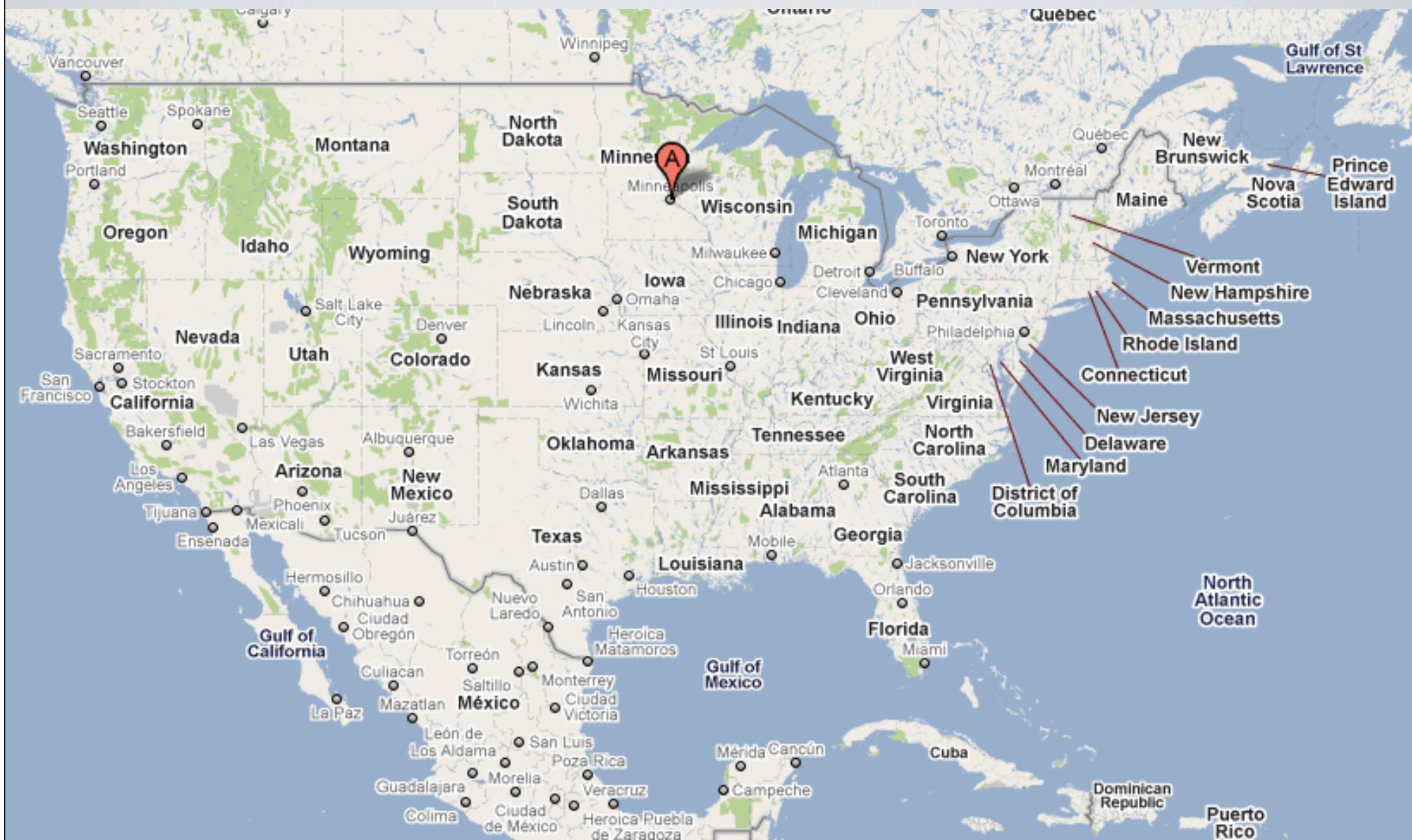


INDEXING AND SEARCHING RDF DATASETS

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

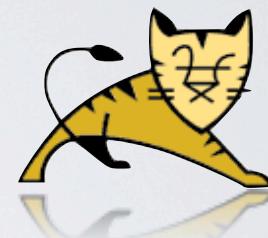
Mike Hugo
Entagen, LLC

slides and sample code can be found at
<https://github.com/mjhugo/rdf-lucene-siren-presentation>





ORACLE®



ENTAGEN OVERVIEW

- Founded in May 2008
- 17 Team Members located in two onshore offices: Boston & Minneapolis
- Life sciences & biomedical expertise, extensive network and focused technical specialization in enterprise class search, integration and analysis systems
- Custom Software Development
- Proprietary Software Applications: TripleMap, BrightScan, ExaSeq
- Professional Services/Strategic Consulting

