

# Apache Lucene/Solr Internals

# About me

Java and all around

Principal Software Engineer  
at Grid Dynamics

Kharkiv

[asokolenko@griddynamics.com](mailto:asokolenko@griddynamics.com)



**Grid Dynamics**

Scalable eCommerce Platform Solutions

the magic of  
 macy's<sup>®</sup>  
.com

# Apache Lucene/Solr Internals



June 2013 database

14.630.209 records

VM

16 CPU cores

16 GB memory

4 nodes

× 12GB disk space

Indexing took 5 hours  
in 100 threads  
1000 batch

**lightweight  
performant  
search  
library**





elasticsearch.



# Data Model

- document oriented
- flat
- store
- index

# Data Model

- document oriented
- flat
- store
- index

boost = 1.1

Document  
docID = 23

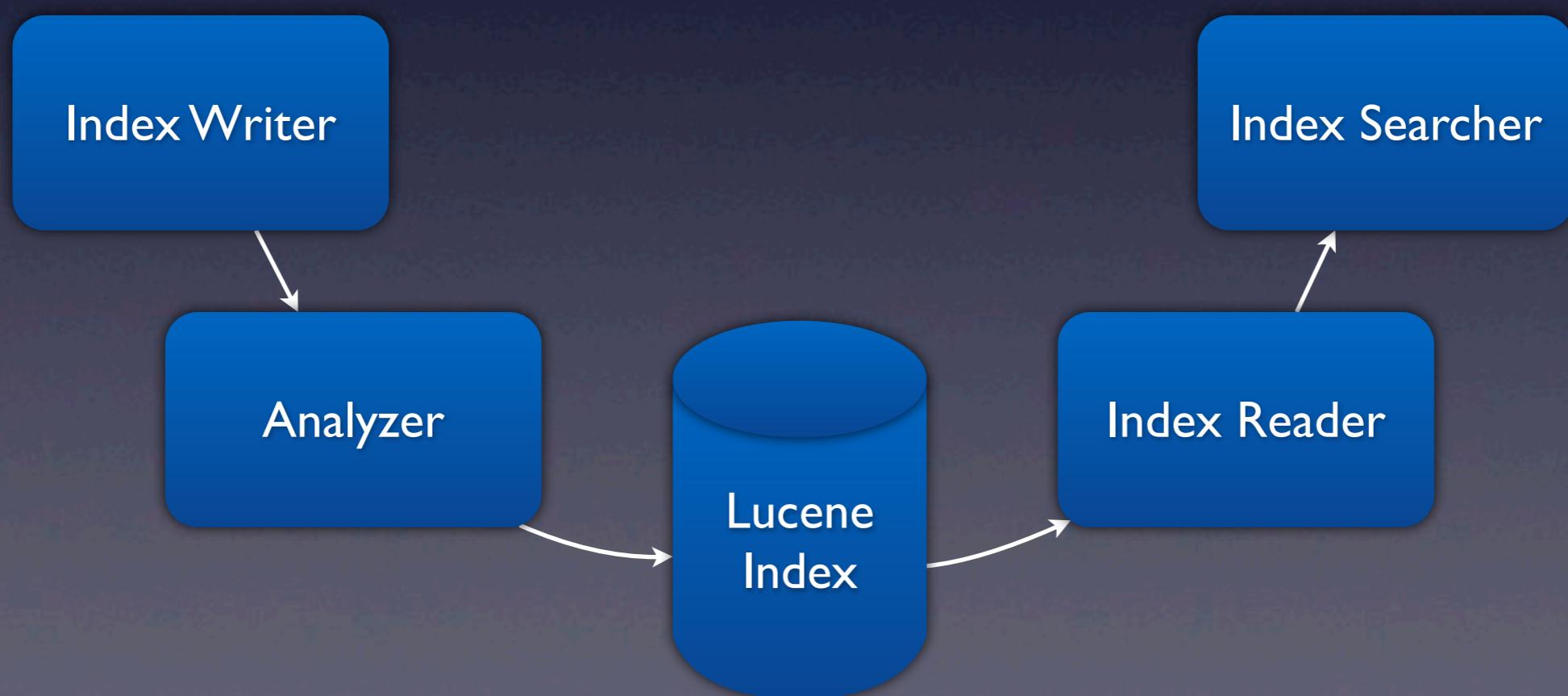
score:1

tag:java

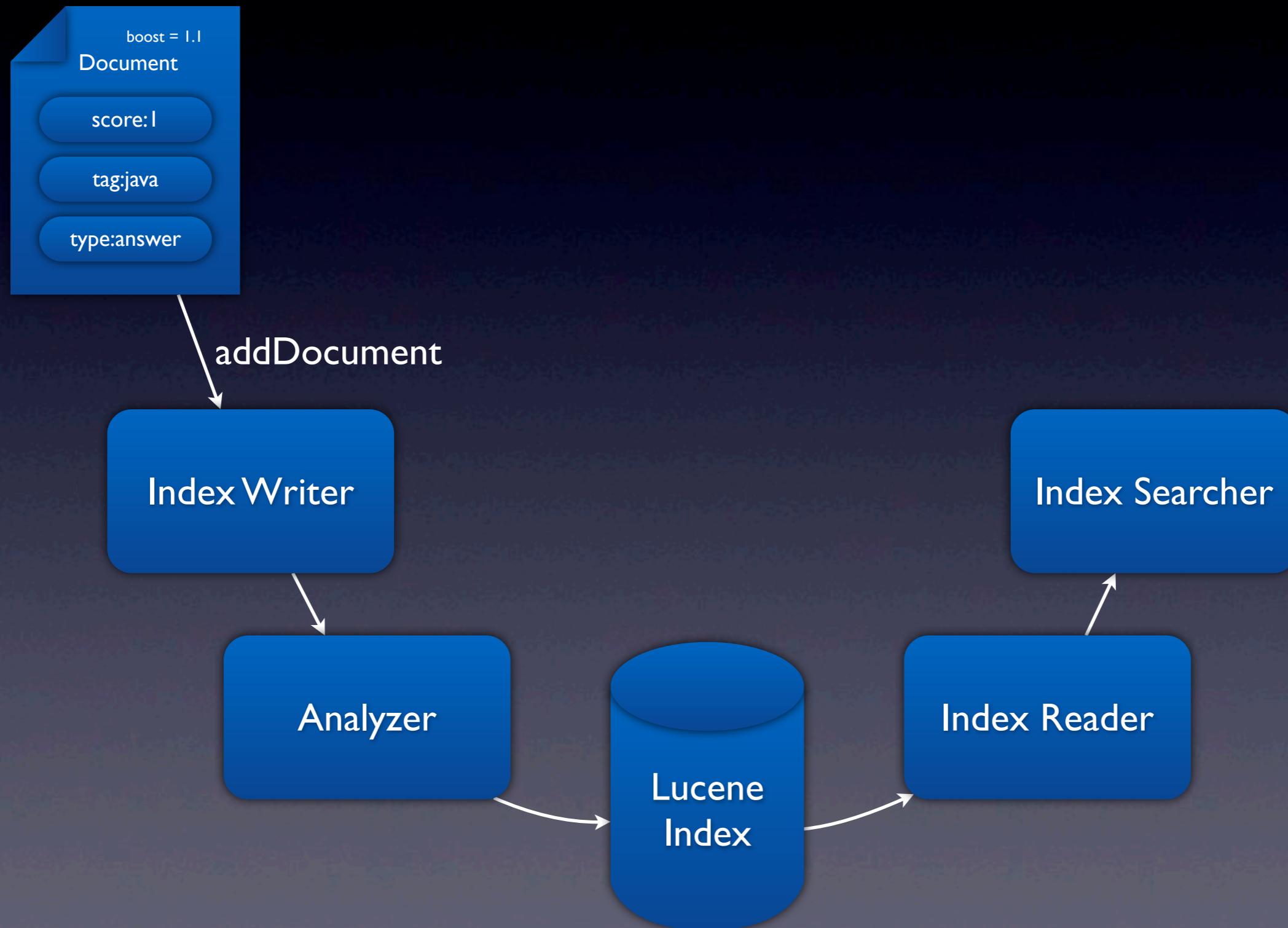
type:answer

# Showcase

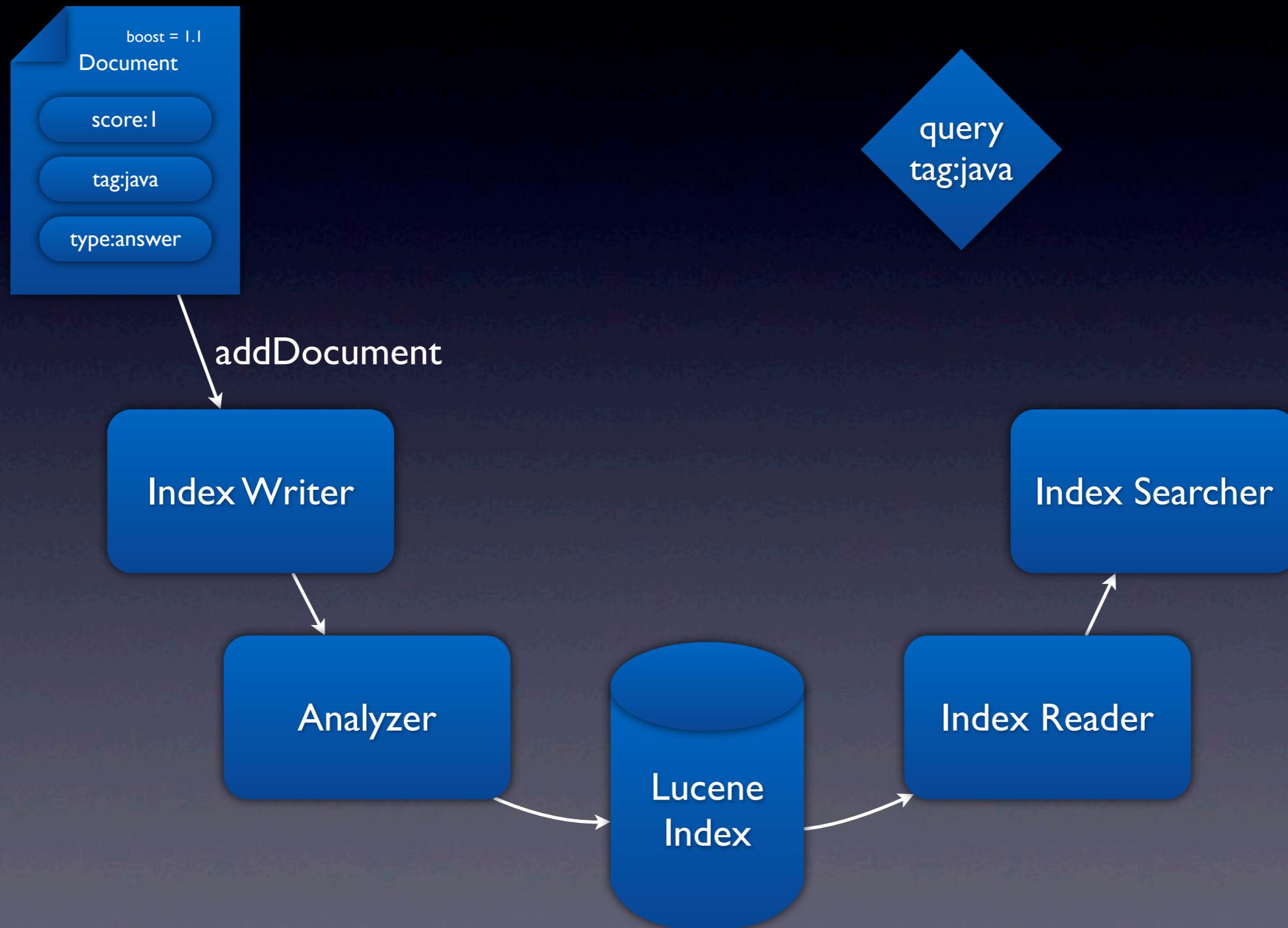
# Basic Flow



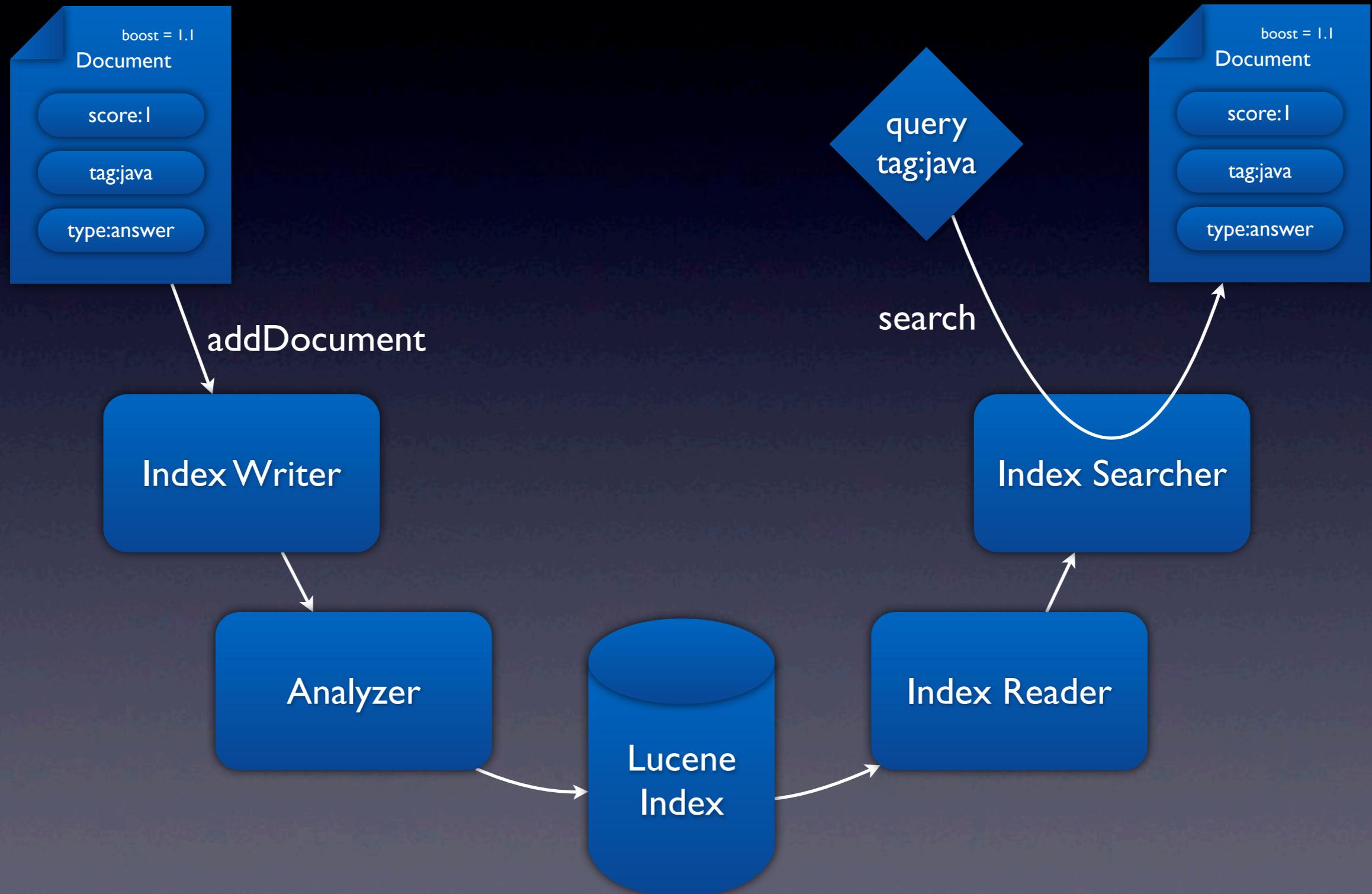
# Basic Flow



# Basic Flow

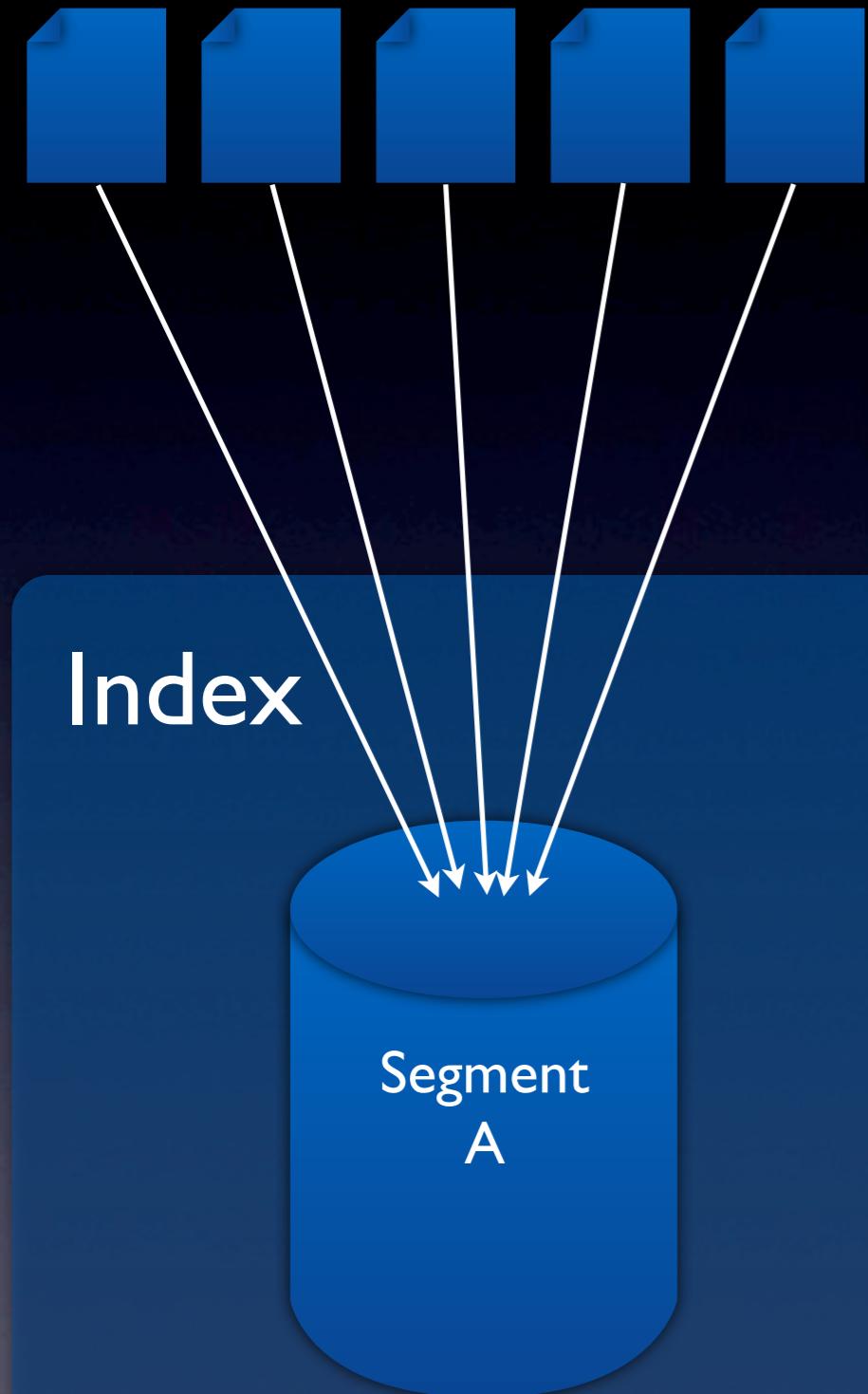


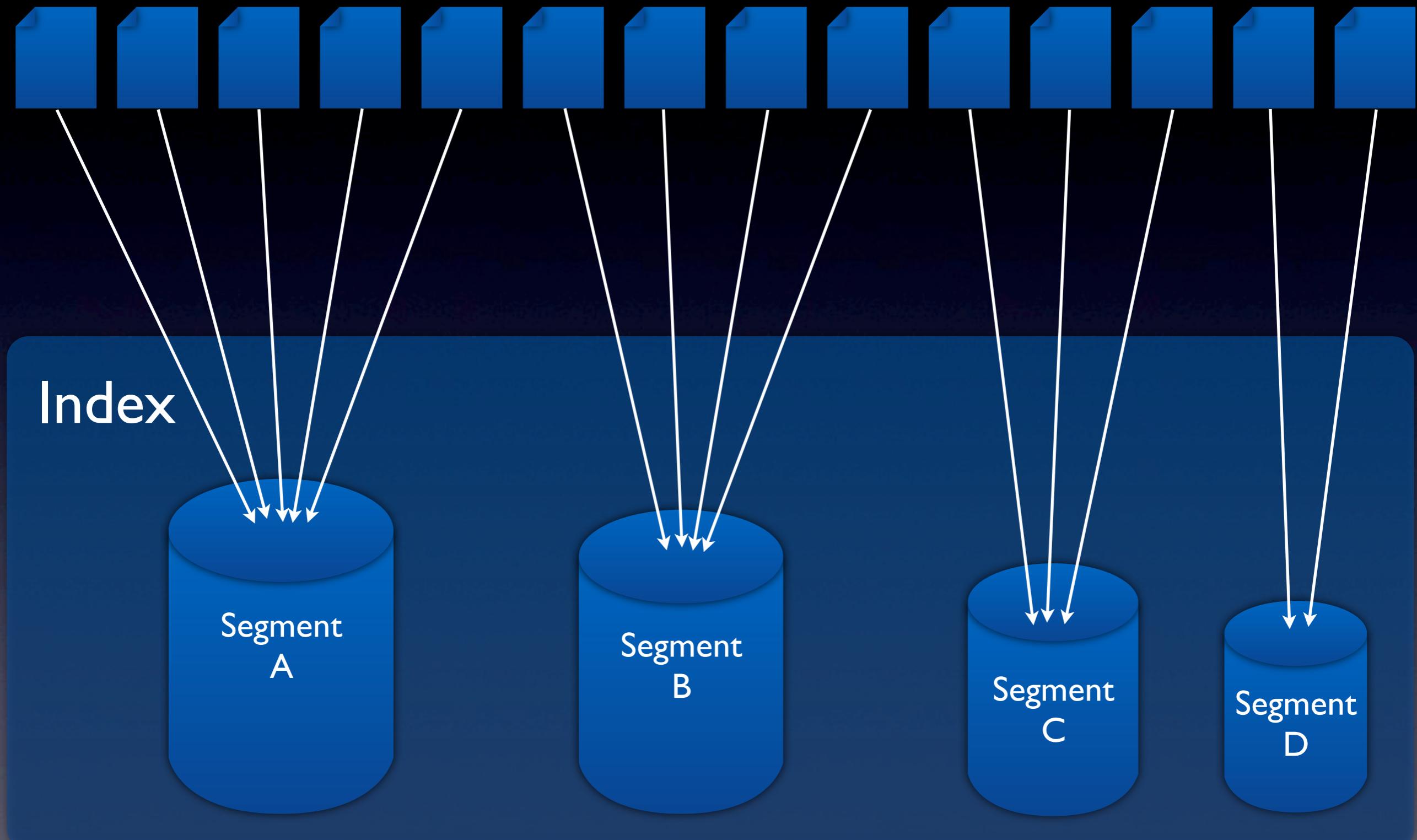
# Basic Flow



# Lucene Index Structure

# Index







# Term Infos

score:0

3

score:1

4

score:5

2

...

tag:java

2

tag:mysql

3

tag:css

4

...

type:answer

3

type:question

2

## Term Infos

score:0	3
score:1	4
score:5	2
...	
tag:java	2
tag:mysql	3
tag:css	4
...	
type:answer	3
type:question	2

## Term Frequencies

3 <sup>1</sup>	+1 <sup>1</sup>	+2 <sup>1</sup>
10 <sup>1</sup>	+3 <sup>1</sup>	+1 <sup>1</sup>
+7 <sup>1</sup>		
4 <sup>1</sup>	+11 <sup>1</sup>	
5 <sup>1</sup>	+2 <sup>1</sup>	
6 <sup>1</sup>	+52 <sup>1</sup>	+1 <sup>1</sup>
<sup>1</sup>	+30 <sup>1</sup>	+27 <sup>1</sup>
		+2 <sup>1</sup>
3 <sup>1</sup>	+7 <sup>3</sup>	+1 <sup>5</sup>
5 <sup>2</sup>	+2 <sup>1</sup>	

## Term Info Index

## Term Infos

## Term Frequencies

score:0 3

...

tag:mysql 3

...

type:question 2

score:0 3

score:1 4

score:5 2

...

tag:java 2

tag:mysql 3

tag:css 4

...

