

JUDCon

JBoss Users & Developers Conference

London:2011

Who am I?

Sanne Grinovero

Software Engineer at Red Hat

- Hibernate, especially Search
- Infinispan, focus on Query and Lucene
- Hibernate OGM
- Apache Lucene
- JGroups



Our Index

- Searching in Infinispan
 - Map/Reduce
 - Fulltext indexing
- Infinispan Query engine
- Clustering a Lucene index
- Cloud deployed applications
- Future

Infinispan

- An advanced multi-node cache
- A transactional scalable datagrid targeting high performance and cloud
- A “NoSQL database”, a *key-value* store
 - How do you query a key value store?

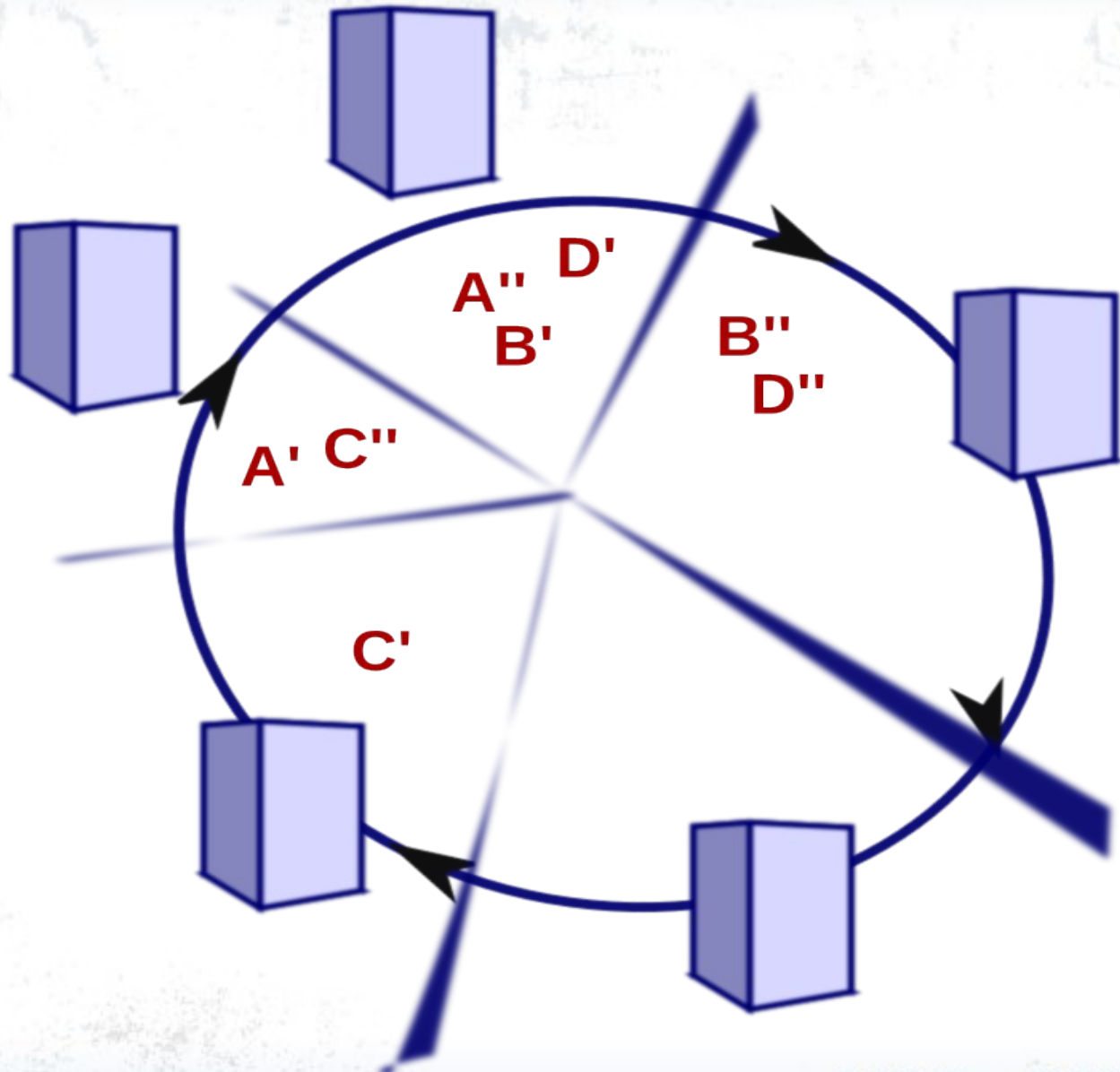
```
SELECT * FROM GRID
```


To Query a Grid

- What's in C7 ?

```
Object v =  
    cache.get("c7");
```





If you don't know the key, no
way to find the value

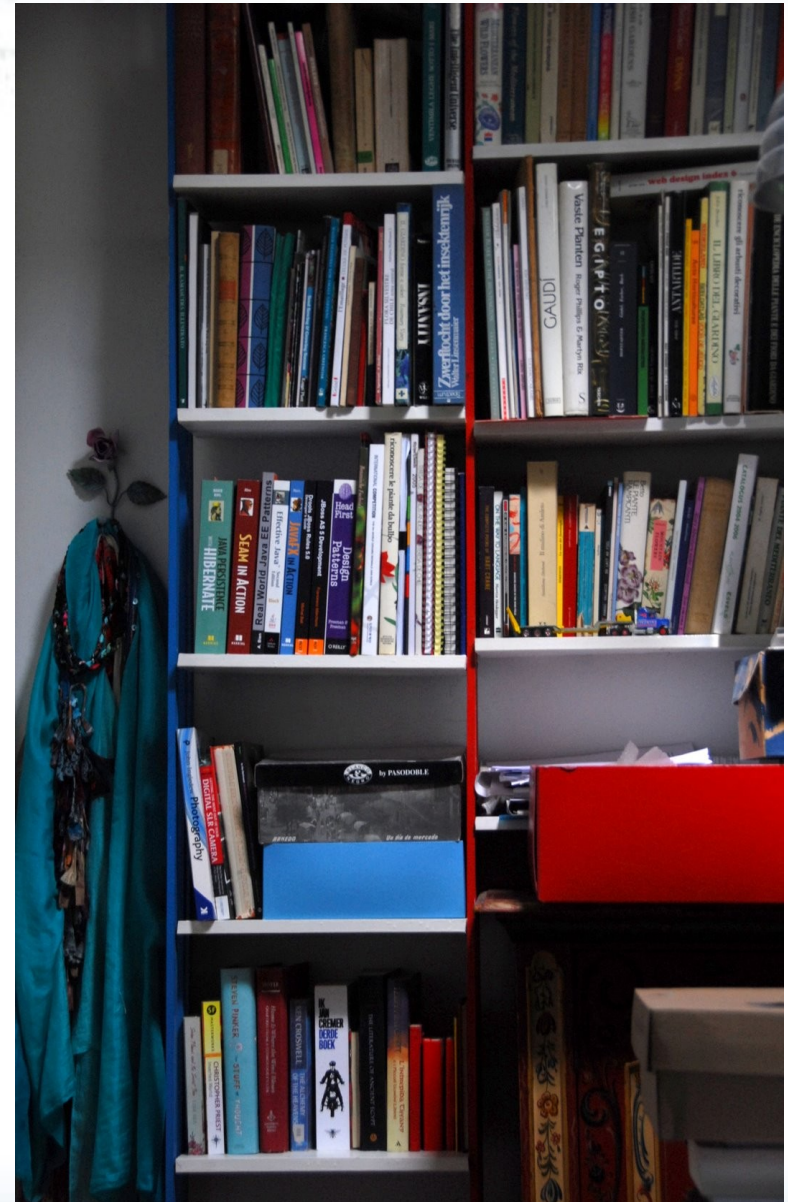
Some services have no chance:

Enter ISBN code:

Find Book

Let's test my bookshelf

- Where's Hibernate Search in Action?
- Could you hand me ISBN 978-1-933988-17-7 ?
- How many books about Gaudí ?



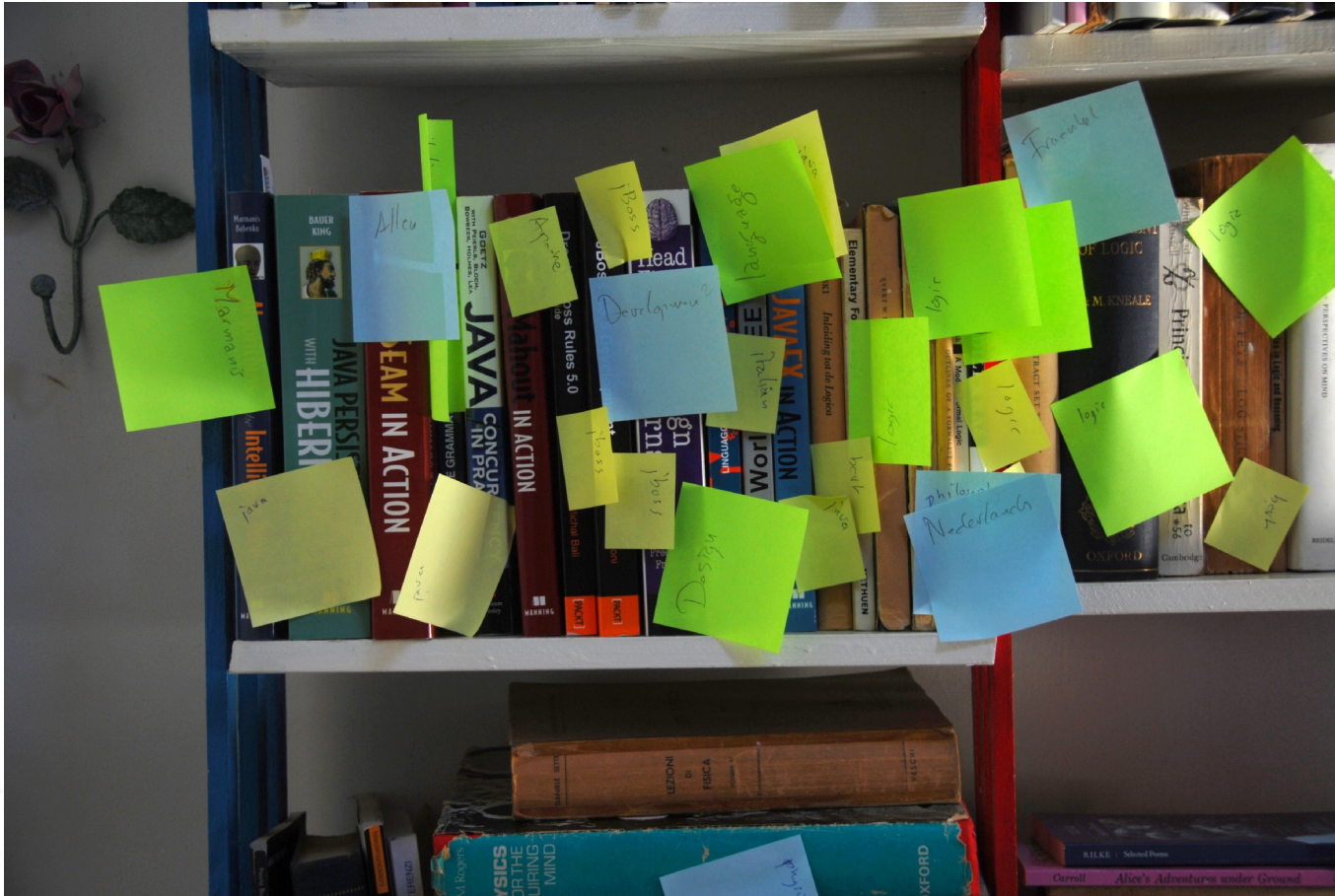
A *real-world* example



A *real-world* example



A *real-world* example



Bookshelves don't scale

How to implement the bookshelf features on a k/v?

- Where's "Hibernate Search in Action"?
- Can you hand me "ISBN 978-1-933988-17-7" ?
- How many books about Gaudí ?