SELECT ENTAGEN CLIENTS

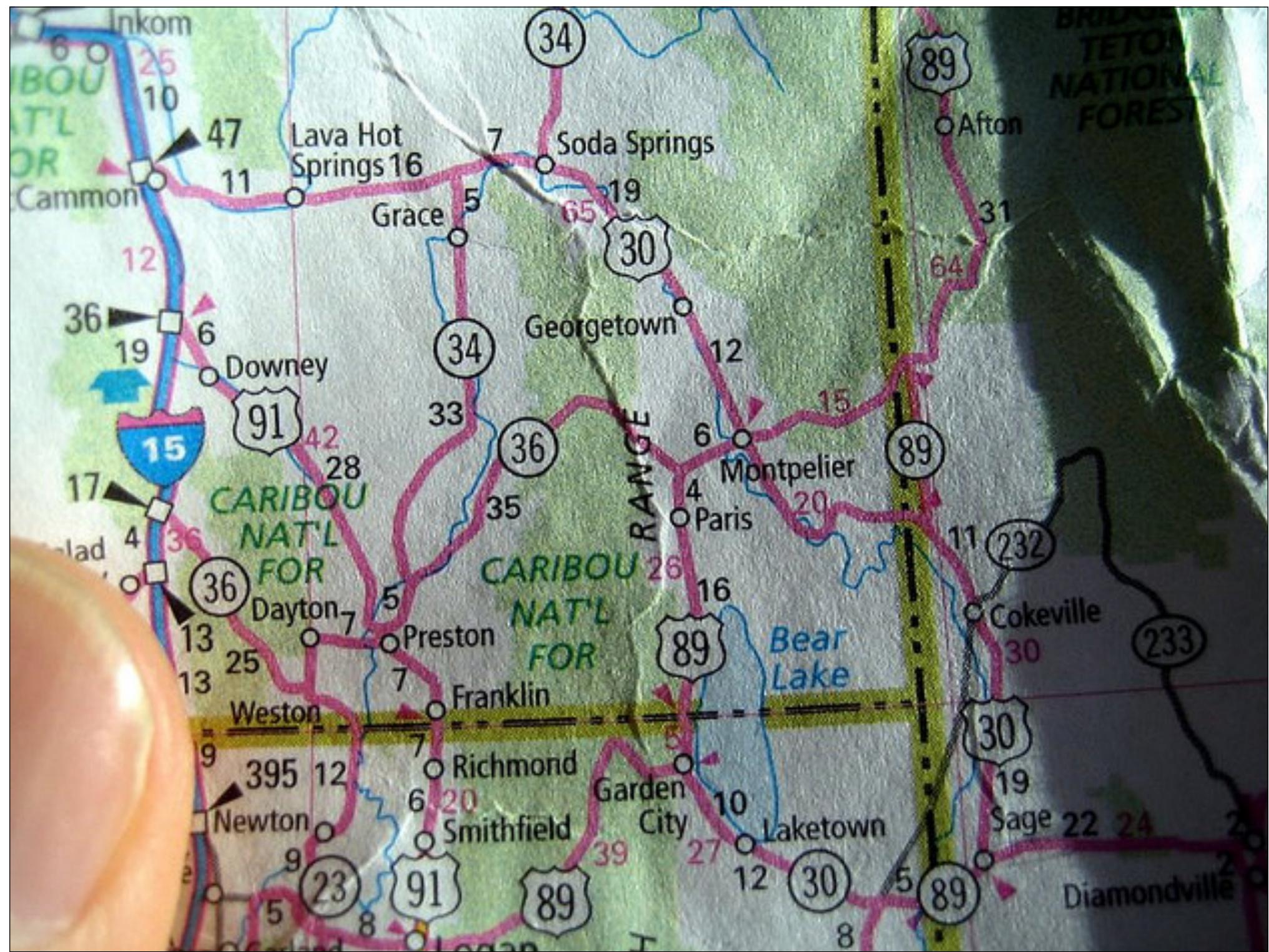


Boehringer
Ingelheim



DANA-FARBER
CANCER INSTITUTE

AGENDA



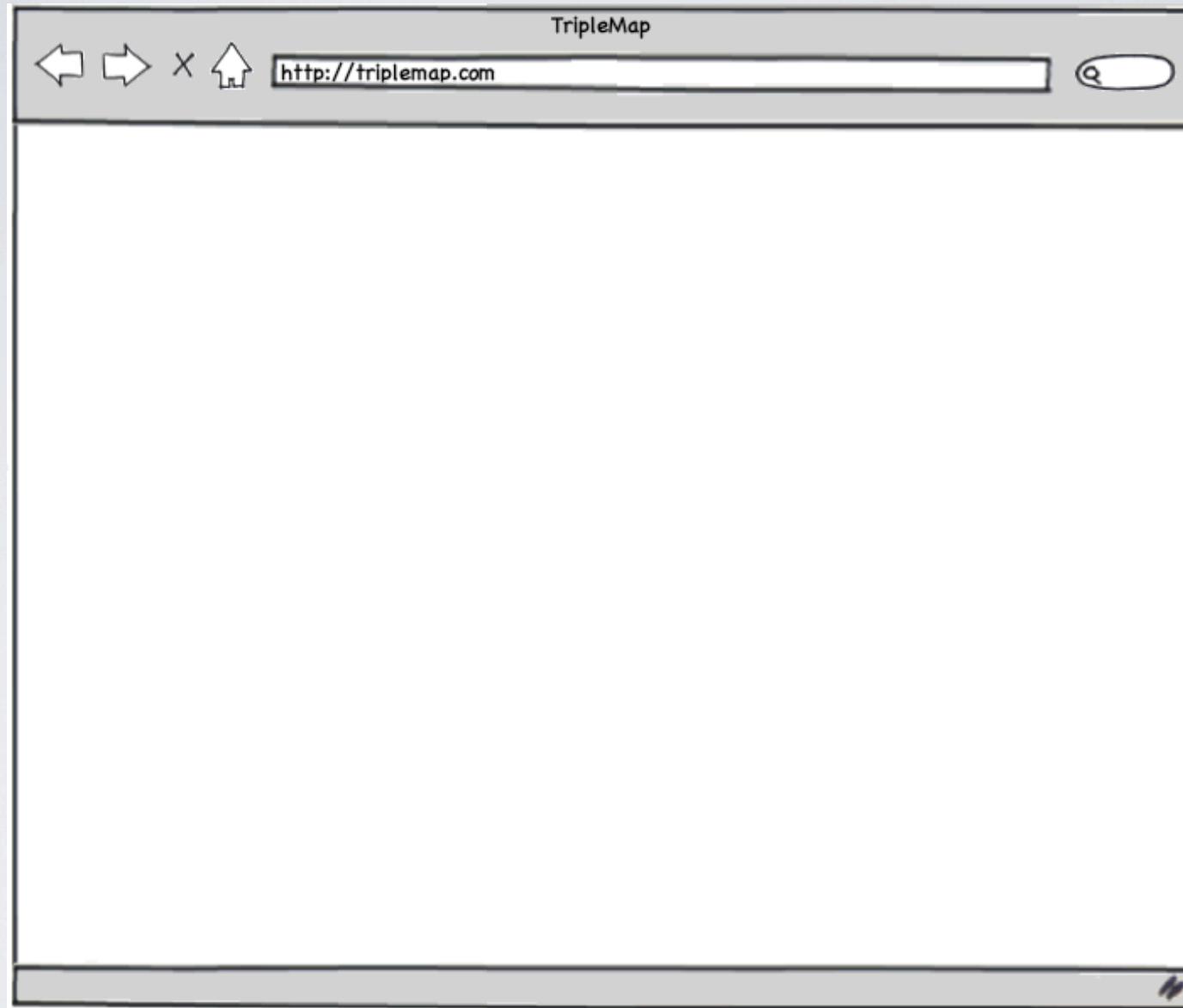
SPARQL

LUCENE

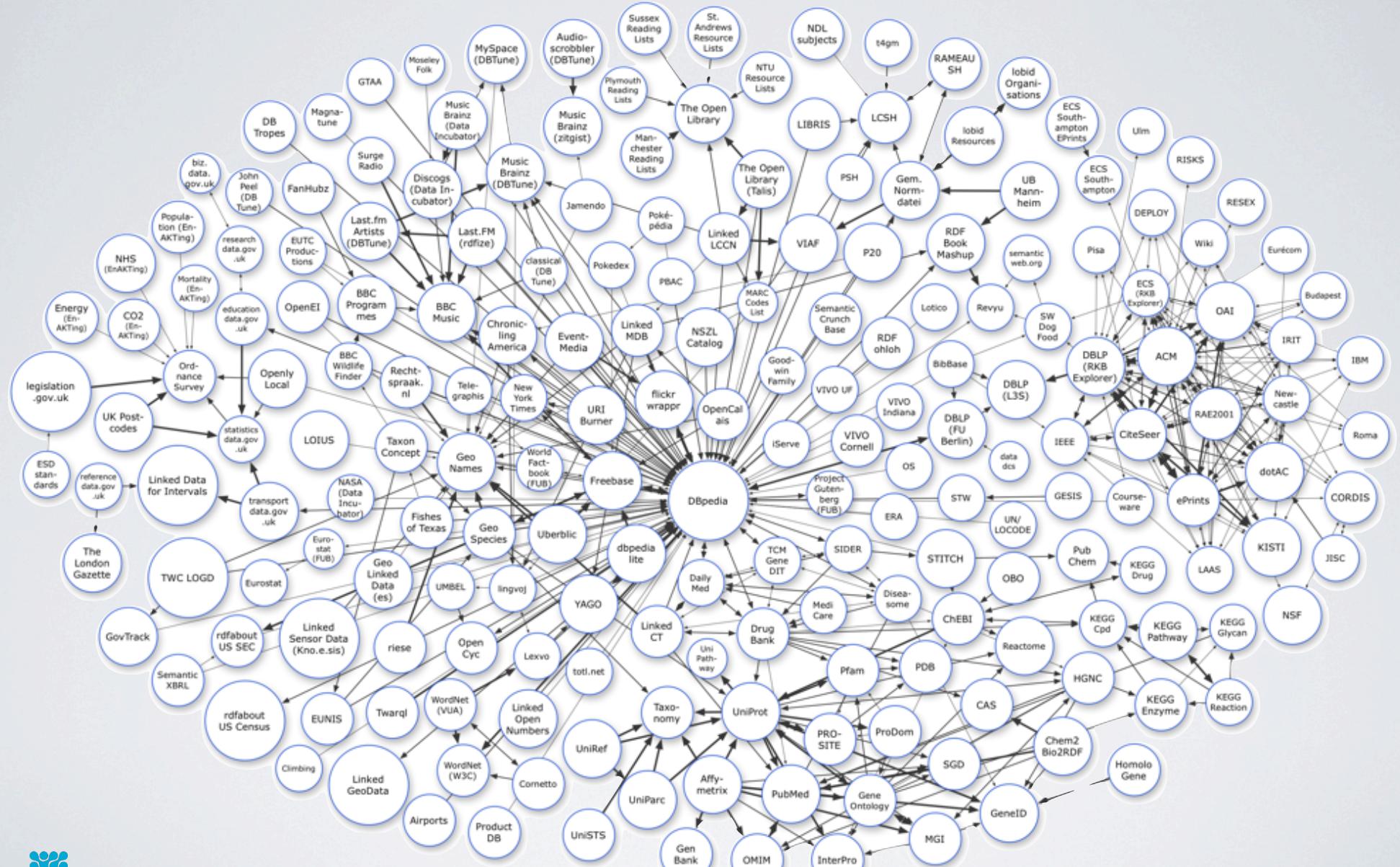
SIREN

TripleMap.com



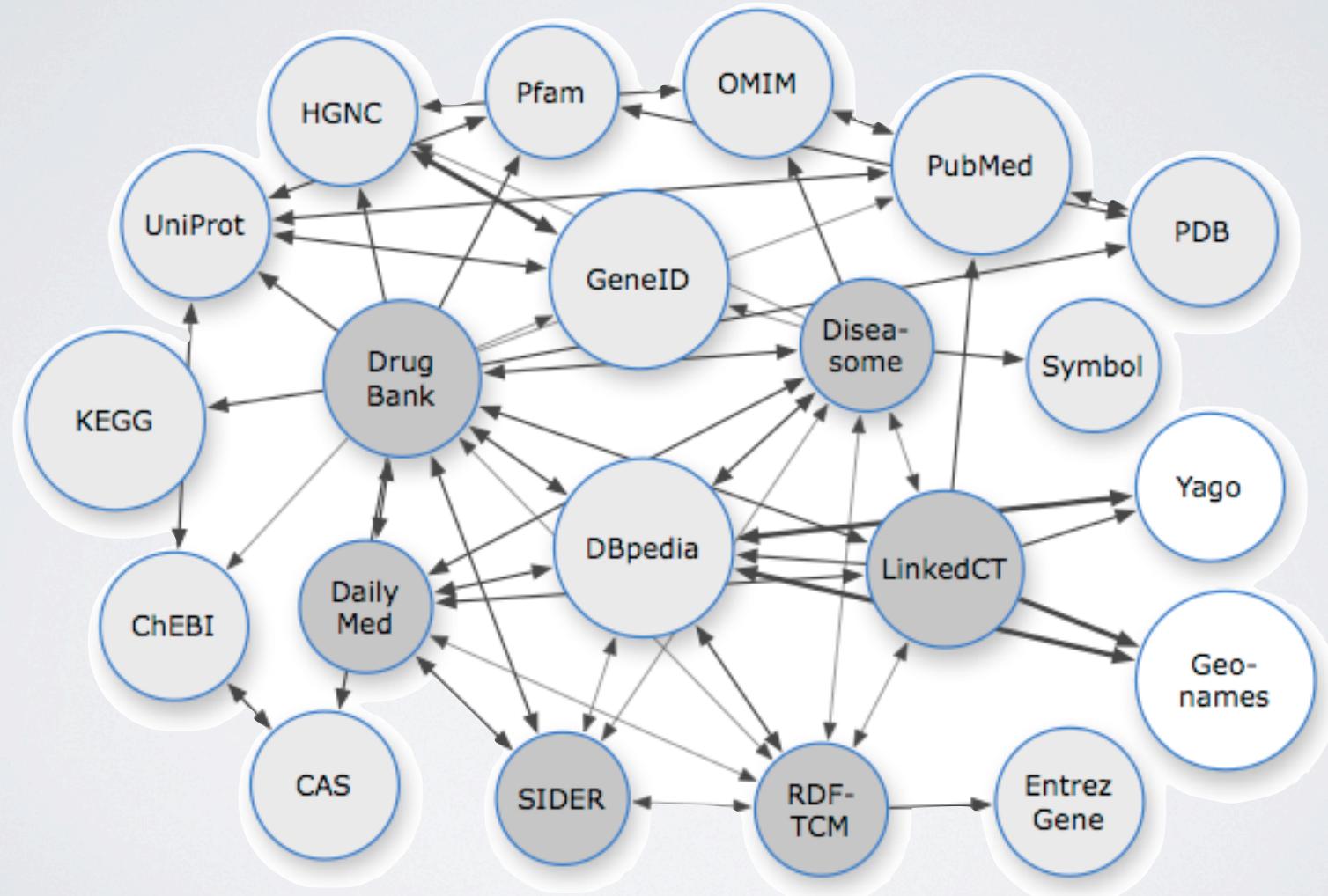


LINKING OPEN DATA



LINKING OPEN DRUG DATA

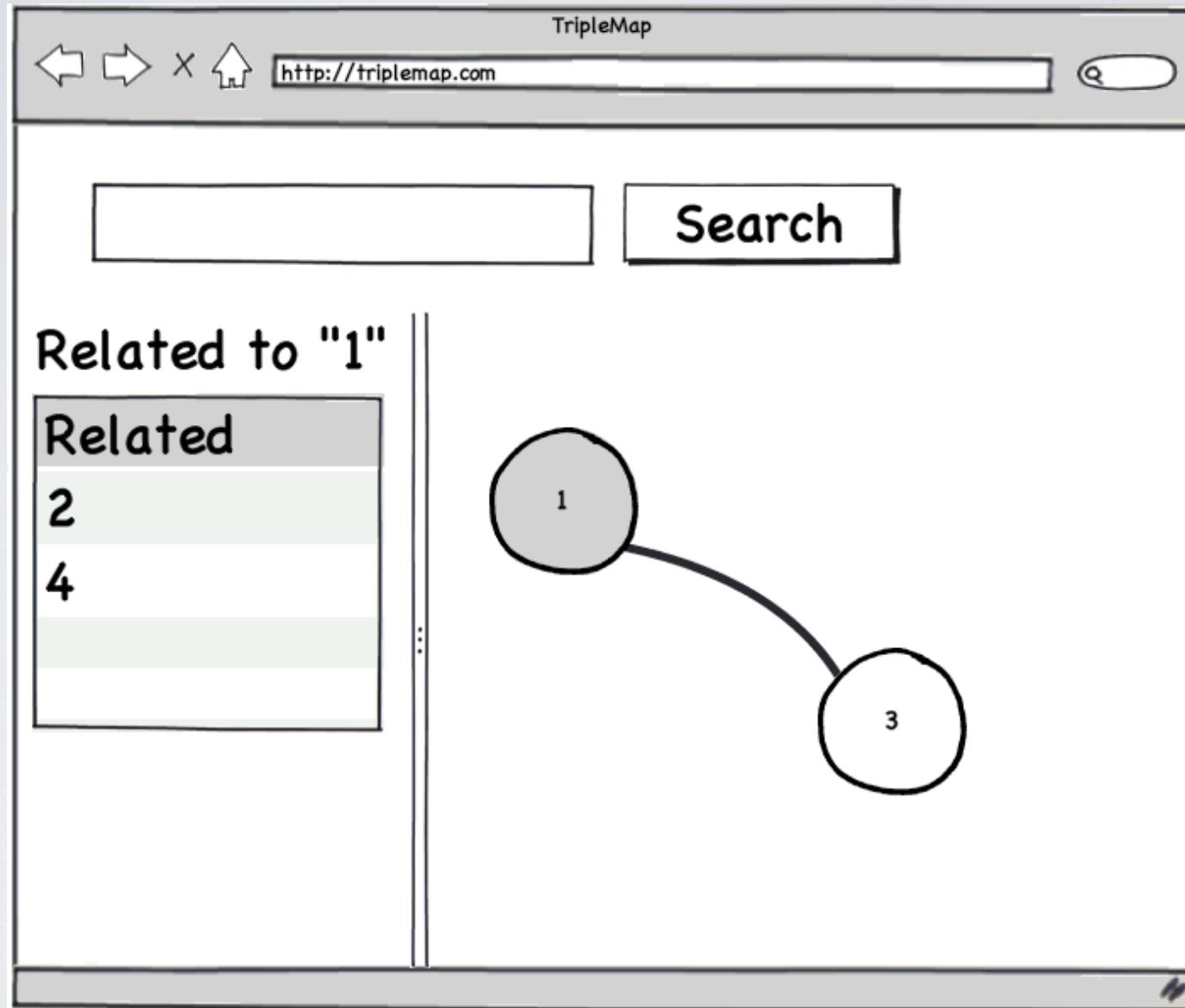
<http://esw.w3.org/HCLSIG/LODD>

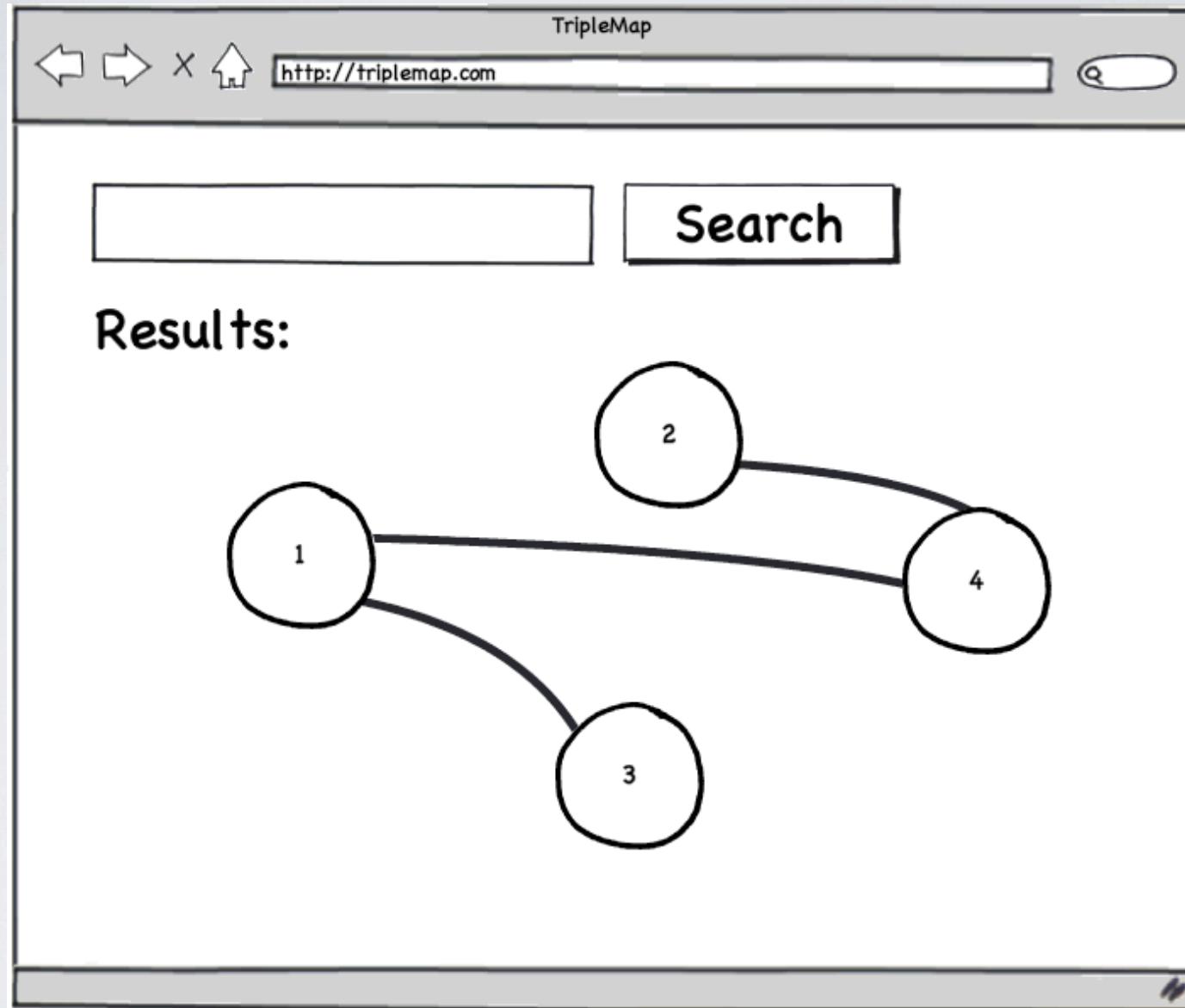


The diagram illustrates a web browser window with the following components:

- Title Bar:** TripleMap
- Address Bar:** http://triplemap.com
- Search Bar:** A large input field followed by a "Search" button.
- Results Section:** A heading "Results:" followed by a table.
- Table:** A data grid with columns labeled "ID" and "Label". The rows contain the following data:

ID	Label
1	Lorem
2	Ipsum
3	Dolor
4	Sit Amet





Ready, GO!

