Java API for adding search and indexing capability to applications
- Widely used in many applications
- Works well with structured and unstructured data
- Complements existing storage layers (RDBS, Triple Store, etc.).



- Index of Content
- Document
- Fields

## Sample Document

<b>field</b>	<b>value</b>
name	“mike hugo”
company	“Entagen”
bio	“lorem ipsum...”



- Searching
- Keyword
  - name = mike
- Phrases
  - name = “Mike Hugo”
- Wildcards
  - name = “Mike H\*”

# LUCENE INDEX

- First pass - a lucene index of searchable labels
- For each RDF entity, create a lucene document with fields:
  - subject uri
  - label
  - rdf:type

# DEMO

Lucene Index of Searchable Labels

```
String queryLabels = """
    SELECT ?uri ?label
    WHERE {
        ?uri rdfs:label ?label .
    }
"""

sparqlQueryService.executeForEach(repository, queryLabels) {
    def doc = new Document()
    String uri = it.uri.stringValue()
    String label = it.label.stringValue()

    doc.add(new Field(SUBJECT_URI_FIELD, uri,
                      Field.Store.YES, Field.Index.ANALYZED))
    doc.add(new Field(LABEL_FIELD, label,
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
String queryLabels = """
    SELECT ?uri ?label
    WHERE {
        ?uri rdfs:label ?label .
    }
"""

sparqlQueryService.executeForEach(repository, queryLabels) {
    def doc = new Document()
    String uri = it.uri.stringValue()
    String label = it.label.stringValue()

    doc.add(new Field(SUBJECT_URI_FIELD, uri,
                      Field.Store.YES, Field.Index.ANALYZED))
    doc.add(new Field(LABEL_FIELD, label,
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
String queryLabels = """
    SELECT ?uri ?label
    WHERE {
        ?uri rdfs:label ?label .
    }
"""

sparqlQueryService.executeForEach(repository, queryLabels) {
    def doc = new Document()
    String uri = it.uri.stringValue()
    String label = it.label.stringValue()

    doc.add(new Field(SUBJECT_URI_FIELD, uri,
                      Field.Store.YES, Field.Index.ANALYZED))
    doc.add(new Field(LABEL_FIELD, label,
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
def query = {
    Query query = new QueryParser(
        Version.LUCENE_CURRENT,
        LABEL_FIELD,
        new StandardAnalyzer(Version.LUCENE_CURRENT)).parse(params.query);

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}

public List executeQuery(Query query) {
    IndexSearcher searcher = luceneSearcherManager.get()
    ScoreDoc[] scoreDocs = searcher.search(query, 10).scoreDocs
    List results = []
    def connection = repository.connection
    scoreDocs.each {
        Document doc = searcher.doc(it.doc)
        String uri = doc[SUBJECT_URI_FIELD]
        Map labelAndType = sparqlQueryService.getLabelAndType(uri, connection)
        results << [uri: uri, type: labelAndType.type, label: labelAndType.label]
    }
    connection.close()
    return results
}
```

```
def query = {
    Query query = new QueryParser(
        Version.LUCENE_CURRENT,
        LABEL_FIELD,
        new StandardAnalyzer(Version.LUCENE_CURRENT)).parse(params.query);

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}

public List executeQuery(Query query) {
    IndexSearcher searcher = luceneSearcherManager.get()
    ScoreDoc[] scoreDocs = searcher.search(query, 10).scoreDocs
    List results = []
    def connection = repository.connection
    scoreDocs.each {
        Document doc = searcher.doc(it.doc)
        String uri = doc[SUBJECT_URI_FIELD]
        Map labelAndType = sparqlQueryService.getLabelAndType(uri, connection)
        results << [uri: uri, type: labelAndType.type, label: labelAndType.label]
    }
    connection.close()
    return results
}
```

```
def query = {
    Query query = new QueryParser(
        Version.LUCENE_CURRENT,
        LABEL_FIELD,
        new StandardAnalyzer(Version.LUCENE_CURRENT)).parse(params.query);

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}

public List executeQuery(Query query) {
    IndexSearcher searcher = luceneSearcherManager.get()
    ScoreDoc[] scoreDocs = searcher.search(query, 10).scoreDocs
    List results = []
    def connection = repository.connection
    scoreDocs.each {
        Document doc = searcher.doc(it.doc)
        String uri = doc[SUBJECT_URI_FIELD]
        Map labelAndType = sparqlQueryService.getLabelAndType(uri, connection)
        results << [uri: uri, type: labelAndType.type, label: labelAndType.label]
    }
    connection.close()
    return results
}
```

```
def query = {
    Query query = new QueryParser(
        Version.LUCENE_CURRENT,
        LABEL_FIELD,
        new StandardAnalyzer(Version.LUCENE_CURRENT)).parse(params.query);

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}

public List executeQuery(Query query) {
    IndexSearcher searcher = luceneSearcherManager.get()
    ScoreDoc[] scoreDocs = searcher.search(query, 10).scoreDocs
    List results = []
    def connection = repository.connection
    scoreDocs.each {
        Document doc = searcher.doc(it.doc)
        String uri = doc[SUBJECT_URI_FIELD]
        Map labelAndType = sparqlQueryService.getLabelAndType(uri, connection)
        results << [uri: uri, type: labelAndType.type, label: labelAndType.label]
    }
    connection.close()
    return results
}
```