Most document based NoSQLs support Map/Reduce

- Infinispan does not focus on documents
 - That won't stop you from using any format JSON, XML, YAML, Java:

```
public class Book implements Serializable {  
  
    final String title;  
    final String author;  
    final String editor;  
  
    public Book(String title, String author, String editor) {  
        this.title = title;  
        this.author = author;  
        this.editor = editor;  
    }  
  
}
```

Iterate & collect

```
class TitleBookSearcher implements
    Mapper<String, Book, String, Book> {
    final String title;
    public TitleBookSearcher(String t) { title = t; }
    public void map(String key, Book value, Collector collector){
        if ( title.equals( value.title ) )
            collector.emit( key, value );
    }
}
```

```
class BookReducer implements
    Reducer<String, Book> {
    public Book reduce(String reducedKey, Iterator<Book> iter) {
        return iter.next();
    }
}
```

How to implement the bookshelf features on a k/v?

- ✓ Where's "Hibernate Search in Action"?
- ✓ Can you hand me "ISBN 978-1-933988-17-7" ?
- ✗ How many books *about* "Shakespeare" ?
 - To properly score fulltext results we need to consider relative term frequencies on the whole corpus
 - Pre-tagging is a poor choice

Apache Lucene

- Open source Apache™ top level project
- Countless products and sites use it
- Integrates in Hibernate via Hibernate Search
- Clusterable via Infinispan



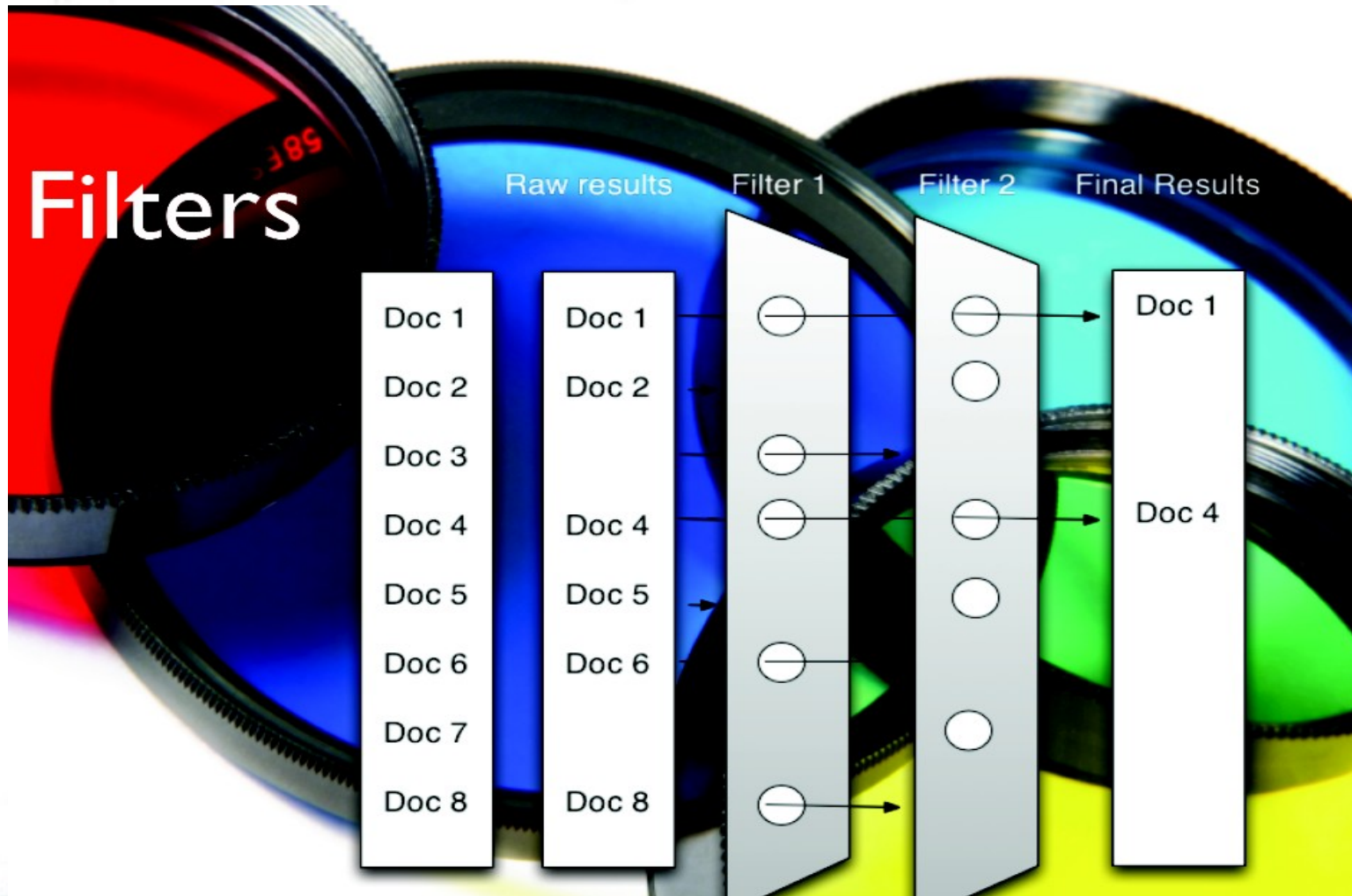
What does Lucene get us?

- Similarity scoring searches
- Advanced text analysis
 - Synonyms, Stopwords, Stemming, ...
- Reusable declarative Filters
- TermVectors
- MoreLikeThis
- Faceted Search
- Speed!

Lucene: Stopwords

a, able, about, across, after, all, almost, also, am, among, an, and, any, are, as, at, be, because, been, but, by, can, cannot, could, dear, did, do, does, either, else, ever, every, for, from, get, got, had, has, have, he, her, hers, him, his, how, however, i, if, in, into, is, it, its, just, least, let, like, likely, may, me, might, most, must, my, neither, no, nor, not, of, off, often, on, only, or, other, our, own, rather, said, say, says, she, should, since, so, some, than, that, the, their, them, then, there, these, they, this, tis, to, too, twas, us, wants, was, we, were, what, when, where, which, while, who, whom, why, will, with, would, yet, you, your

Filters



Faceted Search

[Shop All Departments](#) Search

Books [Advanced Search](#) [Browse Subjects](#) [New Releases](#) [Bestsellers](#) [TI](#)

Department

- < Any Department
 - < Books
 - Computers & Internet**
 - Programming (14)
 - Computer Science (4)
 - Databases (2)
 - Software (2)
 - Web Development (2)
 - Networking (1)
 - Home Computing (1)

Format

- Paperback (15)

Author

- Any Author**
 - Joe Vitale (1)

Shipping Option [\(What's this?\)](#)

- Any Shipping Option**
 - Free Super Saver Shipping

Books > Computers & Internet > "Hibernate Search"

Showing 1 - 12 of 15 Results

- LOOK INSIDE!**



Hibernate Search in Action I

★★★★★ (3 customer reviews)

Formats

Paperback

Order in the next **2 hours** to get it by **Monday, Apr 18.** ~~\$49.95~~

Only 1 left in stock - order soon.