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

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

```
SELECT ?uri ?type ?label WHERE {  
    ?uri rdfs:label ?label .  
    ?uri rdf:type ?type .  
    FILTER (?label = '${params.query}')  
} LIMIT 10
```

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

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

query as literal value

case insensitive

DEMO

Baseline SPARQL Query Performance

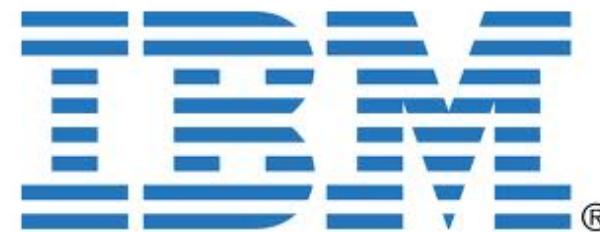
FASTER!

Lucene



Java API

Indexing and Searching Text

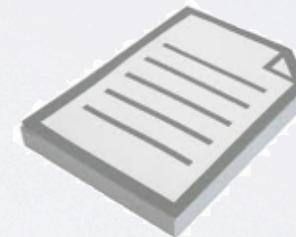


<http://wiki.apache.org/lucene-java/PoweredBy>

indexing

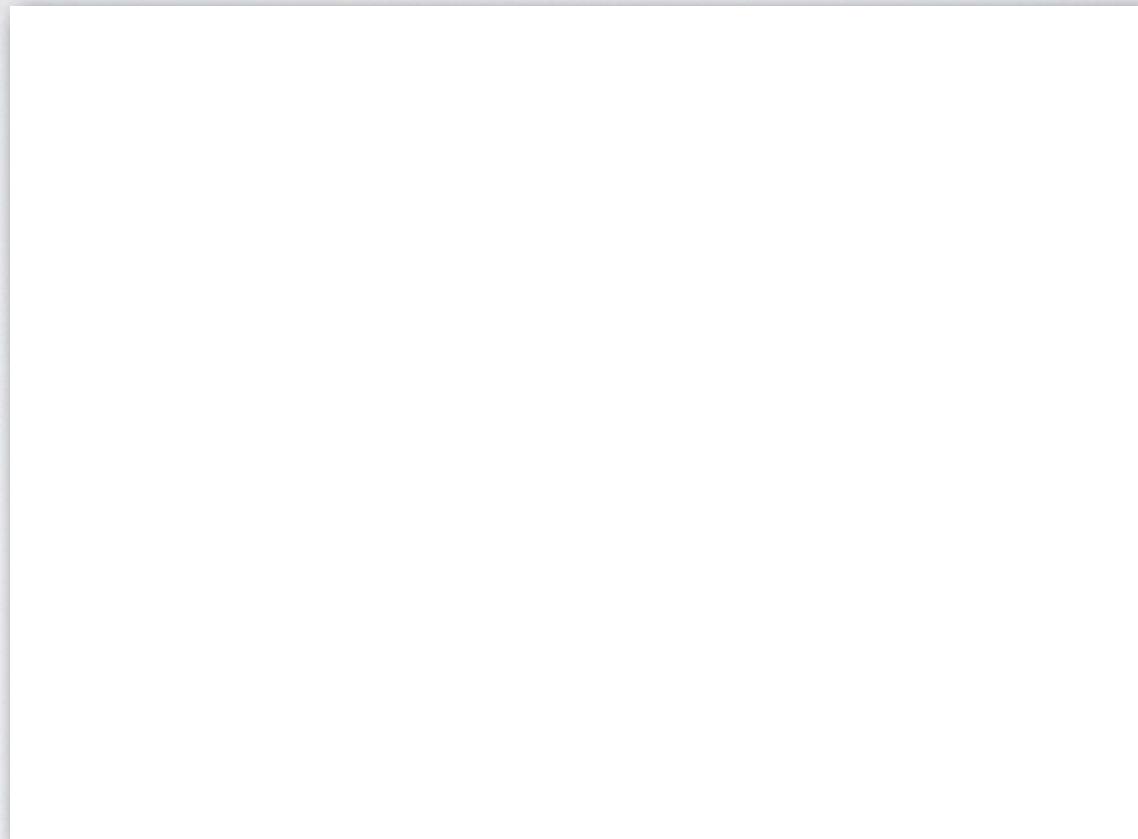
Lucene

storage





Document





Document

field	value
ID	2
name	“Mike Hugo”
company	“Entagen”
bio	“lorem ipsum dolor sum etc...”



Index

field	value
id	2
name	“Mike Hugo”
company	“Entagen”
bio	“lorem ipsum dolor sum etc..”

**Indexed
not
Stored**



Query:

name: mike



Query:

name: mike

Matching
Documents:

field	value
field	value
field	value
field	value
id	2



field	value
id	2





field	value
id	2



Lucene

field	value
id	2



field	value
ID	2
name	“Mike Hugo”
company	“Entagen”
bio	“lorem ipsum dolor sum etc...”

Simplest Solution

Lucene index of rdfs:label

Build the Index

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

```

**Build a SPARQL
query to find all the
rdfs:label properties**

```
sparqlQueryService.executeForEach(repo,
    def doc = new Document()
    String uri = it.uri.stringValue()
    String label = it.label.stringValue()
    doc.add(new Field(SUBJECT_URI_FIELD, Field.Store.YES, Field.Indexed)
            .setString(it.uri.stringValue()))
    doc.add(new Field(SUBJECT_LABEL_FIELD, Field.Store.YES, Field.Indexed)
            .setString(it.label.stringValue()))
    repo.add(doc)
}
```

```
sparqlQueryService.executeForEach  
(repository, queryLabels) {  
    String uri = it.uri.stringValue()  
    String label = it.label.stringValue()
```

**Execute the
SPARQL query**

```
new Document()
```

```
doc.add(new Field(SUBJECT_URI_FIELD,  
                  Field.Store.YES, Field.Indexed))
```

```
doc.add(new Field(LABEL_FIELD, label,  
                  Field.Store.NO, Field.Indexed))
```

```
writer.addDocument(doc)
```

```
String uri = it.uri.stringValue()
String label = it.label.stringValue()
```

```
Document doc = new Document()
doc.add(new Field(SUBJECT_FIELD,
                  uri,
                  Field.Store.YES,
                  Field.Index.ANALYZED),
       new Field(LABEL_FIELD,
                 label,
                 Field.Store.NO,
                 Field.Index.ANALYZED))
```

**Instantiate a
new Lucene
Document**

```
writer.addDocument(doc)
```



```
Document doc = new Document()  
doc.add(new Field(SUBJECT_URI_FIELD,  
value uri,  
field.Store.YES,  
Field.Index.ANALYZED))
```

```
doc.add(new Field(LABEL_FIELD,  
label,  
Field.Store.NO  
Field.Index.AN
```

**Add the Subject
URI to the
Document**

```
writer.addDocument(doc)
```

```
        Field.Store.YES,  
        Field.Index.ANALYZED));  
doc.add(new Field(LABEL_FIELD, key  
value label,  
        Field.Store.NO,  
        Field.Index.ANALYZED) )
```

```
writer.addDocument
```

**Add the Label field
to the document
(but don't store it)**

```
ly {  
ter.close() // Close index
```

```
        Field.Store.NO,  
        Field.Index.ANALYZED) )
```

```
writer.addDocument(doc)
```

```
}
```

```
finally {
```

```
writer.close() //
```

**Add the document
to the Index**

Query the Index

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

**Create a Lucene
Query from user
input** **for this value**

```
def sime  
List teQuery(query)  
def eime
```

```
render(view: 'index', model: [results:
```

```
IndexSearcher searcher = luceneSearcher
ScoreDoc[ ] scoreDocs =
    searcher.search(query, 10).scoreDocs
List results = []
def connection = r
scoreDocs.each {
    Document doc =
        String uri = doc[ SUBJECT_URI_FIELD ]
        Map labelAndType = sparqlQueryServ
    results << [uri: uri, type: labelA
}
connection.close()
return results
```

**Search the index
(limit 10) for
matching
documents**

```
list results = []
def connection = repository.connection
scoreDocs.each {
    Document doc = searcher.doc(it.doc)
    String uri = doc[SUBJECT_URI_FIELD]
    Map labelAndType =
        [span style="border: 1px solid black; padding: 5px;">
For each matching document, get the doc and extract the Subject URI

        type: labelAndType.type,
        label: labelAndType.label])
}
connection.close()
return results
```



```
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.add([  
        uri: uri,  
        type: labelAndType.type,  
        label: labelAndType.label  
    ]) }  
connection.close()  
return results
```

**Using the Subject
URI, load properties
from the triplestore])**

```
list results = []

def connection = repository.connection
scoreDocs.each {
    Document doc
    String uri =
    Map labelAndType
    sparqlQueryService.
        getLabelAndType(uri, connection)
    results.add([
        uri: uri,
        type: labelAndType.type,
        label: labelAndType.label])
}

connection.close()
return results
```

**return results
containing Subject
URI, Type, and Label**

DEMO

Lucene Index of Searchable Labels

WHAT ABOUT ENTITY RELATIONSHIPS?

WHAT ABOUT OTHER PROPERTIES?

SIREn

{ *Semantic Information Retrieval Engine* }

Lucene Extension

Indexing and Searching Semi-Structured Data

SIREn
{ Semantic Information Retrieval Engine }

Document

SIREn

{ Semantic Information Retrieval Engine }

Document

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Build the Index