type:answer 3

type:question 2

3 1 + 1 1 + 2 1

10 1 + 3 1 + 1 1 + 7 1

4 1 + 1 1

5 1 + 2 1

6 1 + 52 1 + 1 1

1 1 + 30 1 + 27 1 + 2 1

3 1 + 7 3 + 1 5

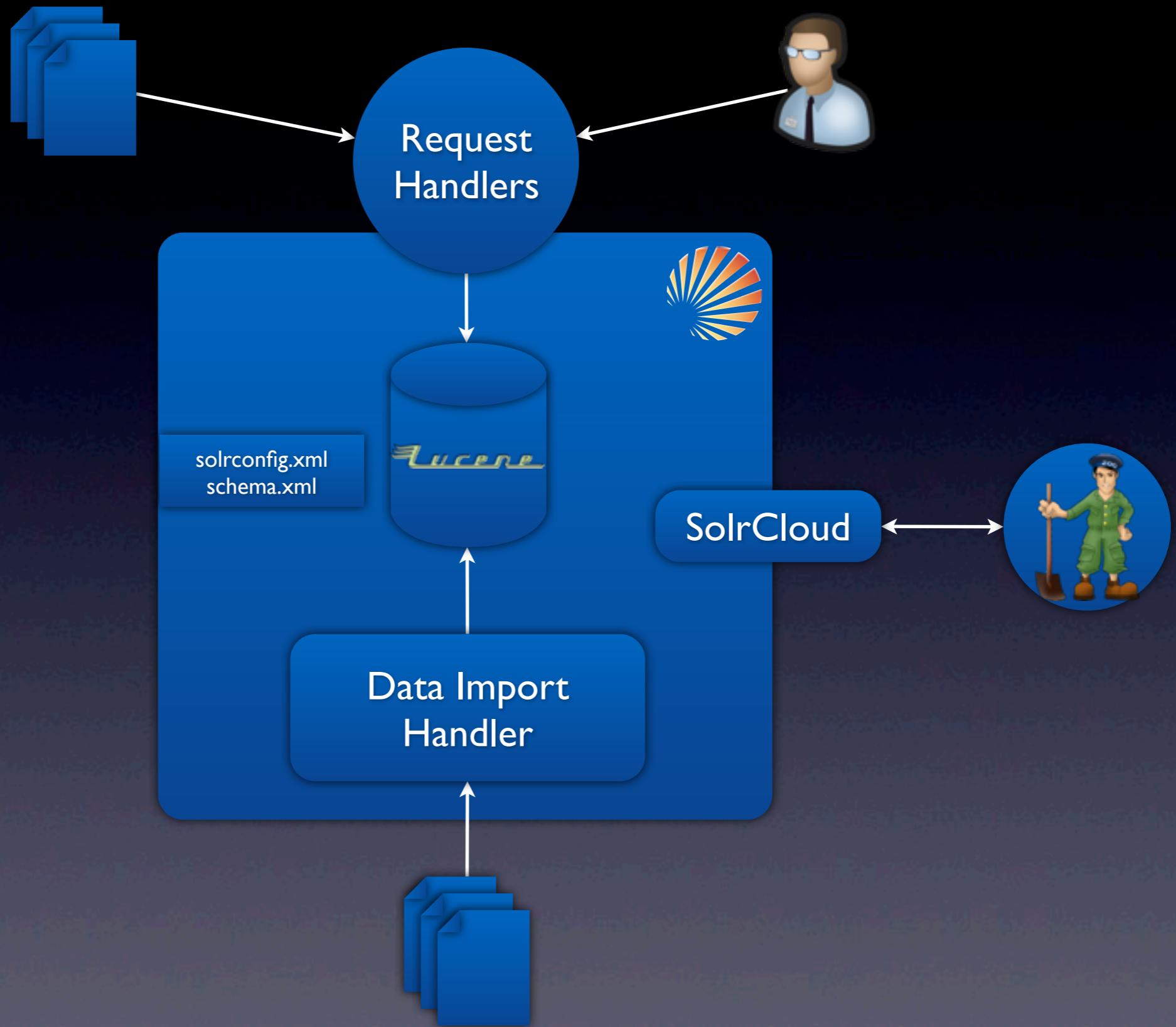
5 2 + 2 1

# Showcase

scalable  
enterprise search  
server

Apache  
**Solr**





# SolrCloud

# Join Cluster



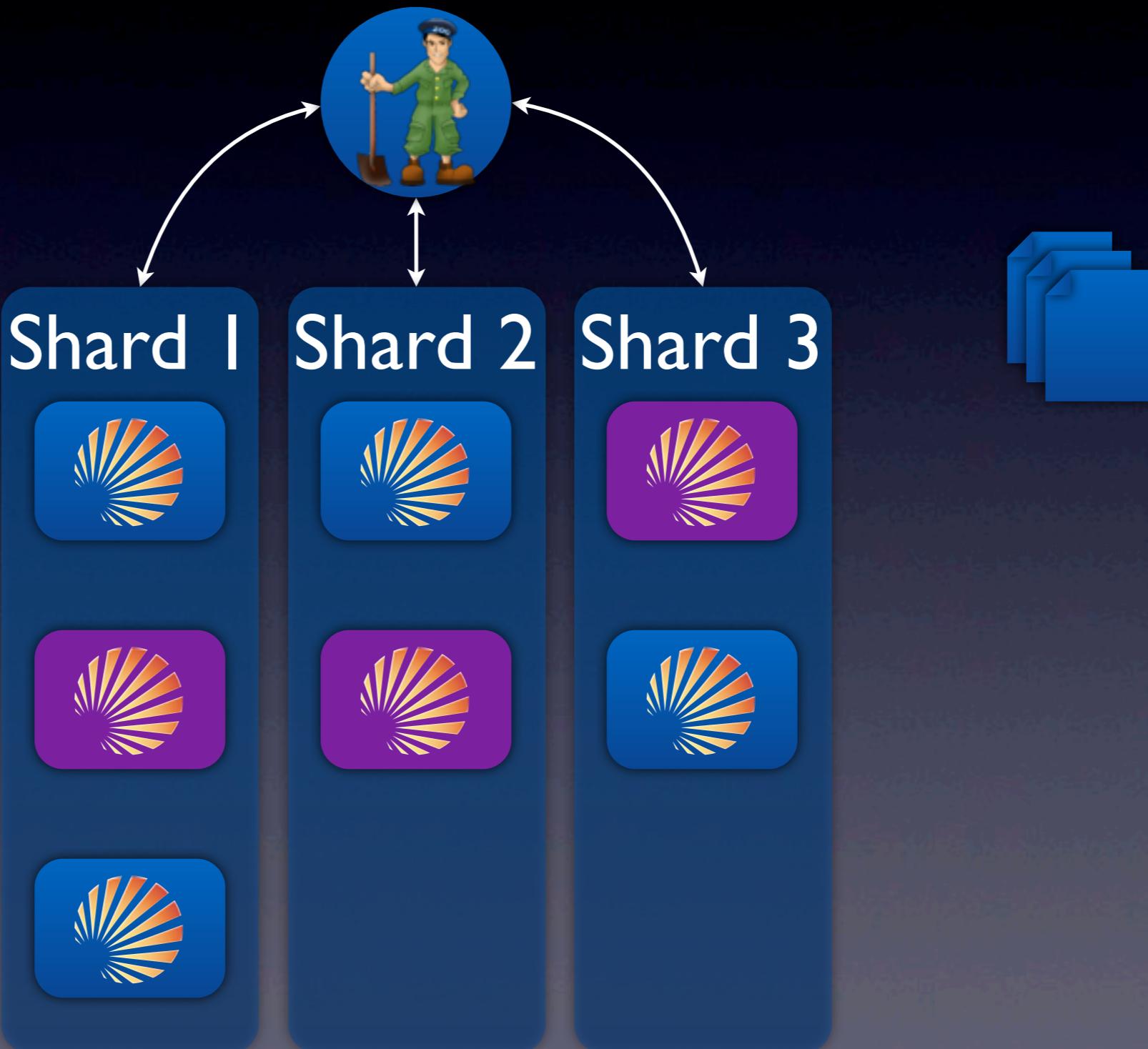
# Join Cluster



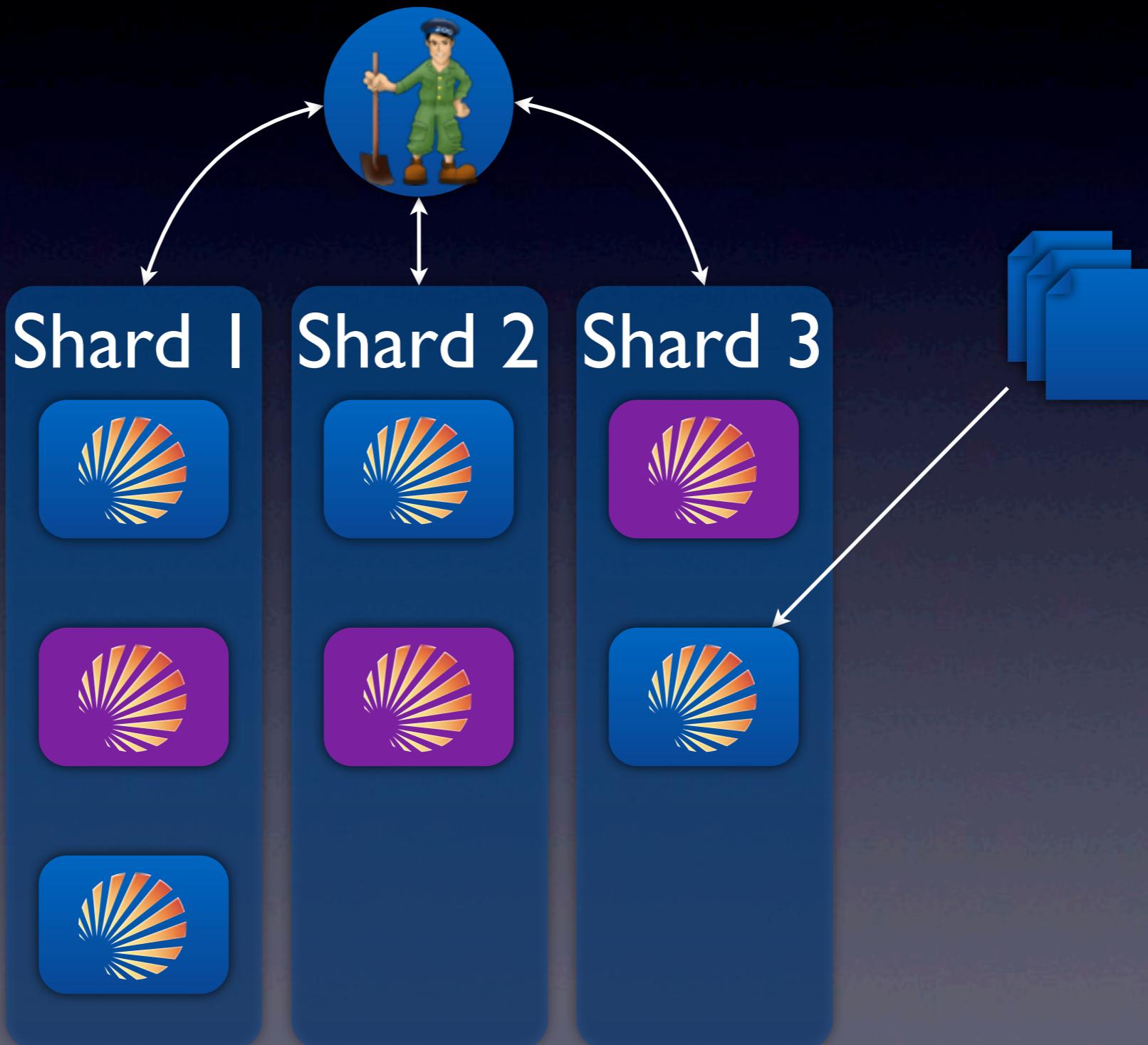
# Indexing



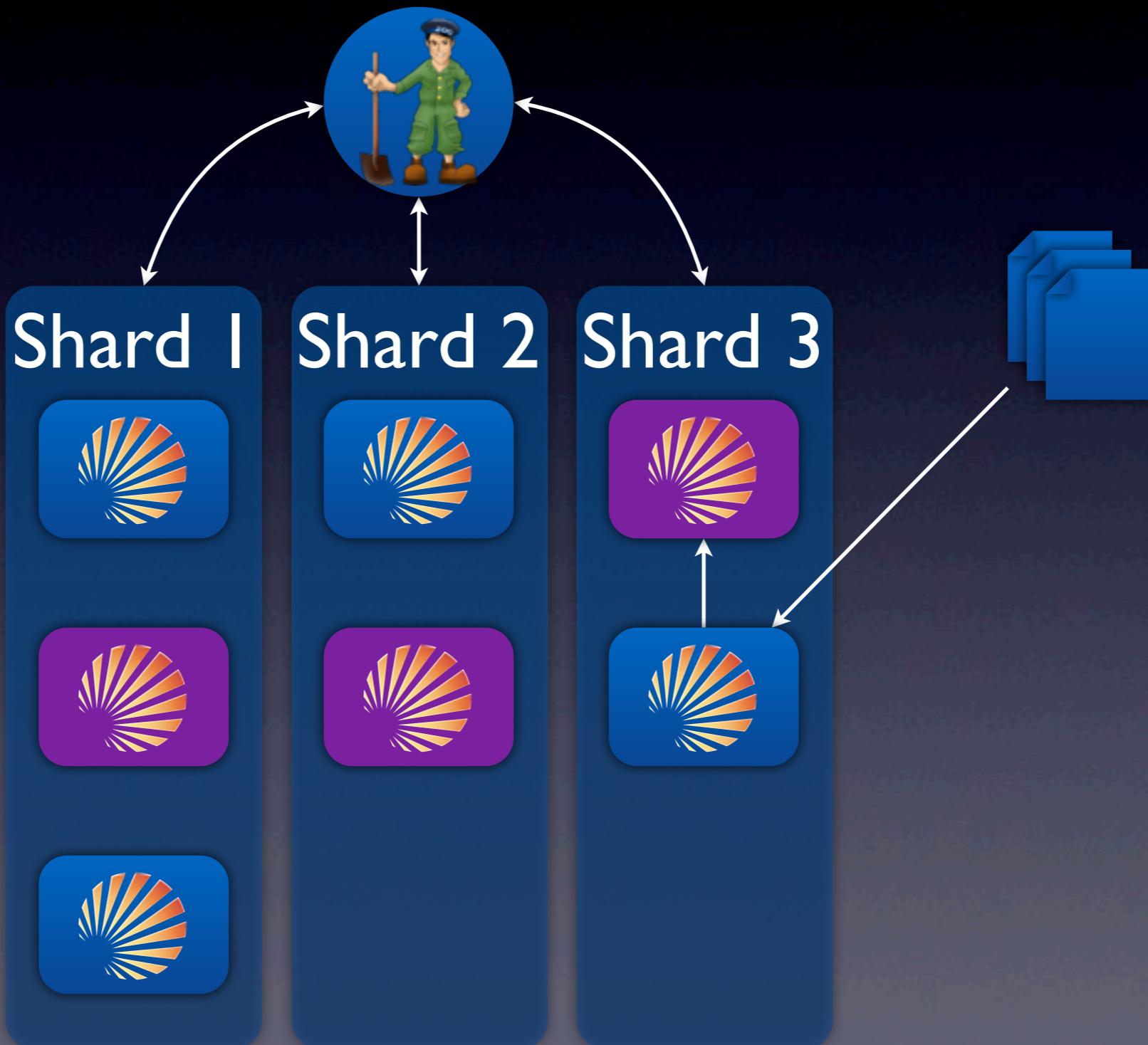
# Indexing



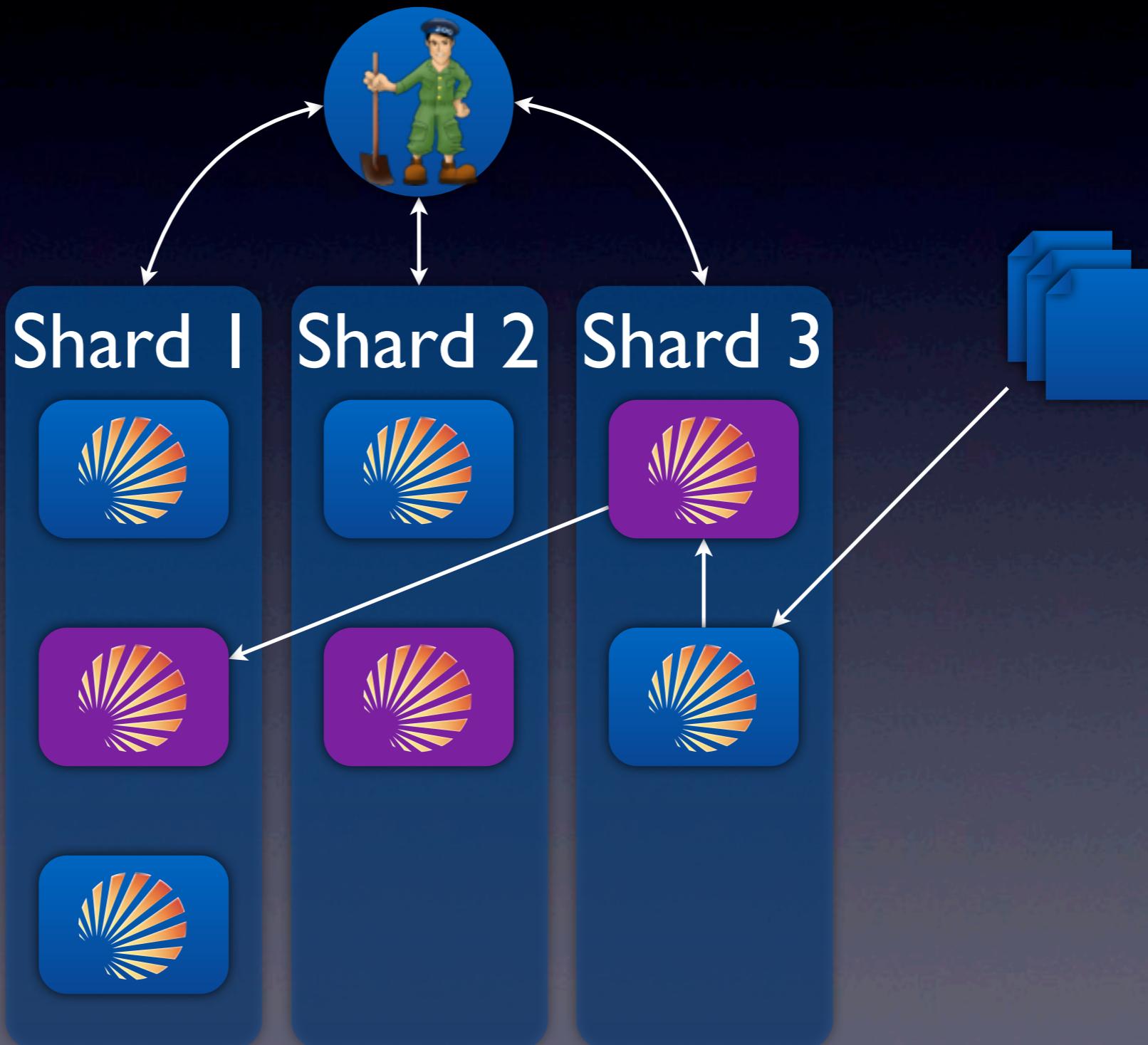
# Indexing



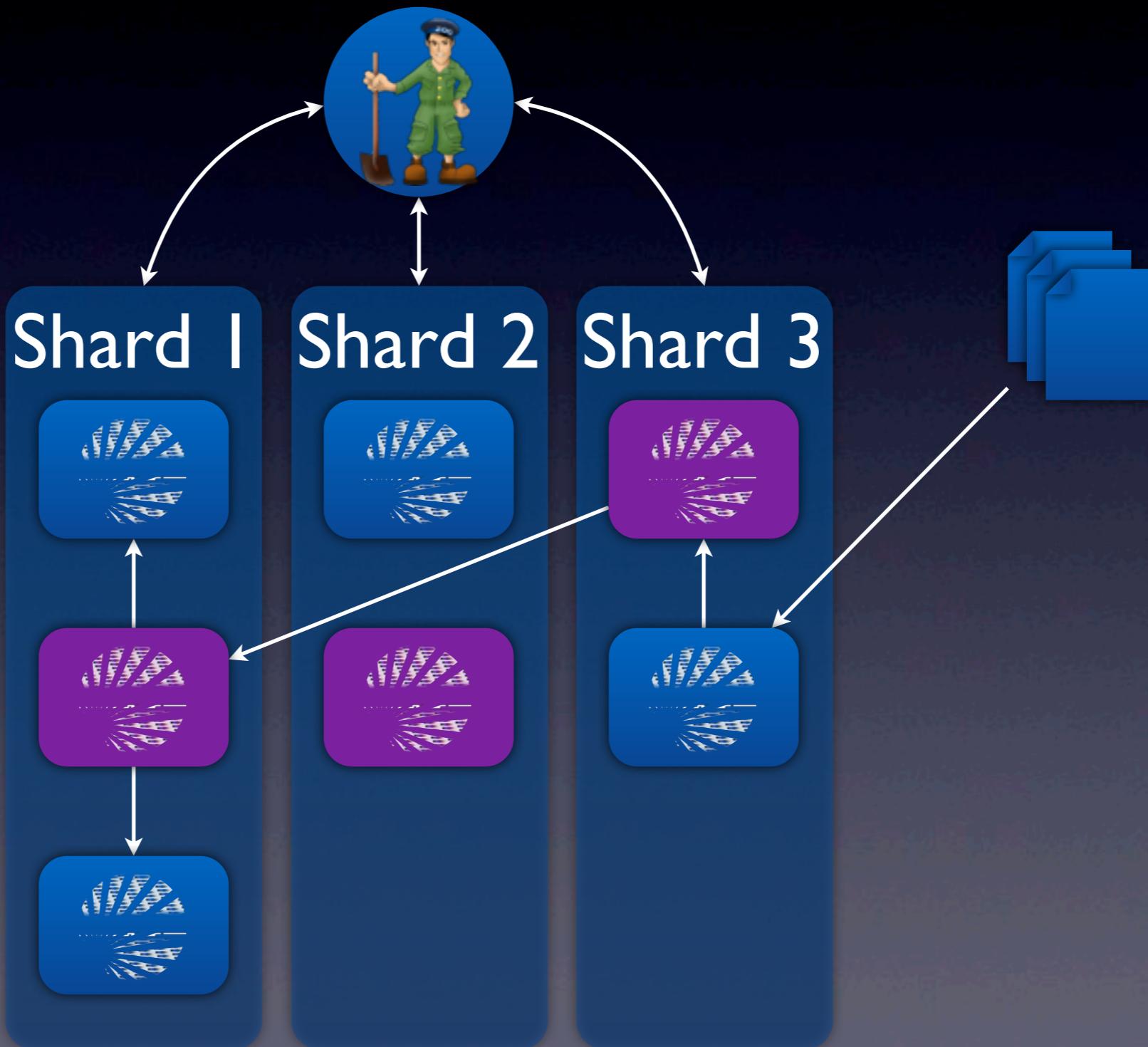
# Indexing



# Indexing



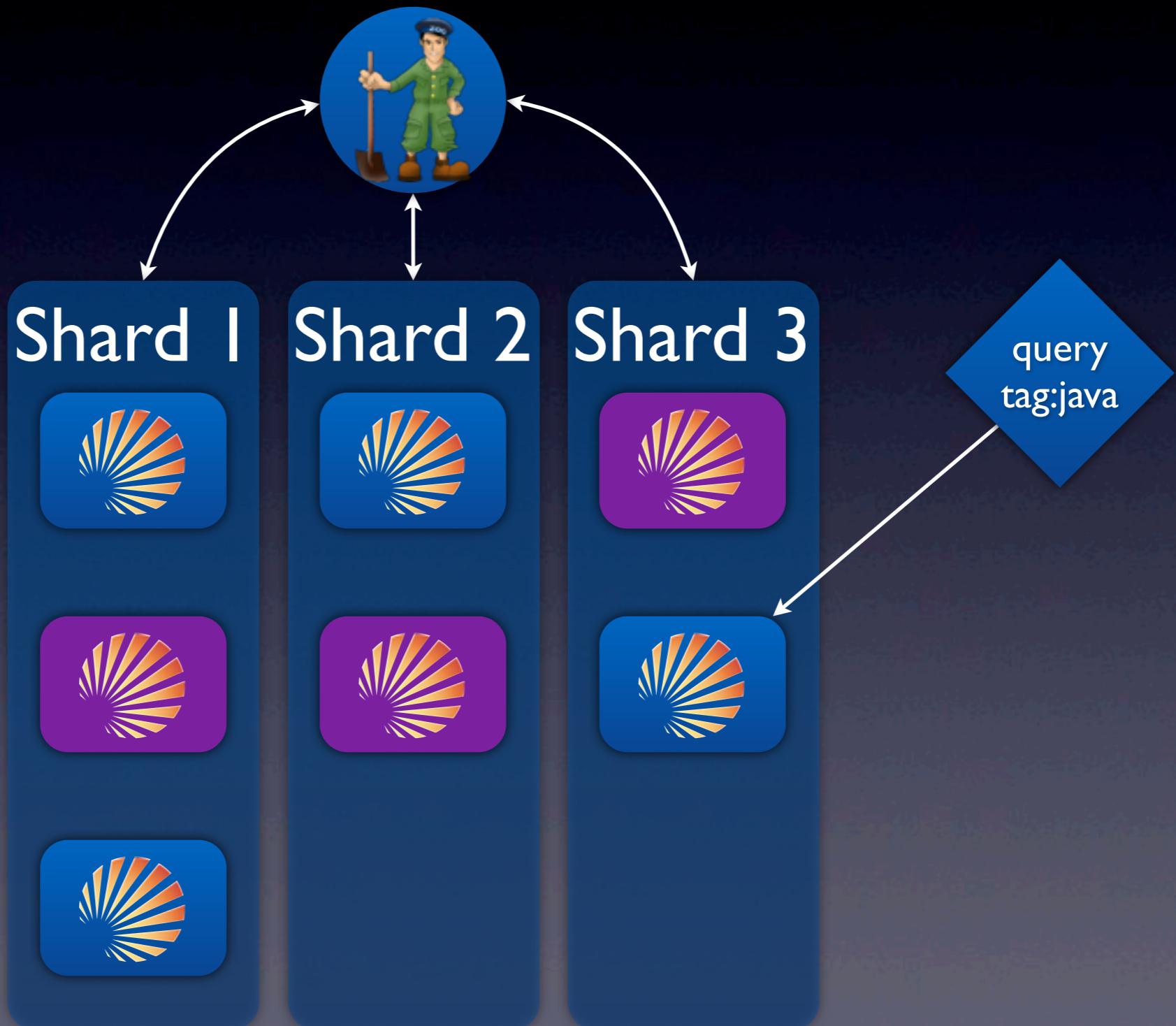
# Indexing



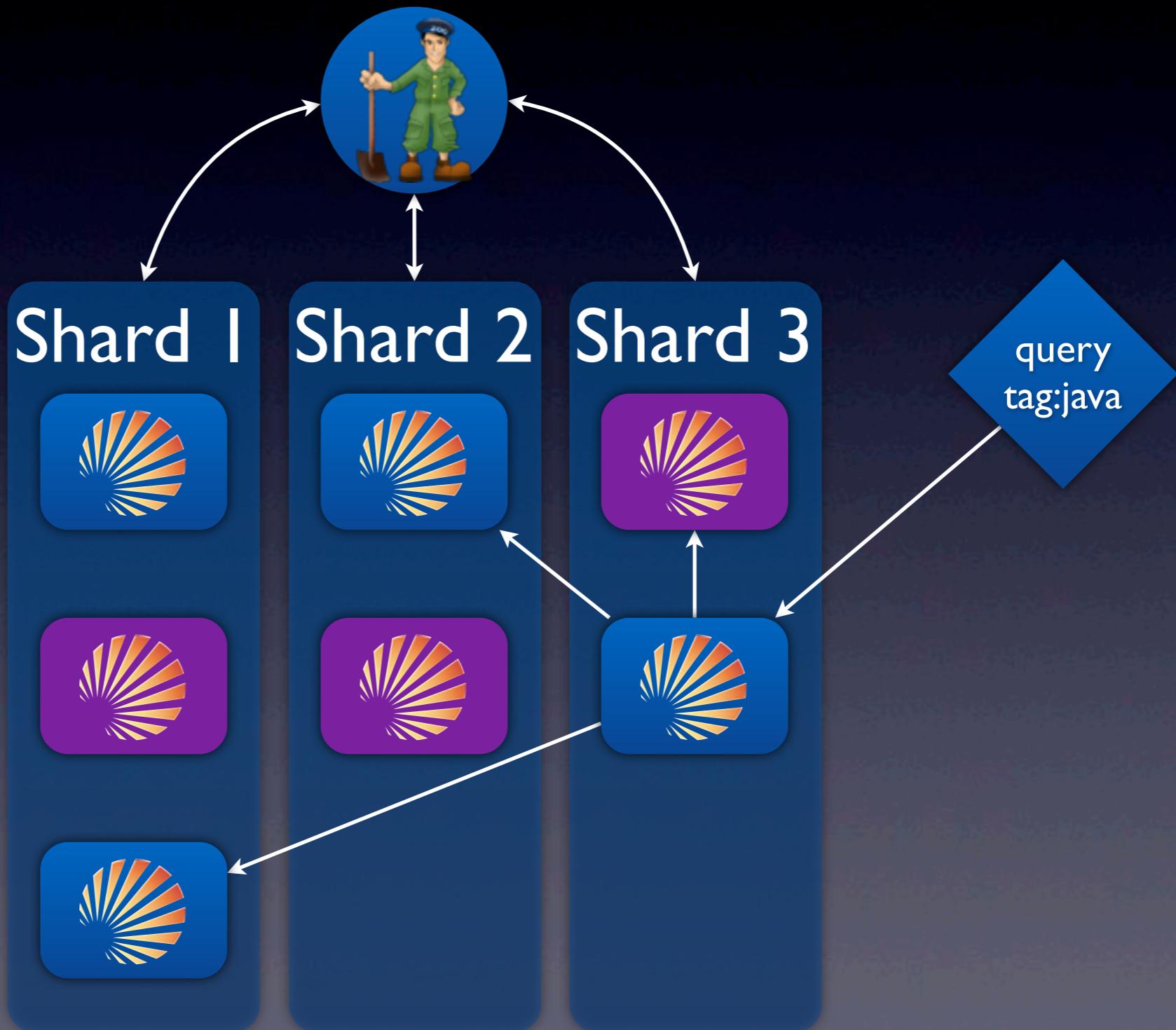
# Query



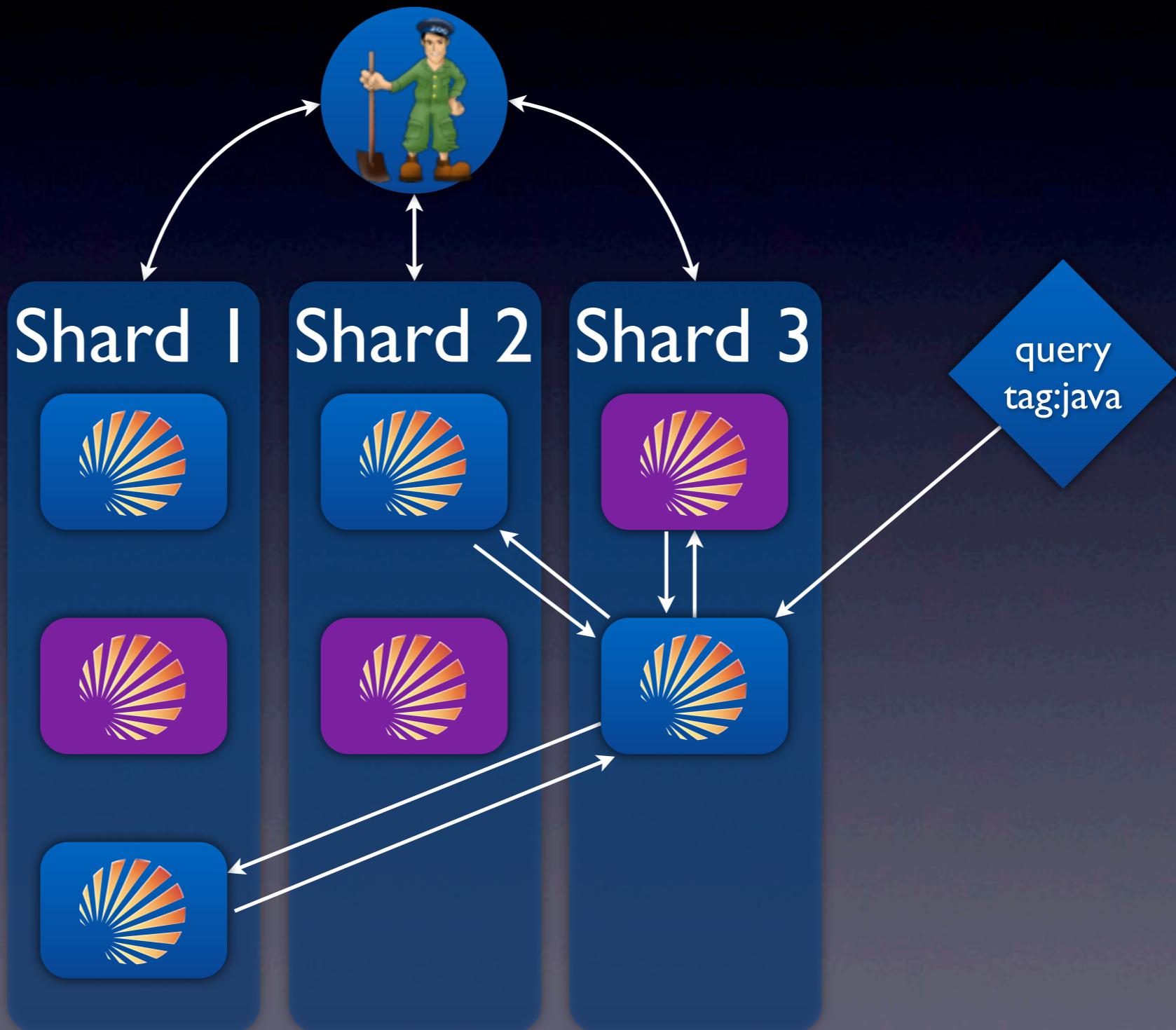
# Query



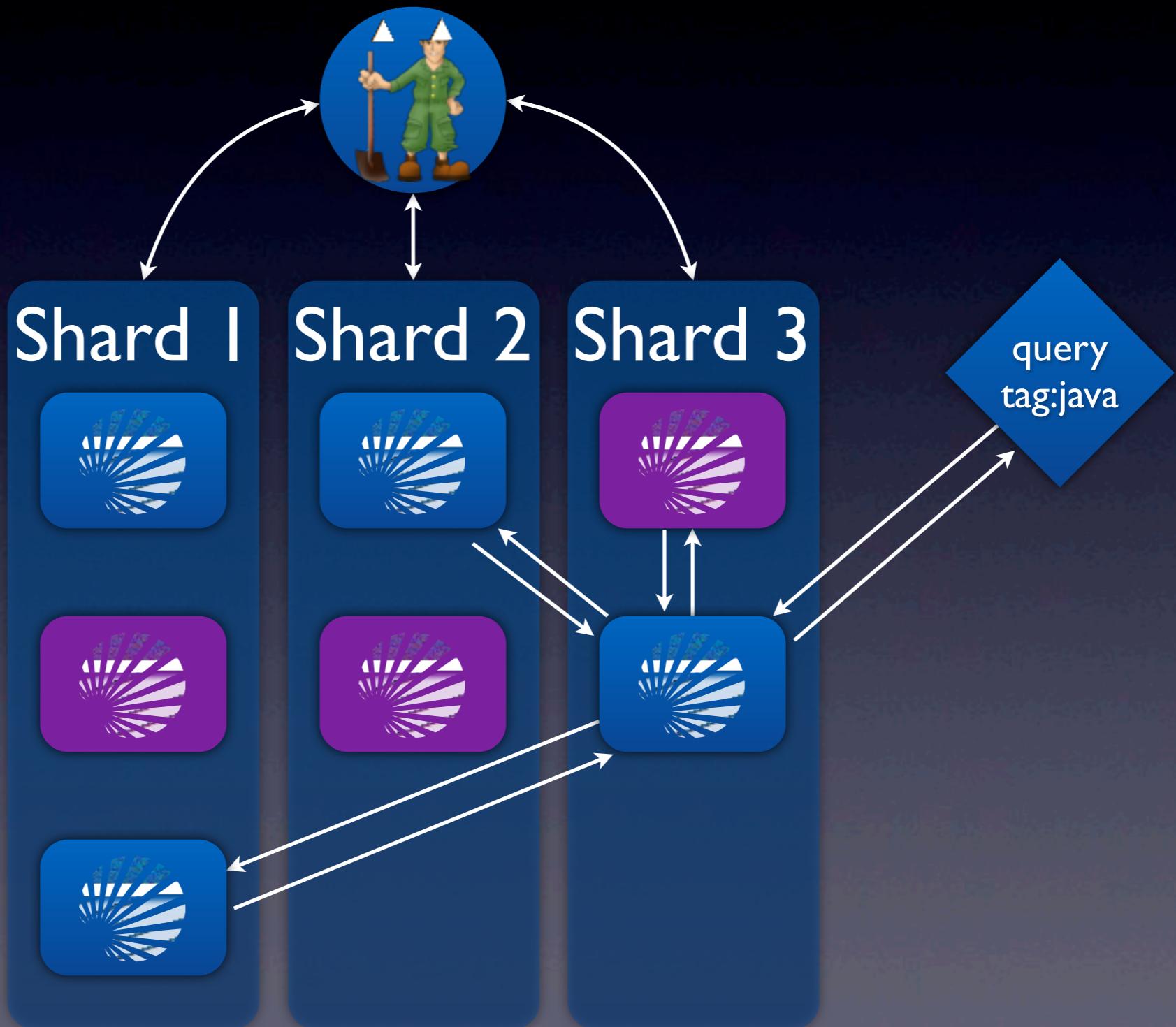
