

# LuceneRDD

## Entity Linkage & (Spatial) Search

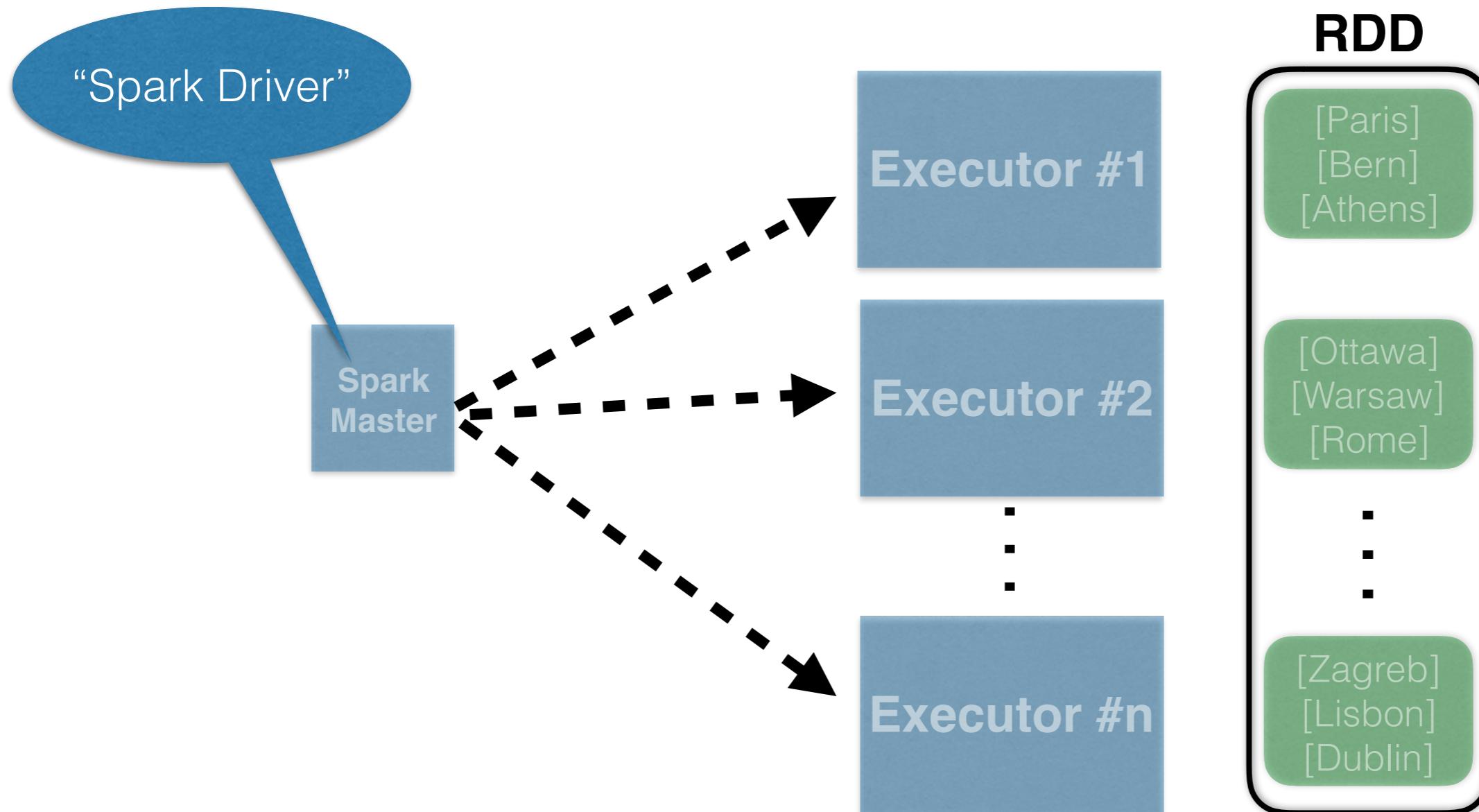
Anastasios Zouzias  
ScalaiO @ Lyon 2016

# About myself

- \* Data Scientist @ Swisscom, CH: Mobility Insights Team
- \* Researcher @ IBM Research, Zurich: Machine Learning & Search Engineer
- \* PhD in Univ. of Toronto: ML & Randomized Algorithms
- \* Scala Developer (day job): 2-3 years
- \* Apache Solr / Elasticsearch developer: 3+ years

How many of you  
have used Spark?

# Apache Spark



**RDD (Resilient Distributed Dataset)**  
**Distributes** data over multiple nodes

How many of you  
have used  
Lucene/Solr/Elastic?