```
Connection connection = repository.getConnection();
String subjectUris = """
    SELECT distinct ?uri
    WHERE {
        ?uri ?p ?o .
    }
"""
sparqlQueryService.executeForEach(repository, query, new SparqlQueryService.QueryProcessor() {
    @Override
    public void process(ResultSet result) throws RepositoryException {
        while (result.hasNext()) {
            Result row = result.next();
            String subjectUri = row.get("uri").toString();
            if (!subjectUris.contains(subjectUri)) {
                subjectUris += subjectUri + " ";
            }
        }
    }
});
```

**Select all Subject
URIs from the
triplestore**

```
String subjectUris = " ";
def doc = new DocumentBuilder()
doc.add(new Field(SUBJECT_URI_FIELD_NAME, subjectUris, Field.Store.YES, Field.Index.NON_TOKENIZED));
repository.add(doc);
System.out.println(subjectUris);
```

```
"""
```

```
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.NO))  
  
    StringWriter triplesStringWriter =  
    NTriplesWriter nTriplesWriter =  
        new NTriplesWriter(triplesStringWriter)  
    connection.exportStatements(  
        doc, nTriplesWriter)
```

**Execute the Sparql Query
For each URI, create a
new Document**

```
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
NTriplesWriter nT
new NTriplesW connection.exportStatements(
    new URIImpl(subjectUri),
    null, null, false,
```

**Add the Subject URI
to the Document**

```
Field.Index.ANALYZED) )
```

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

```
doc.add(new Field(T  
triplesStri  
Field.Store  
Field.Index.ANALYZED) )
```

**Get an NTriples
string from the
triplestore**

```
new URITriple(subjectURI),  
null, null, false,  
nTriplesWriter)
```

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

```
writer.addDocument(doc);
```

**Add the NTriples
string to the
document**

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

```
writer.addDocument(doc)
```

**Add the document
to the index**

Query the Index

```
SirenCellQuery predicate =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                RDFS.LABEL.stringValue()))));  
predicate.constraint = PREDICATE_CELL
```

```
SirenCellQuery object =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                params.query.toLowerCase())));  
object.constraint = OBJECT_CELL
```

```
SirenCellQuery predicate =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                RDFS.LABEL.stringValue())));  
predicate.constraint = PREDICATE_CELL
```

```
SirenCellQuery object =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                params.query.toLowerCase())))  
object.constraint = OBJECT_CELL
```

```
SirenCellQuery predicate =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                RDFS.LABEL.stringValue())));  
predicate.constraint = PREDICATE_CELL
```

of rdfs:label *

SirenCellQuery object

```
new SirenCellQuery(  
    new SirenTermQuery(  
        new Term(TRIPLES_FIELD,
```

*** note: could be any predicate!**

Case())

```
SirenCellQuery object =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                params.query.toLowerCase())))  
object.constraint = OBJECT_CELL
```

```
Query query = new SirenTupleQuery()  
query.add(predicate,  
        SirenTupleClause.  
        query.add(object,  
                SirenTupleClause.Occur.MUST)
```

**query the Triples
field**

List<results> = executeQuery(query)

```
SirenCellQuery object =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                params.query.toLowerCase())))  
object.constraint = OBJECT_CELL
```

```
Query query = new SirenTupleQuery()  
query.add(predicate,  
        SirenTupleClause.Occur.MUST)  
query.add(object,  
        SirenTupleClause.Occur.MUST)
```

for an object

```
SirenCellQuery object =  
    new SirenCellQuery(  
        new SirenTermQuery(  
            new Term(TRIPLES_FIELD,  
                params.query.toLowerCase())))  
object.constraint = OBJECT_CELL
```

```
Query query = new SirenTupleQuery()  
query.add(predicate,  
        SirenTupleClause.match(  
            object,  
            SirenTupleClause.Occur.MUST))
```

**matching the
user input**

```
params.query.toLowerCase( ))  
object.constraint = OBJECT_CELL
```

```
Query query = new SirenTupleQuery()  
query.add(predicate,  
          SirenTupleClause.Occur.MUST)  
query.add(object,  
          SirenTupleClause.Occur.MUST)
```

```
List results = ex
```

**Build query of with
both pieces -
predicate and object**

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query: **“imatinib”**

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:

triples field

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:

predicate = rdfs:label

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib". <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:

predicate = rdfs:label

object = “imatinib”

```
List executeQuery(Query query) {  
    IndexSearcher searcher = sirenSearcherM  
    ScoreDoc[ ] scoreDocs =  
        searcher.search(query, 10).scoreDocs  
    List results = []  
    def connection = repos...  
    scoreDocs.each {  
        Document doc = sear...  
        String uri = doc[S...  
        Map labelAndType = sparqlQueryServ...  
            getLabelAndType(uri, connection)  
        results.add([  
            uri: uri,  
            type: labelAndType.type,  
            label: labelAndType.label])  
    }  
}
```

**Search the index
(limit 10) for
matching
documents**

```
List results = []
def connection = repository.connection
scoreDocs.each {
    Document doc = searcher.doc(it.doc)
    String uri = doc[ SUBJECT_URI_FIELD ]
    Map labelAndType = sparqlQueryService
        .getLabelAndTypeForSubjectUri(uri)
    results.add(
        type: labelAndType.type,
        label: labelAndType.label )
}
connection.close()
return results
```

For each matching document, get the doc and extract the Subject URI

```
connection = repository.connection
reDocs.each {
    Document doc = searcher.doc(it.doc)
    String uri = doc[SUBJECT_URI_FIELD]
    Map labelAndType = sparqlQueryService.
        getLabelAndType(uri, connection)
    results.add([
        uri: uri,
        type: labelAndType[type],
        label: labelAndType[label]])
}
connection.close()
return results
```

**Using the Subject
URI, load properties
from the triplestore**

```
String uri = doc[SUBJECT_ORI_FIELD]  
Map labelAndType = sparqlQueryService.  
    getLabelAndType(uri, connection)  
results.add([  
    uri: uri,  
    type: labelAndType.type,  
    label: labelAndType.label])
```