Eligible for **FREE** Super Saver Shipping.

Excerpt - Page 1: "... breaking the susy
Surprise me! See a random page in the
- LOOK INSIDE!**



Spring Persistence with Hib
(Nov 2, 2010)

★★★★☆ (5 customer reviews)

Formats

Paperback

JUDCon2011:London
JBoss Users & Developers Conference

Would you want a web
Search engine to return hits
in alphabetical order?

[A - Wikipedia](#) ☆

A o a (nome italiano "a" /a/) è la 1ª lettera dell'alfabeto latino e italiano e anche nella maggior parte degli alfabeti derivanti da quello fenicio. ...

[Origine](#) - [Usò nelle lingue](#) - [Codici informatici](#) - [Voci correlate](#)

it.wikipedia.org/wiki/A - [Copia cache](#) - [Simili](#)

[Home ATM, Azienda Trasporti Milanese ATM, Azienda Trasporti Milanese](#) ☆

Immagine campagna e link a nuova pagina sito BikeMi ... Mercoledì di Champions a San Siro. Il servizio Atm per Milan-Real Madrid - 02/11 ...

www.atm-mi.it/ - [Copia cache](#) - [Simili](#)

[A come AMORE - Frasi D'Amore, Poesie D'Amore, Sms D'Amore](#) ☆

A Come Amore, Tutto Sull'Amore: Frasi D'Amore, Poesie D'Amore, SMS Sull'Amore, Forum Amore, Consigli D'Amore, Gelosia, Affirma Di Coppo.

www.icconeanche.com/ - [Copia cache](#) - [Simili](#)

[A sbafo](#) ☆

Prova a contattare l'appresentante del casino online tramite le chat room, telefono o la posta elettronica. Cerca di porre le giuste domande e ...

www.sbafo.com/ - [Copia cache](#) - [Simili](#)

[adieta.it | dieta, alimentazione, fitness, salute e bellezza](#) ☆

Dimagrire, mettersi a dieta, mangiare sano, dieta, in linea, calcolare il peso ideale, ricette light: tutti i modi e ancor di più quello che preferite ...

www.adieta.it/ - [Copia cache](#) - [Simili](#)

[Home Page - Scherzi a parte](#) ☆

Direttore Format, Programmi e Sit-com R.T.I., Fatma Ruffini è nata a Reggio Emilia e si è laureata in Storia dell'Arte. E' approdata a Canale 5, nel 1981, ...

www.scherziaparte.mediaset.it/ - [Copia cache](#) - [Simili](#)

[Dove Andiamo A Cena Stasera? - Guida Ufficiale Ai Ristoranti](#) ☆

ristorante guida, ristoranti guida, guida ristoranti, guida ufficiale ai ristoranti in internet, official restaurants guide of internet, first born,

www.acena.it/ - [Copia cache](#) - [Simili](#)

The downsides

- Requires an Index
 - in memory
 - on filesystem
 - in Infinispan
- Made of immutable segments
 - Optimized for search speed, not for updates
- A world of strings and frequencies

Infinispan Query quickstart

- Enable it in configuration
- Have infinispan-query.jar in your classpath
- Annotate your POJO values to specify what to index

```
<dependency>  
  <groupId>org.infinispan</groupId>  
  <artifactId>infinispan-query</artifactId>  
  <version>5.1.0.BETA3</version>  
</dependency>
```

Enable Infinispan Query, programmatically

```
Configuration c = new Configuration()  
    .fluent()  
    .indexing()  
    .addProperty(  
"hibernate.search.default.directory_provider",  
"ram")  
    .build();  
  
CacheManager manager = new DefaultCacheManager(c);
```

Enable Query in Infinispan XML configurations

```
<?xml version="1.0" encoding="UTF-8"?>
<infinispan
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="urn:infinispan:config:5.0
http://www.infinispan.org/schemas/infinispan-config-5.0.xsd"
  xmlns="urn:infinispan:config:5.0">
<default>
  <indexing enabled="true" indexLocalOnly="true">
    <properties>
      <property name="hibernate.search.option1" value="..." />
      <property name="hibernate.search.option2" value="..." />
    </properties>
  </indexing>
</default>
```


Annotate your model

```
@ProvidedId @Indexed
public class Book implements Serializable {

    @Field String title;
    @Field String author;
    @Field String editor;

    public Book(String title, String author, String editor) {
        this.title = title;
        this.author = author;
        this.editor = editor;
    }
}
```

Run a Query

```
SearchManager qf = Search.getSearchManager(cache);

Query query = qf.buildQueryBuilderForClass(Book.class)
    .get()
    .phrase()
        .onField("title")
        .sentence("in action")
    .createQuery();

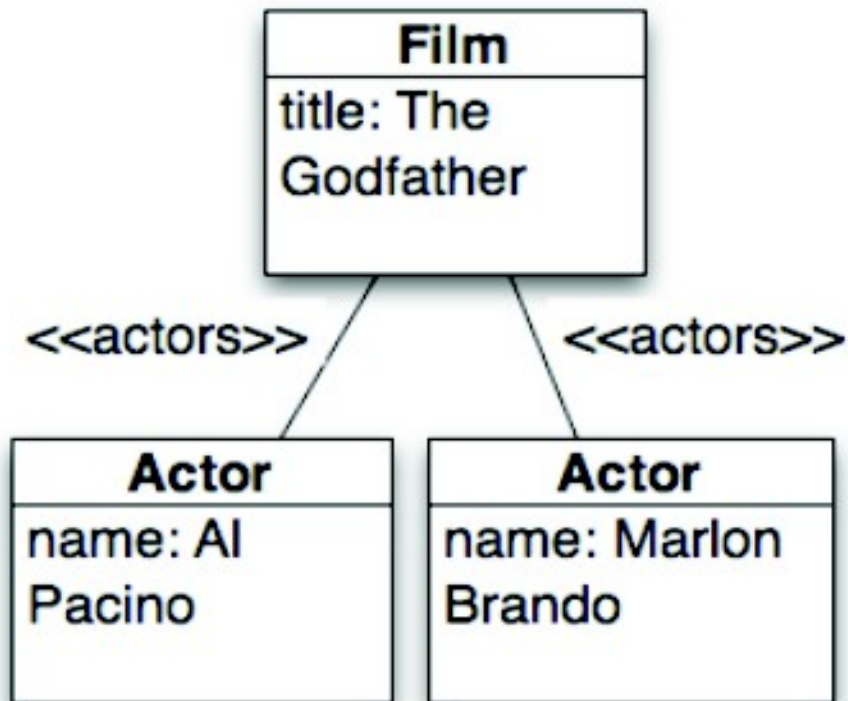
List<Object> list = qf.getQuery(query).list();
```