# Query



# Query



# Query



# Failure



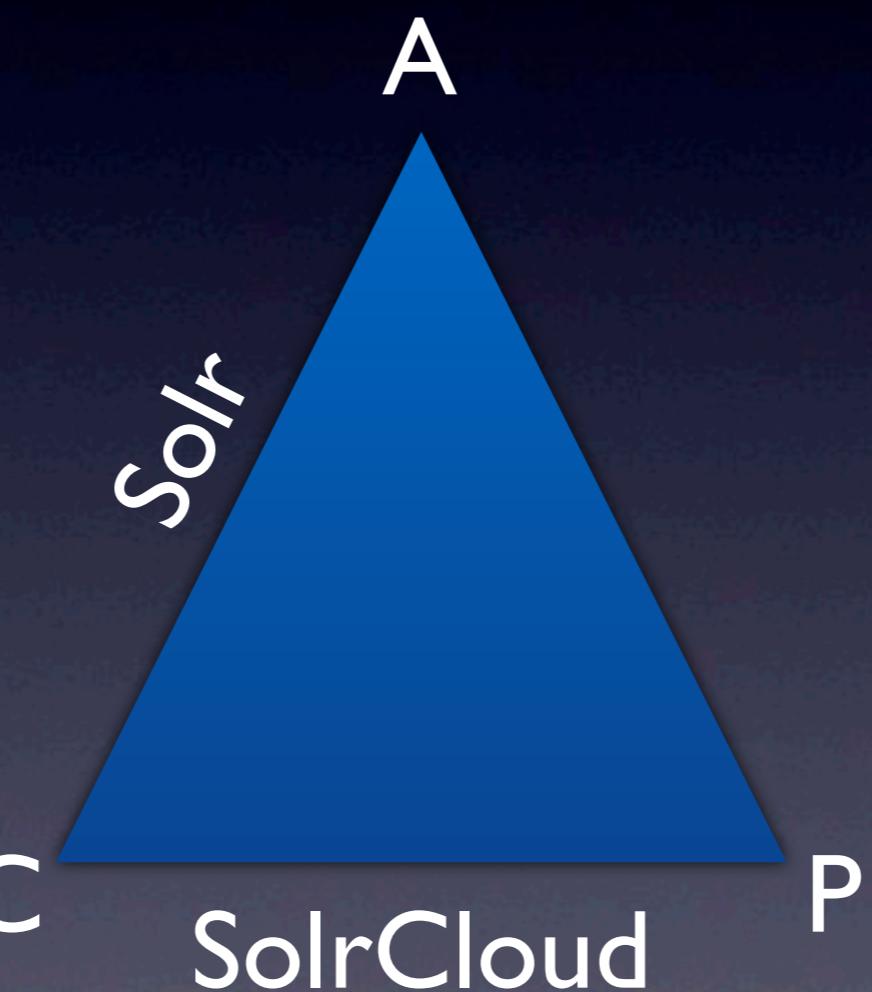
# Failure



# Failure



# CAP Model



# Showcase

# Faceted Navigation

# Showcase

# Algorithm

Query Result

7 31 58 59

Index

tag:java

5 +2

tag:mysql

6 +52 +1

tag:css

1 +30 +27 +2

# Algorithm

Query Result

7 31 58 59

Index

tag:java

5 7

tag:mysql

6 58 59

tag:css

1 31 58 60

# Algorithm

Query Result

7 31 58 59

Index

tag:java

5 7

Facet

1

tag:mysql

6 58 59

2

tag:css

1 31 58 60

2

# Showcase

# Text Analysis

# Analyzer



Char filter



Tokenizer



Filter

# Analyzer

Index time

Char filter

Tokenizer

Filter



# Analyzer

Index time

<strong>There are no pointers in  
Java!</strong>

Char filter

Tokenizer

Filter

# Analyzer

Index time

<strong>There are no pointers in Java!</strong>

Char filter

There are no pointers in Java!

Tokenizer

Filter