```
nection.close()  
urn results
```

**return results
containing Subject
URI, Type, and Label**

DEMO

SIREn Index of RDF Entities

FLEXIBILITY

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib". <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:

predicate = rdfs:label

object = “imatinib”

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib". <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:

object = “imatinib”

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib". <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" <DB00619> drugbank:target <targets/1588> .

Query:

object = “imatinib”

OR

object = “gleevec”

MORE THAN LITERALS

field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:



predicate = brandName

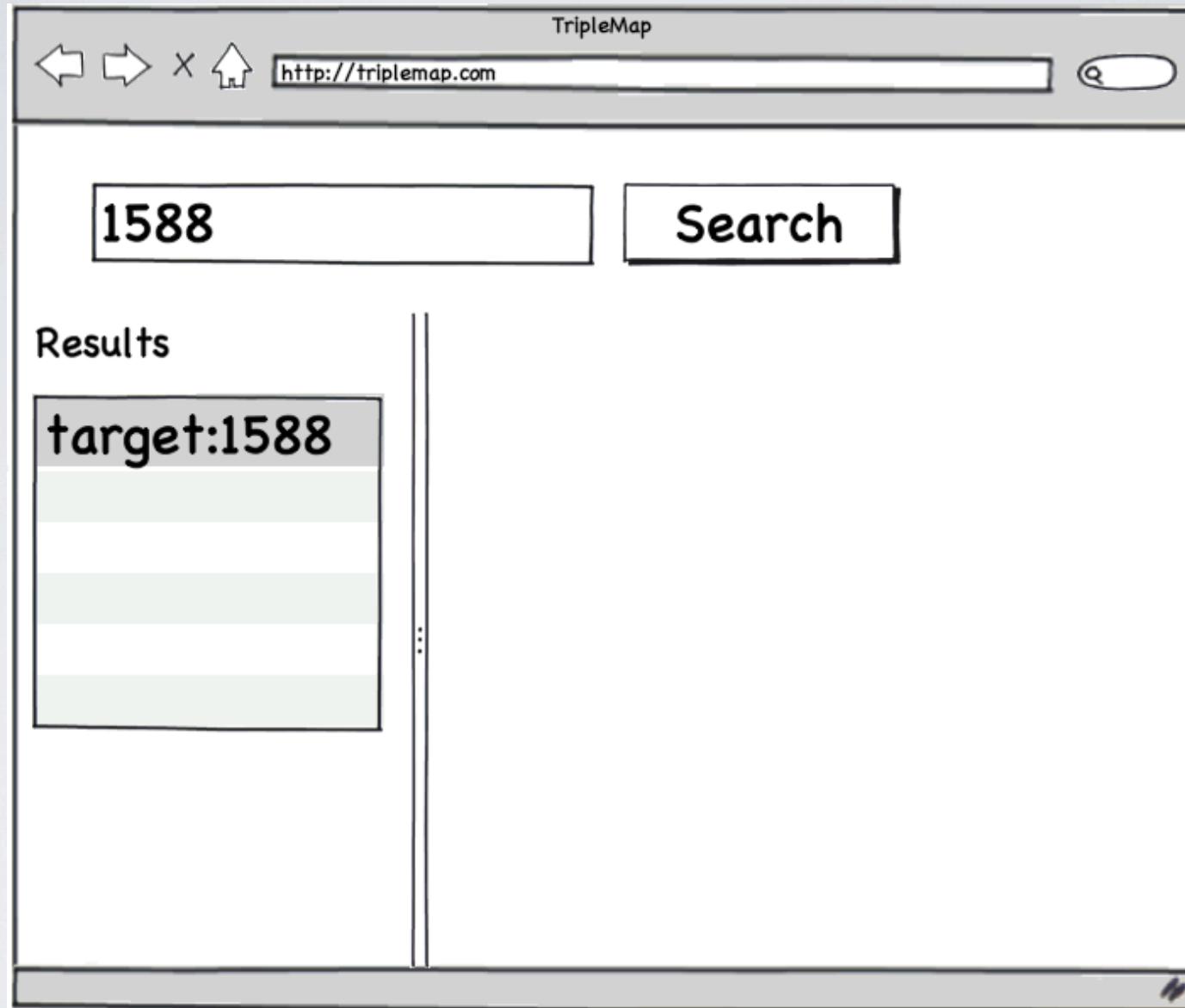
field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