The code

- Integrates Hibernate Search
 - Listen to Hibernate events & transactions
 - Infinispan events & transactions
 - Maps Java types and model graphs to Lucene Documents
 - Thin-layer design

Index mapping

Object world



Index world

Film Document	
title	The Godfather
actor	Al Pacino
actor	Marlon Brando

declarative analyzers

```
@Entity @Indexed
@AnalyzerDef(name = "frenchAnalyzer", tokenizer =
    @TokenizerDef(factory=StandardTokenizerFactory.class), filters = {
        @TokenFilterDef(factory = LowerCaseFilterFactory.class),
        @TokenFilterDef(factory = SnowballPorterFilterFactory.class,
            params = {@Parameter(name = "language", value = "French")})
    })
public class Book {

    @Field(index=Index.TOKENIZED, store=Store.NO)
    @Analyzer(definition = "frenchAnalyzer")
```


Query test

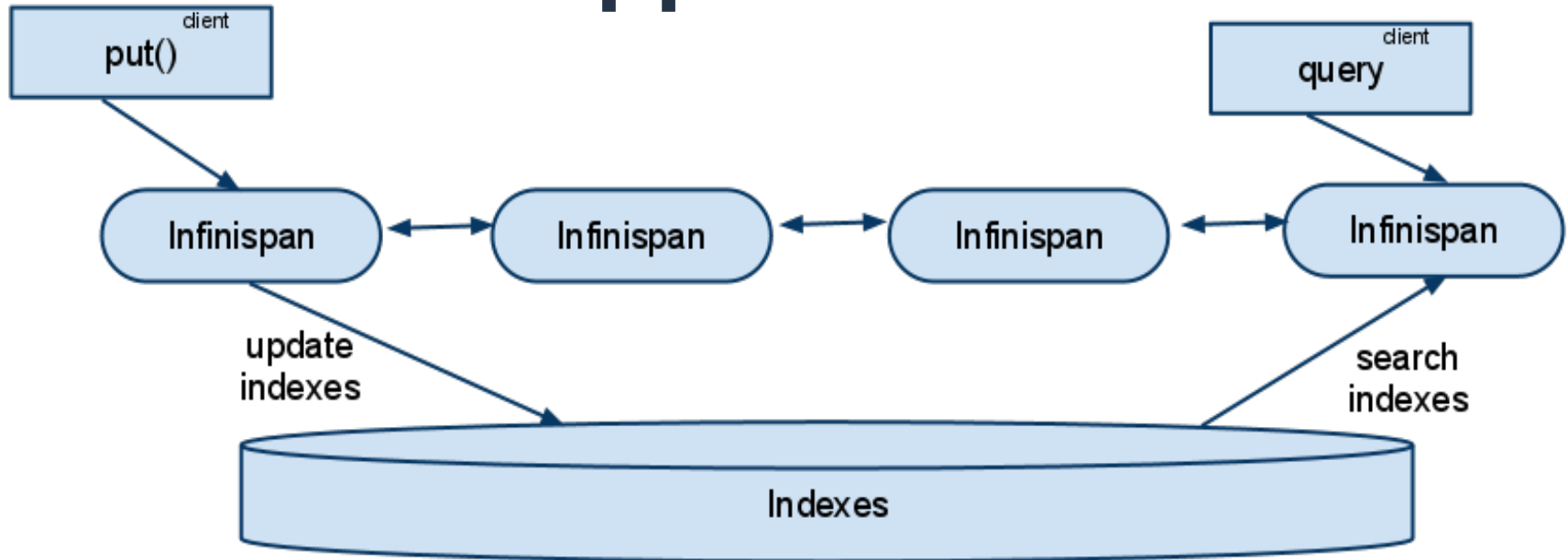
<https://github.com/infinispan/infinispan>

I'm not going in details..

```
org.apache.lucene.search.Query luceneQuery =  
    queryBuilder.phrase()  
        .onField( "description" )  
        .andField( "title" )  
        .sentence( "a book on highly scalable query engines" )  
        .enableFullTextFilter( "ready-for-shipping" )  
        .createQuery();
```

```
CacheQuery cacheQuery =  
    searchManager.getQuery( luceneQuery, Book.class);  
List<Book> objectList = cacheQuery.list();
```