# Analyzer

Index time

<strong>There are no pointers in Java!</strong>

Char filter

There are no pointers in Java!

Tokenizer

There no in  
are pointers Java

Filter

# Analyzer

Index time

<strong>There are no pointers in Java!</strong>

Char filter

There are no pointers in Java!

Tokenizer

There no in  
are pointers Java

Filter

? ? ?  
? pointer java

# Analyzer

Index time

Query time

<strong>There are no pointers in Java!</strong>

Char filter

There are no pointers in Java!

Tokenizer

There no in  
are pointers Java

Filter

? ? ?  
? pointer java

# Analyzer

Index time

Query time

<strong>There are no pointers in Java!</strong>

pointers in Java

Char filter

There are no pointers in Java!

Tokenizer

There no in  
are pointers Java

Filter

? ? ?  
? pointer java

# Analyzer

Index time

Query time

<strong>There are no pointers in Java!</strong>

pointers in Java

Char filter

There are no pointers in Java!

pointers in Java

Tokenizer

There no in  
are pointers Java

Filter

? ? ?  
? pointer java

# Analyzer

Index time

Query time

<strong>There are no pointers in Java!</strong>

pointers in Java

Char filter

There are no pointers in Java!

pointers in Java

Tokenizer

There no in  
are pointers Java

pointers in Java  
in

Filter

? ? ?  
? pointer java

# Analyzer

Index time

Query time

<strong>There are no pointers in Java!</strong>

pointers in Java

Char filter

There are no pointers in Java!

pointers in Java

Tokenizer

There no in  
are pointers Java

pointers in Java

Filter

? ? ?  
? pointer java

pointer java  
?

# Showcase

# Spell Suggestions

# Levenshtein Distance

Levenshtein  
distance = 1

html

htmm

Levenshtein  