WHAT ABOUT ENTITY  
RELATIONSHIPS?

# SIREn

{ *Semantic Information Retrieval Engine* }

- Java API / Lucene Extension
- Specifically built to handle semi-structured data like RDF
- Provides a special Lucene Analyzer (TupleAnalyzer) to index NTriples formatted RDF data
- Can also handle other types of semi-structured data (not limited to RDF)
- Can also work with Solr (Lucene search server)

## Sample (Lucene)Document

field	value
uri	<u><a href="http://www4.wiwiss.fu-berlin.de/drugbank/resource/drugs/DB00619">http://www4.wiwiss.fu-berlin.de/drugbank/resource/drugs/DB00619</a></u>
ntriples	<DB00619> rdfs:label "Imatinib". <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevac" . <DB00619> drugbank:target <targets/1588> .

## Querying Triples

subject	predicate	object
<u>DB00619</u>	rdfs:label	<b>“Imatinib”</b>

```
SirenTupleQuery tupleQuery = new SirenTupleQuery()
tupleQuery.add(createCellQuery('rdfs:label',
  SirenTupleConstraintPosition.PREDICATE) ,
  SirenTupleClause.Occur.MUST)
tupleQuery.add(createCellQuery('imatinib',
  SirenTupleConstraintPosition.OBJECT) ,
  SirenTupleClause.Occur.MUST)
```