Architecture: simplest approach

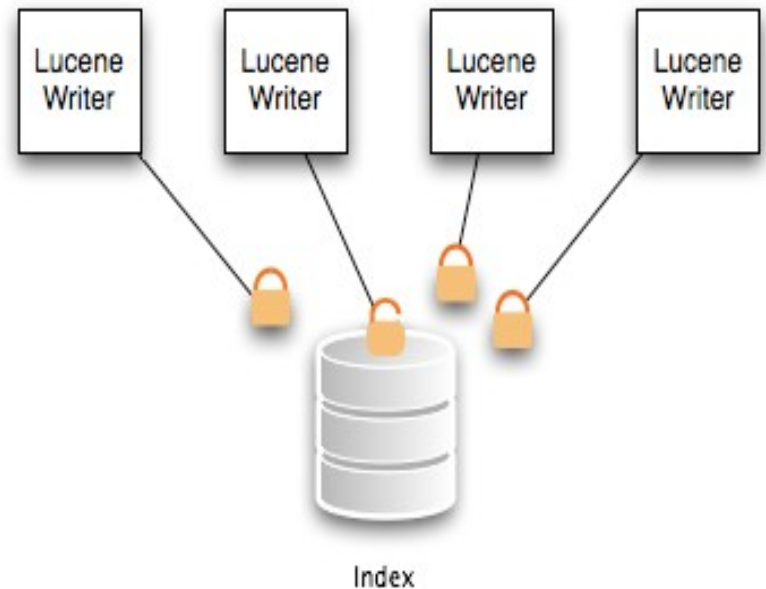


Works for DIST, INVALID and REPL

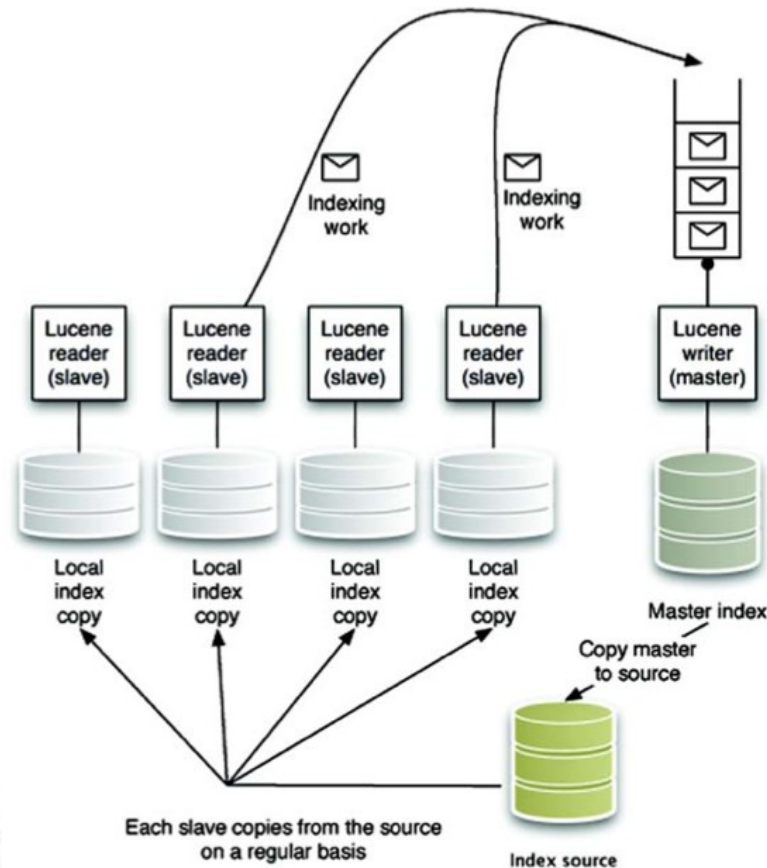
```
sharedIndexes = true  
indexLocalOnly = true
```

Scalability issues

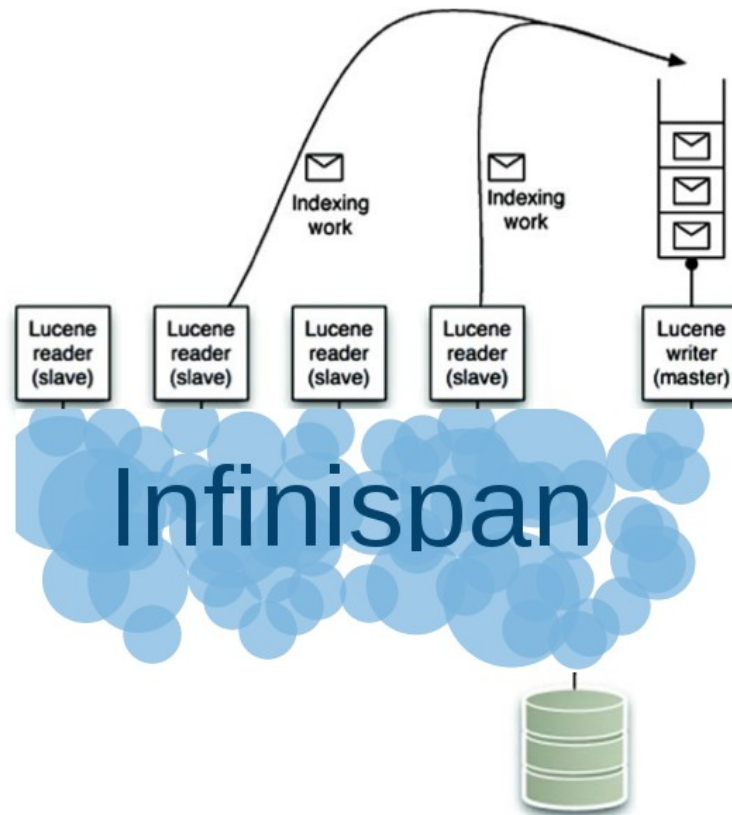
- Global writer locks
- NFS based index sharing *very* tricky



Queue-based clustering (via filesystem)