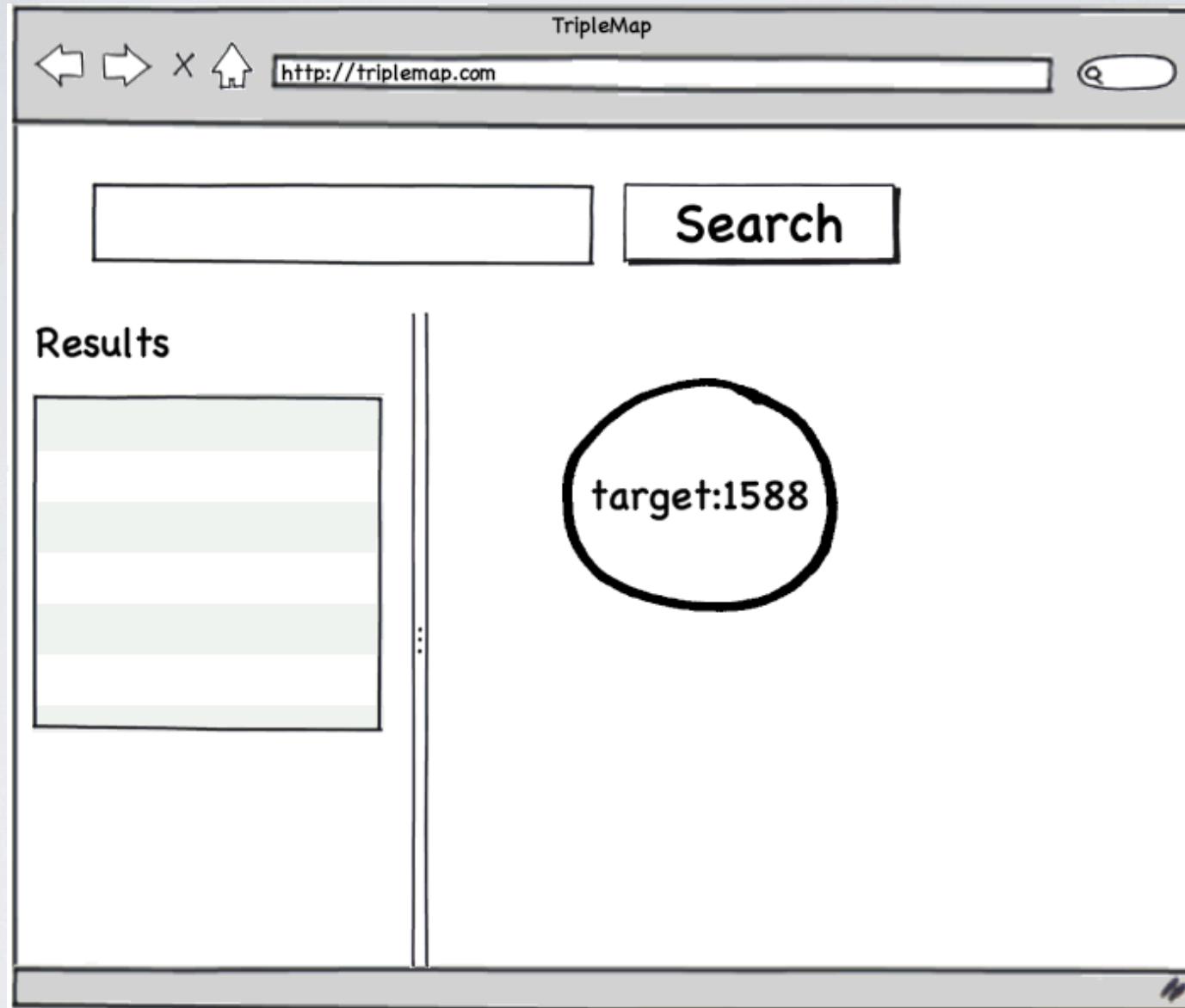
Query:

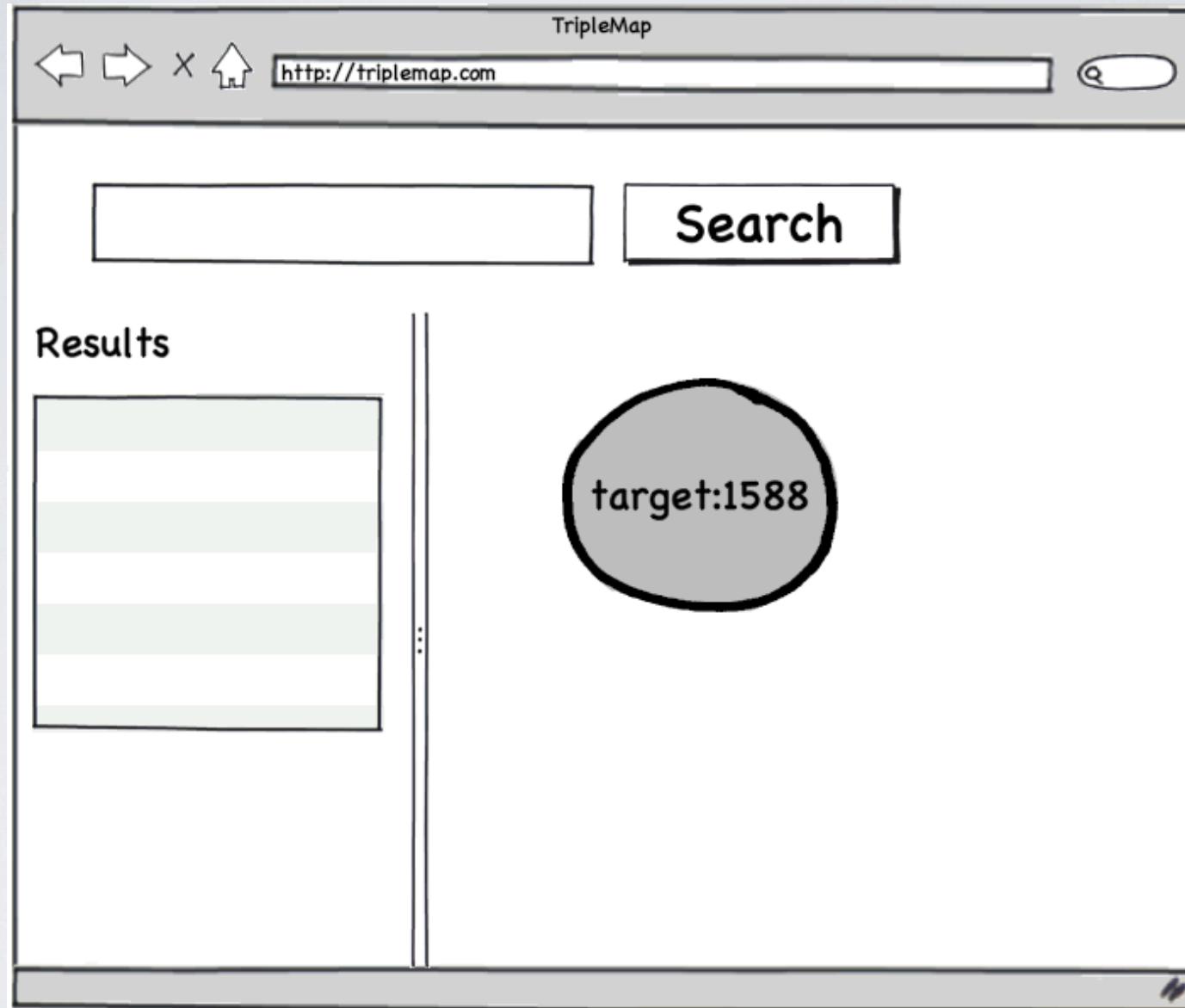


predicate = target

RELATIONSHIPS

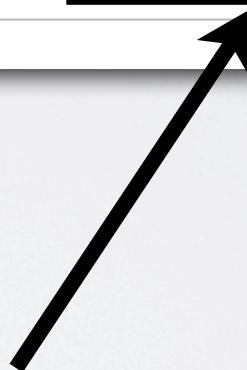




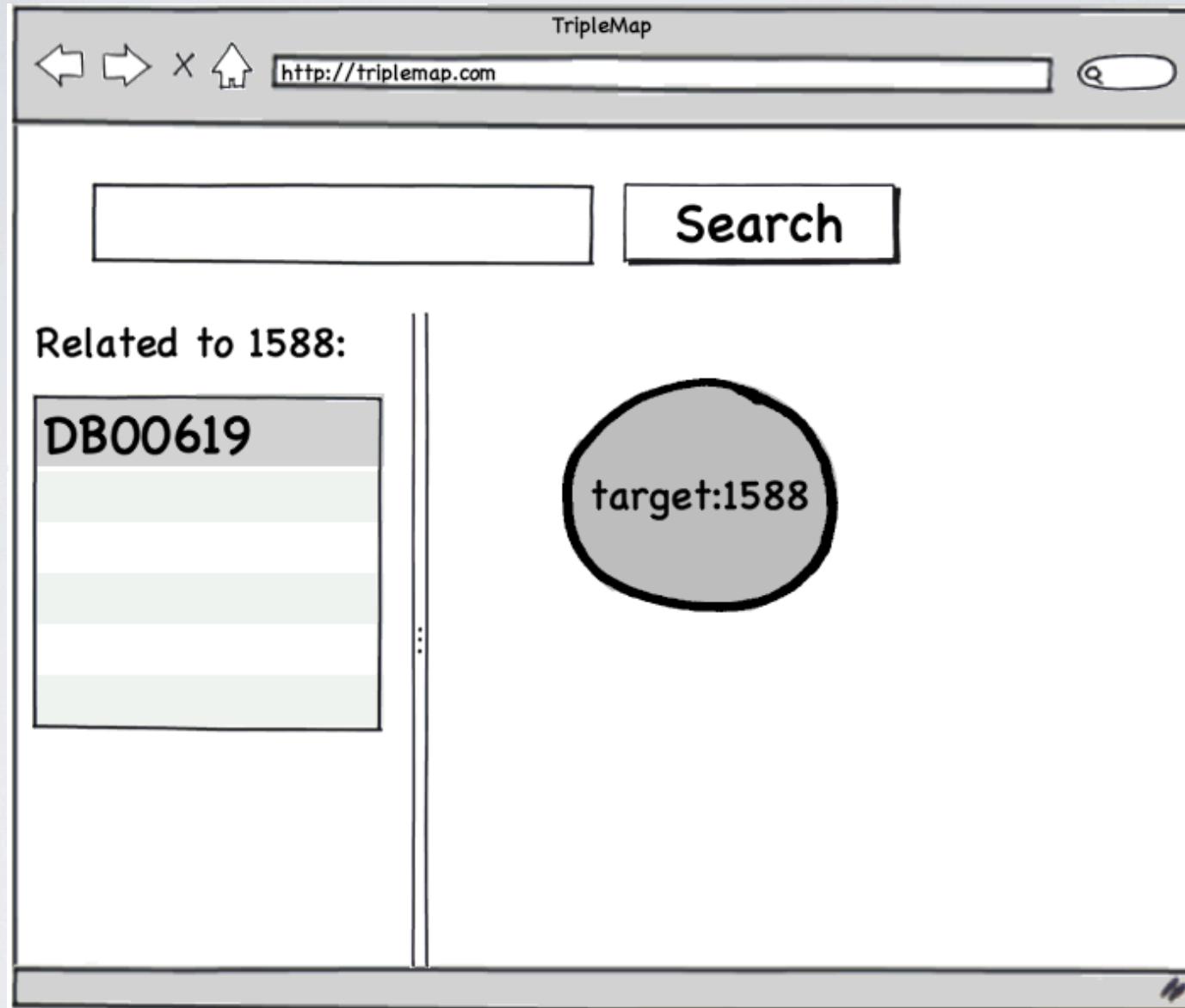


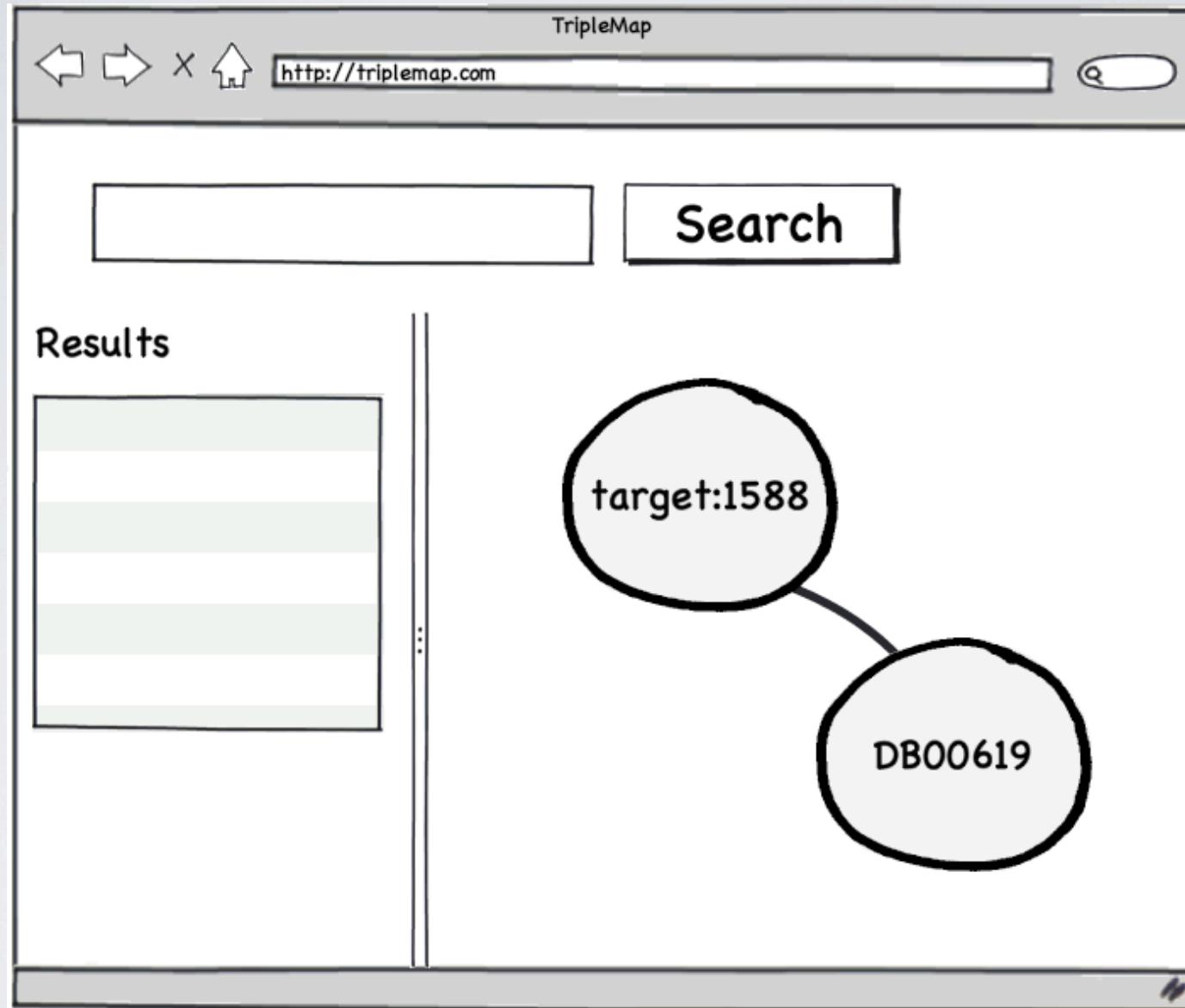
field	value
URI	<DB00619>
triples	<DB00619> rdfs:label "Imatinib" . <DB00619> rdf:type <drugbank:drugs> . <DB00619> drugbank:brandName "Gleevec" . <DB00619> drugbank:target <targets/1588> .

Query:



object = <targets/1588>





DEMO

Searching SIREn Index for Relationships

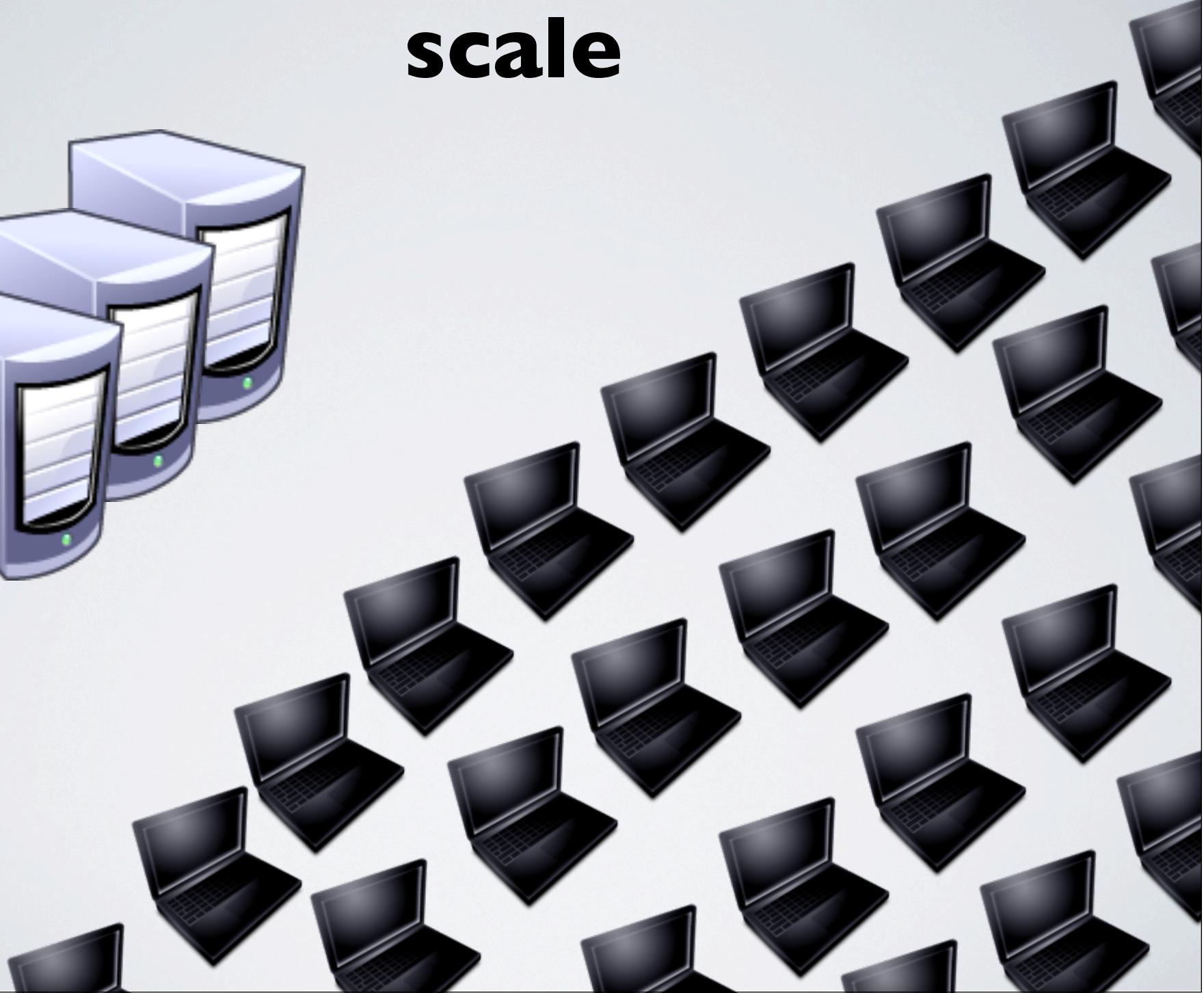
SIREn

{ *Semantic Information Retrieval Engine* }



Distributed Indexing and Searching Semi-Structured Data

scale





Search interface type: [Simple](#) [Advanced](#) [Guru](#) [Query Language](#) [Documentation](#)

keyword(s)

ntriple(s)