# DEMO

SIREn Index of RDF Entities

```
String subjectUris = """
    SELECT distinct ?uri
    WHERE {
        ?uri ?p ?o .
    }
"""

sparqlQueryService.executeForEach(repository, subjectUris) {
    def doc = new Document()

    String subjectUri = it.uri.stringValue()
    doc.add(new Field(SUBJECT_URI_FIELD, subjectUri,
                      Field.Store.YES, Field.Index.ANALYZED))

    StringWriter triplesStringWriter = new StringWriter()
    NTriplesWriter nTriplesWriter = new NTriplesWriter(triplesStringWriter)
    connection.exportStatements(new URIImpl(subjectUri),
                                null, null, false, nTriplesWriter)

    doc.add(new Field(TRIPLES_FIELD, triplesStringWriter.toString(),
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
String subjectUris = """
    SELECT distinct ?uri
    WHERE {
        ?uri ?p ?o .
    }
"""

sparqlQueryService.executeForEach(repository, subjectUris) {
    def doc = new Document()

    String subjectUri = it.uri.stringValue()
    doc.add(new Field(SUBJECT_URI_FIELD, subjectUri,
                      Field.Store.YES, Field.Index.ANALYZED))

    StringWriter triplesStringWriter = new StringWriter()
    NTriplesWriter nTriplesWriter = new NTriplesWriter(triplesStringWriter)
    connection.exportStatements(new URIImpl(subjectUri),
                                null, null, false, nTriplesWriter)

    doc.add(new Field(TRIPLES_FIELD, triplesStringWriter.toString(),
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
String subjectUris = """
    SELECT distinct ?uri
    WHERE {
        ?uri ?p ?o .
    }
"""

sparqlQueryService.executeForEach(repository, subjectUris) {
    def doc = new Document()

    String subjectUri = it.uri.stringValue()
    doc.add(new Field(SUBJECT_URI_FIELD, subjectUri,
                      Field.Store.YES, Field.Index.ANALYZED))

    StringWriter triplesStringWriter = new StringWriter()
    NTriplesWriter nTriplesWriter = new NTriplesWriter(triplesStringWriter)
    connection.exportStatements(new URIImpl(subjectUri),
                                null, null, false, nTriplesWriter)

    doc.add(new Field(TRIPLES_FIELD, triplesStringWriter.toString(),
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
String subjectUris = """
    SELECT distinct ?uri
    WHERE {
        ?uri ?p ?o .
    }
"""

sparqlQueryService.executeForEach(repository, subjectUris) {
    def doc = new Document()

    String subjectUri = it.uri.stringValue()
    doc.add(new Field(SUBJECT_URI_FIELD, subjectUri,
                      Field.Store.YES, Field.Index.ANALYZED))

    StringWriter triplesStringWriter = new StringWriter()
    NTriplesWriter nTriplesWriter = new NTriplesWriter(triplesStringWriter)
    connection.exportStatements(new URIImpl(subjectUri),
                                null, null, false, nTriplesWriter)

    doc.add(new Field(TRIPLES_FIELD, triplesStringWriter.toString(),
                      Field.Store.NO, Field.Index.ANALYZED))

    writer.addDocument(doc)
}
```

```
def query = {
    String userQuery = params.query.toLowerCase()
    SirenCellQuery predicateCellQuery = new SirenCellQuery(
        new SirenTermQuery(new Term(TRIPLES_FIELD,
            RDFS.LABEL.stringValue())))
    predicateCellQuery.constraint = PREDICATE_CELL

    SirenTermQuery objectQuery =
        new SirenTermQuery(new Term(TRIPLES_FIELD, userQuery))
    SirenCellQuery objectCellQuery = new SirenCellQuery(objectQuery)
    objectCellQuery.constraint = OBJECT_CELL

    Query query = new SirenTupleQuery()
    query.add(predicateCellQuery, SirenTupleClause.Occur.MUST)
    query.add(objectCellQuery, SirenTupleClause.Occur.MUST)

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}
```

```
def query = {
    String userQuery = params.query.toLowerCase()
    SirenCellQuery predicateCellQuery = new SirenCellQuery(
        new SirenTermQuery(new Term(TRIPLES_FIELD,
            RDFS.LABEL.stringValue())))
    predicateCellQuery.constraint = PREDICATE_CELL

    SirenTermQuery objectQuery =
        new SirenTermQuery(new Term(TRIPLES_FIELD, userQuery))
    SirenCellQuery objectCellQuery = new SirenCellQuery(objectQuery)
    objectCellQuery.constraint = OBJECT_CELL

    Query query = new SirenTupleQuery()
    query.add(predicateCellQuery, SirenTupleClause.Occur.MUST)
    query.add(objectCellQuery, SirenTupleClause.Occur.MUST)

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}
```



```
def query = {
    String userQuery = params.query.toLowerCase()
    SirenCellQuery predicateCellQuery = new SirenCellQuery(
        new SirenTermQuery(new Term(TRIPLES_FIELD,
            RDFS.LABEL.stringValue())))
    predicateCellQuery.constraint = PREDICATE_CELL

    SirenTermQuery objectQuery =
        new SirenTermQuery(new Term(TRIPLES_FIELD, userQuery))
    SirenCellQuery objectCellQuery = new SirenCellQuery(objectQuery)
    objectCellQuery.constraint = OBJECT_CELL

    Query query = new SirenTupleQuery()
    query.add(predicateCellQuery, SirenTupleClause.Occur.MUST)
    query.add(objectCellQuery, SirenTupleClause.Occur.MUST)

    def s = new Date().time
    List results = executeQuery(query)
    def e = new Date().time

    render(view: 'index', model: [results: results, time: e - s])
}
```

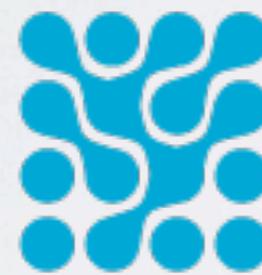


# DEMO

SIREn in action on [TripleMap.com](http://TripleMap.com)

# QUESTIONS?

[mike@entagen.com](mailto:mike@entagen.com) / twitter: @piragua



**ENTAGEN**  
**ACCELERATING INSIGHT**

<http://www.entagen.com>



**TripleMap**

<http://www.triplemap.com>