# Apache Lucene

- \* **Lucene** is a powerful Java search library that lets you easily do **search** or **Information Retrieval (IR)**
- \* Used by LinkedIn, Twitter, and many more...
- \* Scalable and High-performance indexing
- \* Powerful, Accurate and Efficient Search algorithms



# LuceneRDD

What is it?

## Spark:

Partitioning / distribution  
of queries / data

## LuceneRDD:

Query / dispatch /  
aggregation of data

## Lucene:

Indexing / querying  
data (single partition)

Open-source project  
<https://github.com/zouzias/spark-lucenerdd>

# Motivation

Why LuceneRDD exists?

- \* **Natively** support of **full-text / spatial** search in Spark (without external Elastic/Solr cluster)
- \* **Scalable** Entity Linkage (approx. join) with Spark & Lucene
- \* Personal: Better understanding of Spark's Internals (RDDs)
- \* Personal: Better understanding of Scala (implicits)

# LuceneRDD: RDD with search

**LuceneRDD**

Full-text search & Entity Linkage

**FacetedLuceneRDD**

LuceneRDD + faceted search

**ShapeLuceneRDD**

LuceneRDD + Spatial search

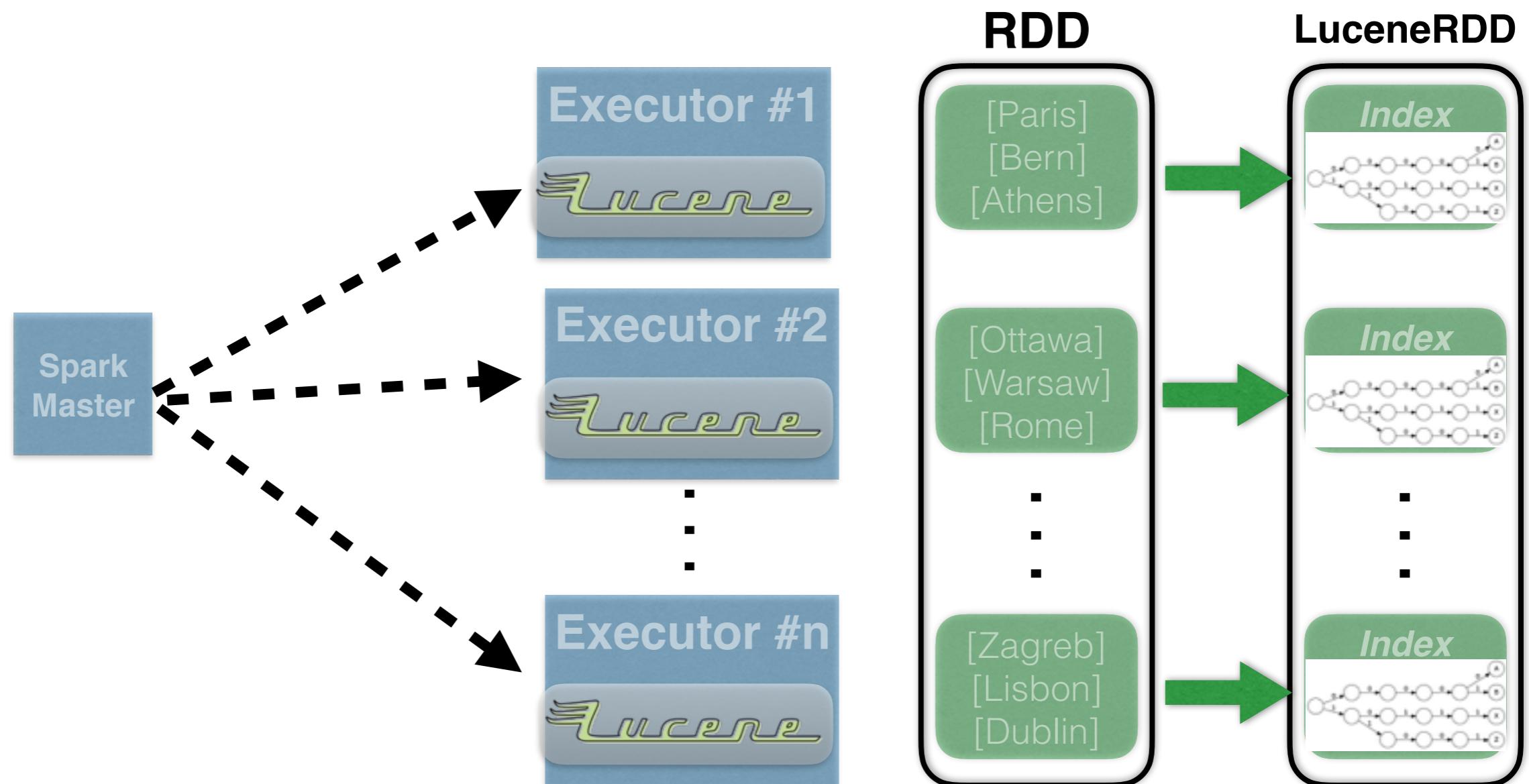
```
spark-shell --packages org.zouzias:spark-lucenerdd:0.2.2
```