distance = 2

html

hlmz

tag:apache

tag:c#

tag:html

tag:java

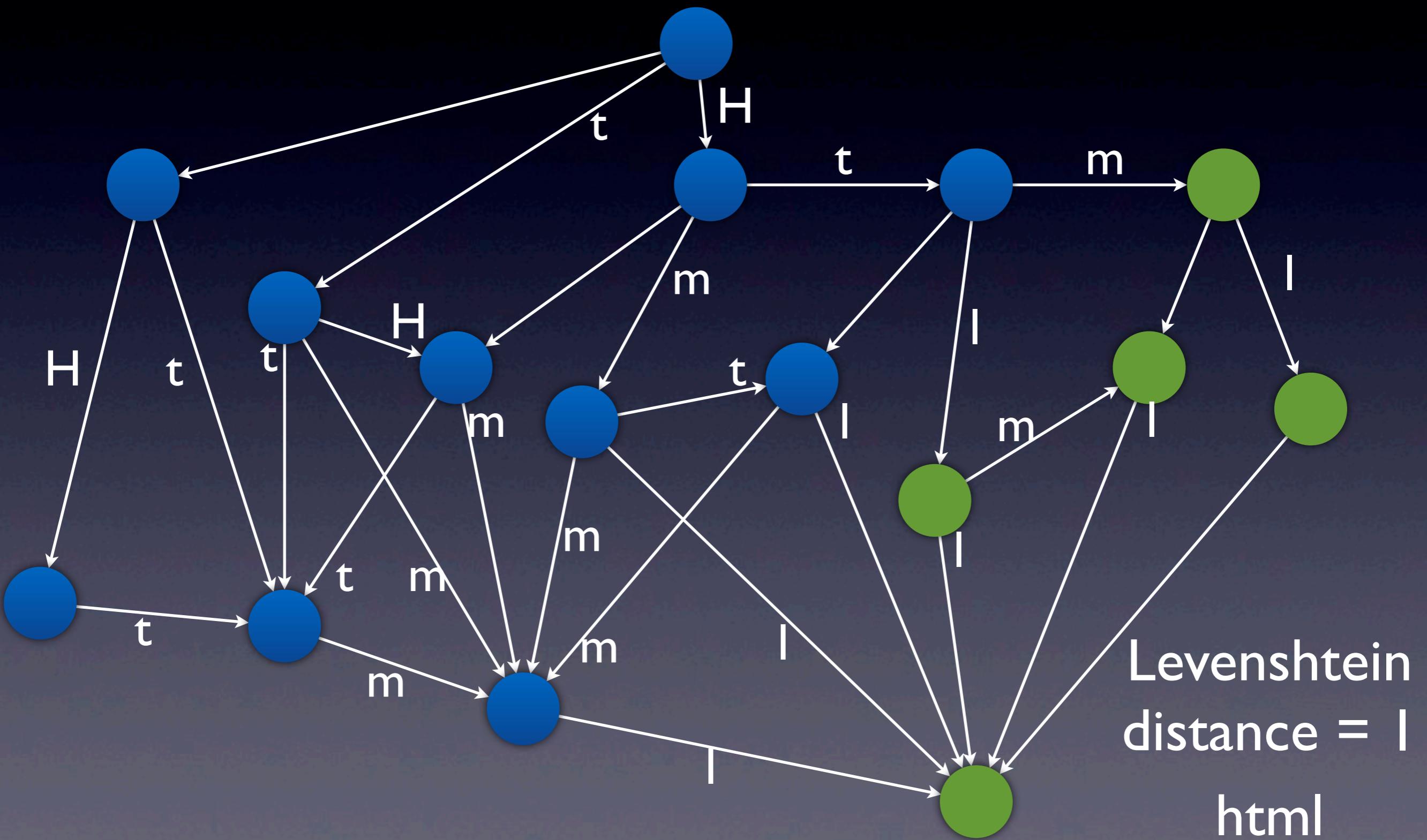
tag:jquery

tag:json

tag:osx

tag:php

# Levenshtein Automaton



# Showcase

# Solr is...

- enterprise level search engine
- vertically scalable
- horizontally scalable, but...
- tunable
- poorly documented
- with active community

# References



- <http://blog.mikemccandless.com>
- [http://lucene.apache.org/core/4\\_3\\_1/index.html](http://lucene.apache.org/core/4_3_1/index.html)
- Introduction to Information Retrieval  
<http://nlp.stanford.edu/IR-book/>
- <http://wiki.apache.org/solr/>
- <https://cwiki.apache.org/confluence/display/solr/Apache+Solr+Reference+Guide>



# Q&A



<http://flip.it/whFqy>



[asokolenko@griddynamics.com](mailto:asokolenko@griddynamics.com)



<http://twitter.com/AnatolSokolenko>

The End