Index stored in Infinispan

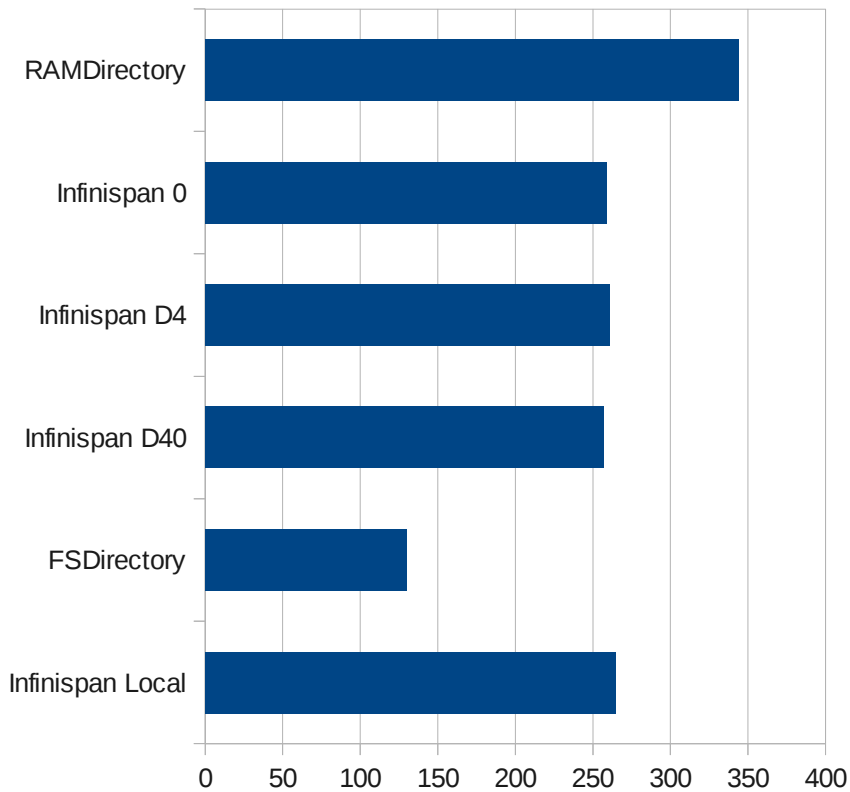


Clustering “native” Lucene access

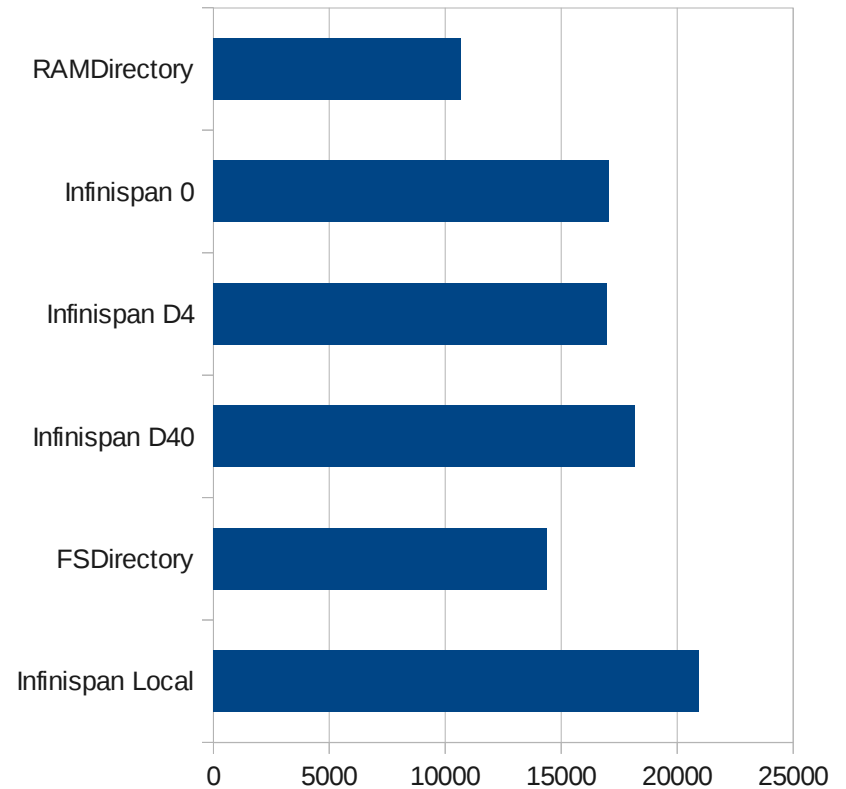
- Using org.apache.lucene directly
 - Distributed on multiple nodes
 - On any cloud