```
* <name> "Tim Berners Lee"
```

filter(s)

SEARCH

Group By Dataset:

Sorted by:

relevance



Quick filters ([All options](#))

Time range:

[Any date](#)

[Today](#) [Yesterday](#) [Last week](#)

[Last month](#) [Last year](#)

Format:

[Any format](#)

[RDF](#) [RDFA](#) [MICROFORMAT](#) [XFN](#)

[HCARD](#) [HCALENDAR](#) [HLISTING](#)

[HRESUME](#) [LICENSE](#) [GEO](#) [ADR](#)

Sindice search:
* <name> "Tim Berners Lee" found 699 documents

<http://dig.csail.mit.edu/2007/01/camp/data> (RDF)

+ 2009-08-12 - 387 triples in 43.2 kB

<http://dig.csail.mit.edu/2007/01/camp/data> ([Search](#)) [Inspect](#): ([Cache](#)) ([Live](#))

http://dbpedia.org/data/Tim_Berners-Lee (RDF)

- 2011-05-18 - 171 triples in 38.8 kB

SIGMA

FULL CONTENT

GRAPH

TRIPLES(890)

See all we know about "Tim Berners Lee"

Query Parser

[Home](#)[About](#)[Se](#)

Search interface type: [Simple](#) [Advanced](#) [Guru](#) [Query Language Documentation](#)

keyword(s)

ntriple(s)