Single node *performance idea*

Write ops/sec

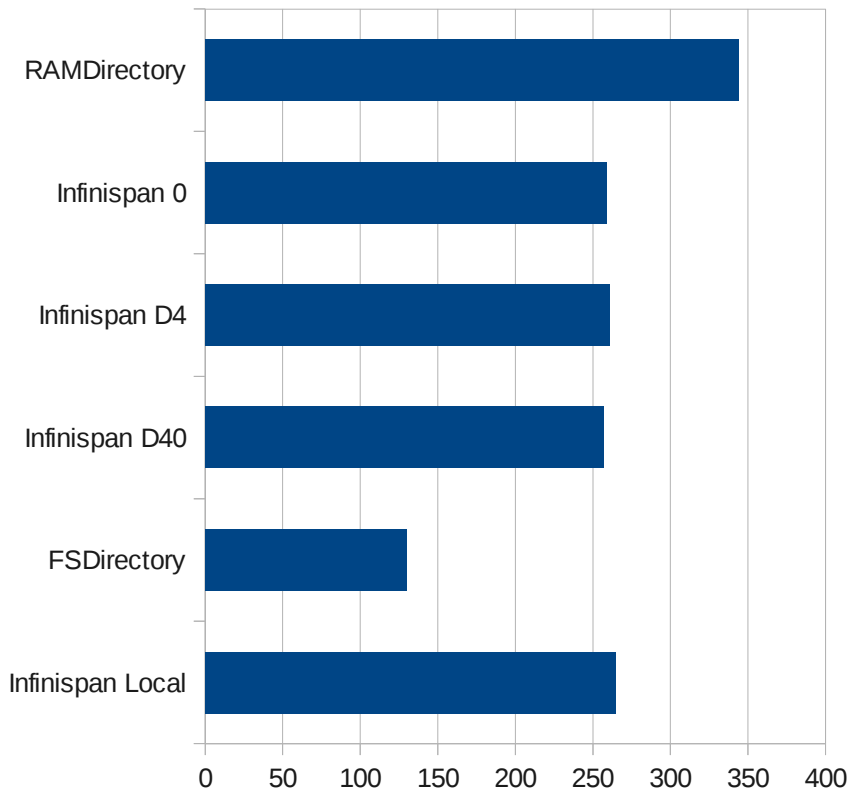


Queries/sec

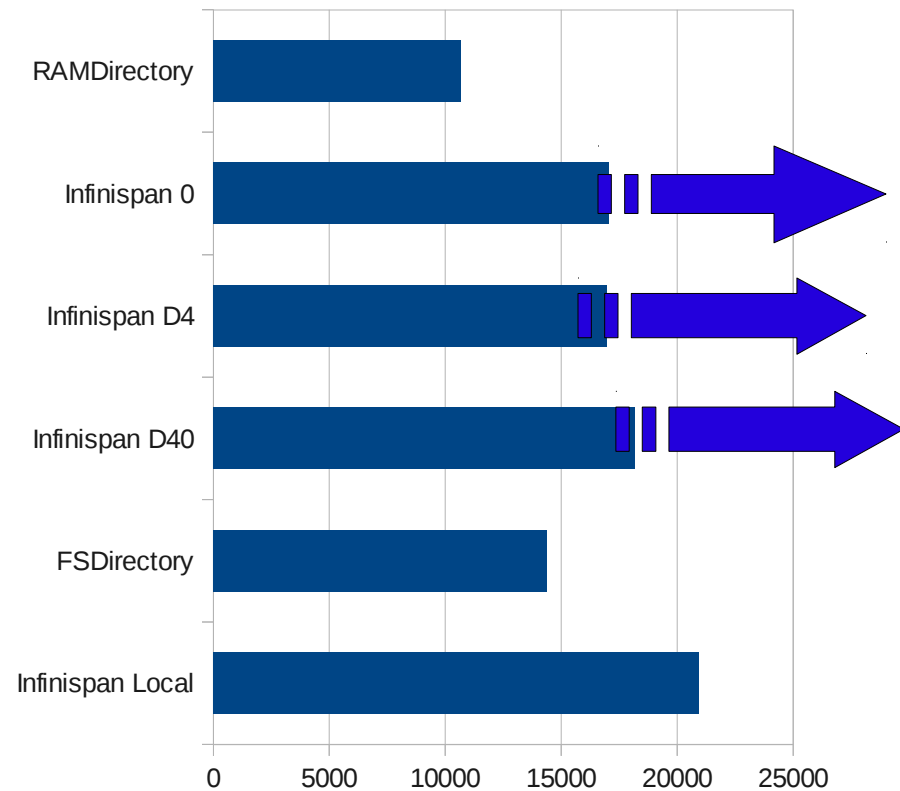


multi-node *performance idea*

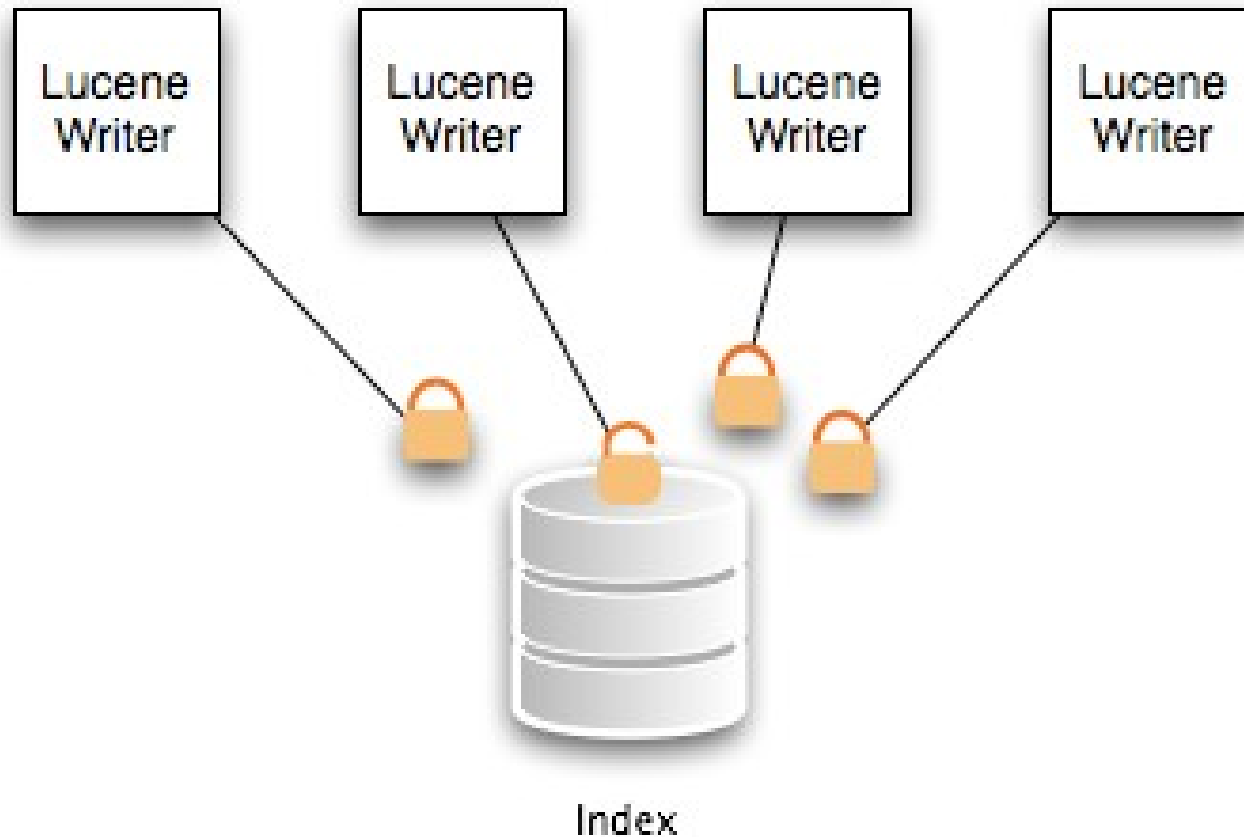
Write ops/sec



Queries/sec



Why does writing not scale?



Performance hints

- Setting Lucene's maximum segment size to fit in LuceneDirectory *chunk_size* will avoid readlocks
- Verify blob sizes fit in JGroups network packets, tune JGroups
- Check for CacheStores “sweet spot” size

Memory requirements

- RAMDirectory: all must fit in a single VM's memory
- FSDirectory: OS does a great caching job – but if it doesn't fit in memory
- Infinispan: comparable to FSDirectory
 - Flexible
 - Fast
 - Network vs. disk

Ingredients for a cloud

One Infinispan to rule them all

- Store Lucene indexes
- Hibernate second level cache
- Application managed cache
- Datagrid
- EJB, session replication in AS7
- As a JPA “store” via Hibernate OGM

Ingredients for a cloud

- JGroups discovery protocol
 - MPING
 - TCP_PING
 - JDBC_PING
 - S3_PING
- Choose a CacheLoader
 - Database based
 - Jclouds (S3, ...)
 - Cassandra

jclouds™

What's next

- Facilitate writing scalability
- Ease configuration aspects for clustering – ergonomics!
- Parallel searching
- A component of
 - <http://www.cloudtm.eu>



Related talks at JUDCon

15:15 – JPA applications in the era of NoSQL and Clouds: **Introducing OGM**



Q&A

Infinispan

<http://infinispan.org>

<http://in.relation.to>

<http://jboss.org>

@Infinispan

@Hibernate

@SanneGrinovero