```
* <name> "Tim Berners Lee"
```

filter(s)

SEARCH

Group By Dataset:

Sorted by:

relevance



Quick filters ([All options](#))

Time range:

[Any date](#)

[Today](#) [Yesterday](#) [Last week](#)

[Last month](#) [Last year](#)

Format:

Sindice search: * <name> "Tim Berners Lee" found 699 documents

<http://dig.csail.mit.edu/2007/01/camp/data> (RDF)

[+] 2009-08-12 - 387 triples in 43.2 kB

<http://dig.csail.mit.edu/2007/01/camp/data> ([Search](#)) Inspect: ([Cache](#)) ([Live](#))

http://dbpedia.org/data/Tim_Berners-Lee (RDF)

[-] 2011-05-18 - 171 triples in 38.8 kB

DEMO

SIREn in action on TripleMap.com

Search Maps Analyze

 Filter Sort

Select - Total count: 2

< Back Indication

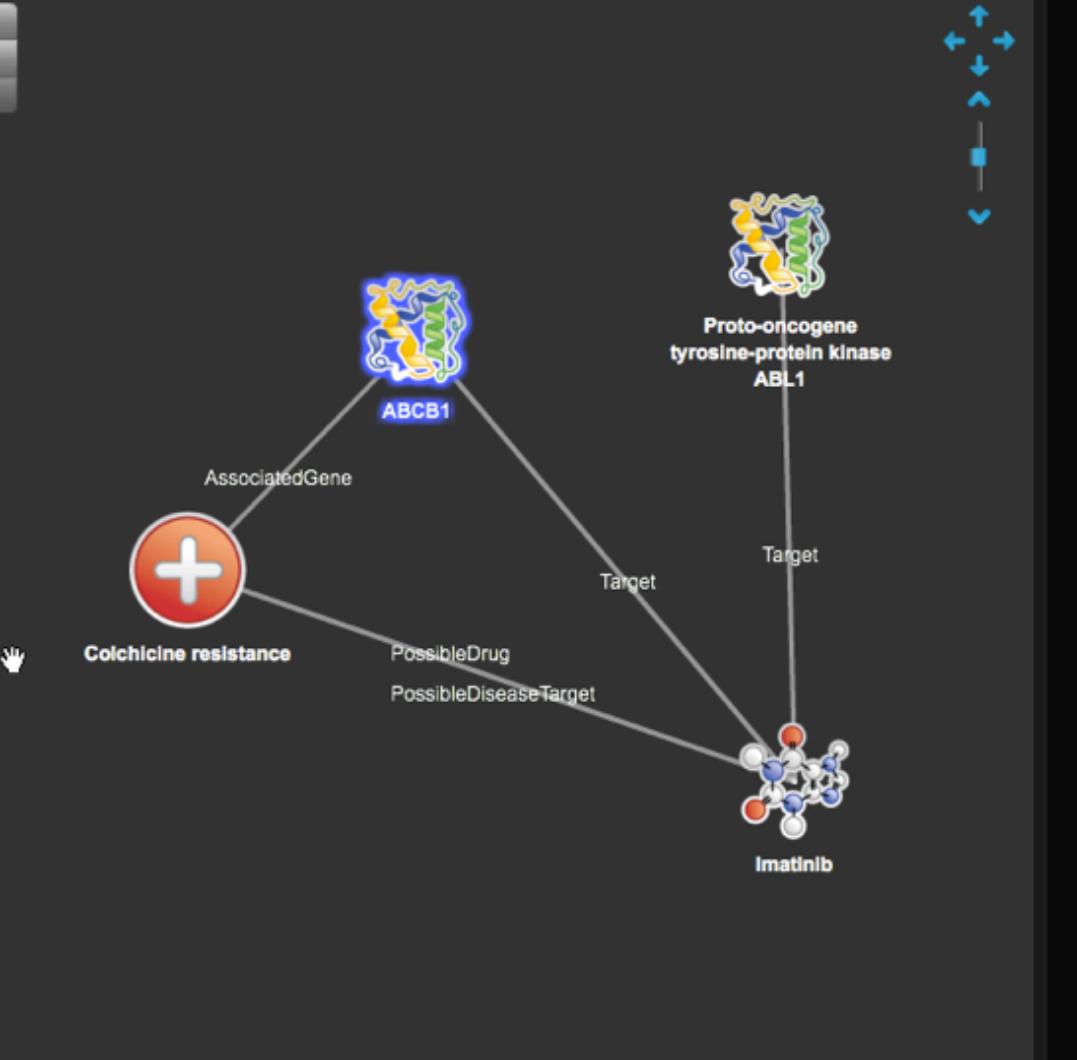
- Colchicine resistance
- Colchicine_resistance

144

2

16

Map View



SPARQL

LUCENE

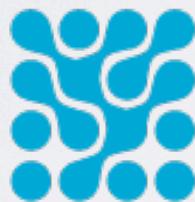
SIREN

TripleMap.com



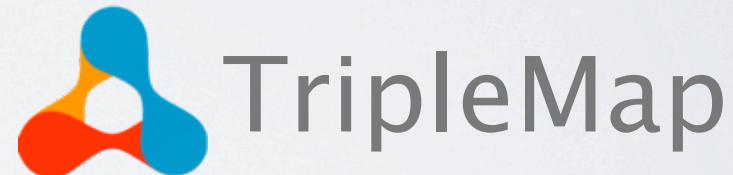
QUESTIONS?

mike@entagen.com / twitter: @piragua



ENTAGEN
ACCELERATING INSIGHT

<http://www.entagen.com>



TripleMap

<http://www.triplemap.com>