or

```
libraryDependencies += "org.zouzias" %% "spark-lucenerdd" % "0.2.2"
```

- \* Main development: Spark **2.x** (supports Spark  $\geq 1.4$ )
- \* Lucene **5.5.3** (Lucene **6.2.2** JVM 8)
- \* Released on maven central & spark-packages (Scala 2.10 / 2.11)

# LuceneRDD

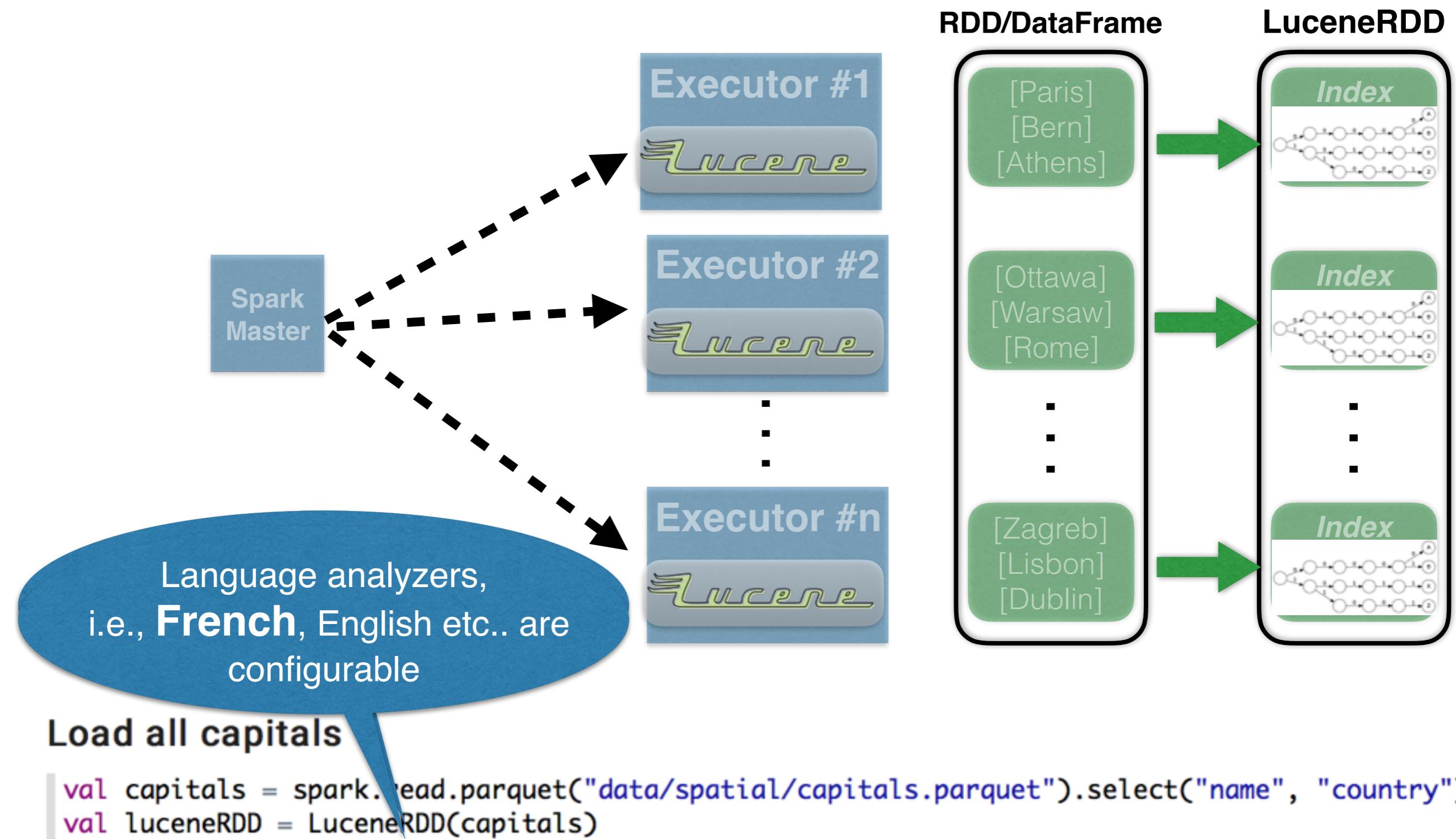


## Main idea

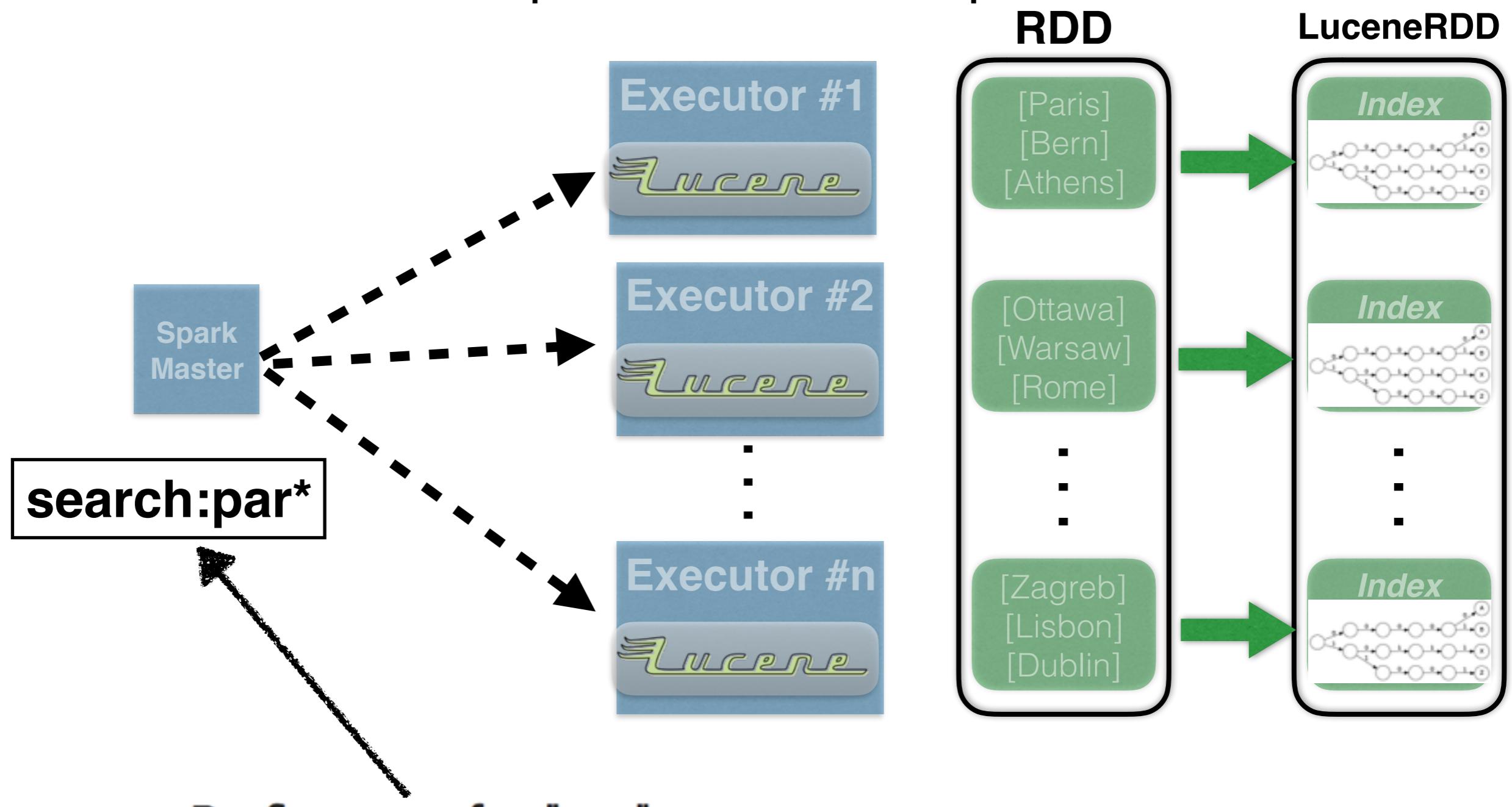
Index RDD partition data with distributed Lucene index  
**(inverted index per partition)**

Index storage  
**disk** or **memory**

# DataFrame to LuceneRDD\*



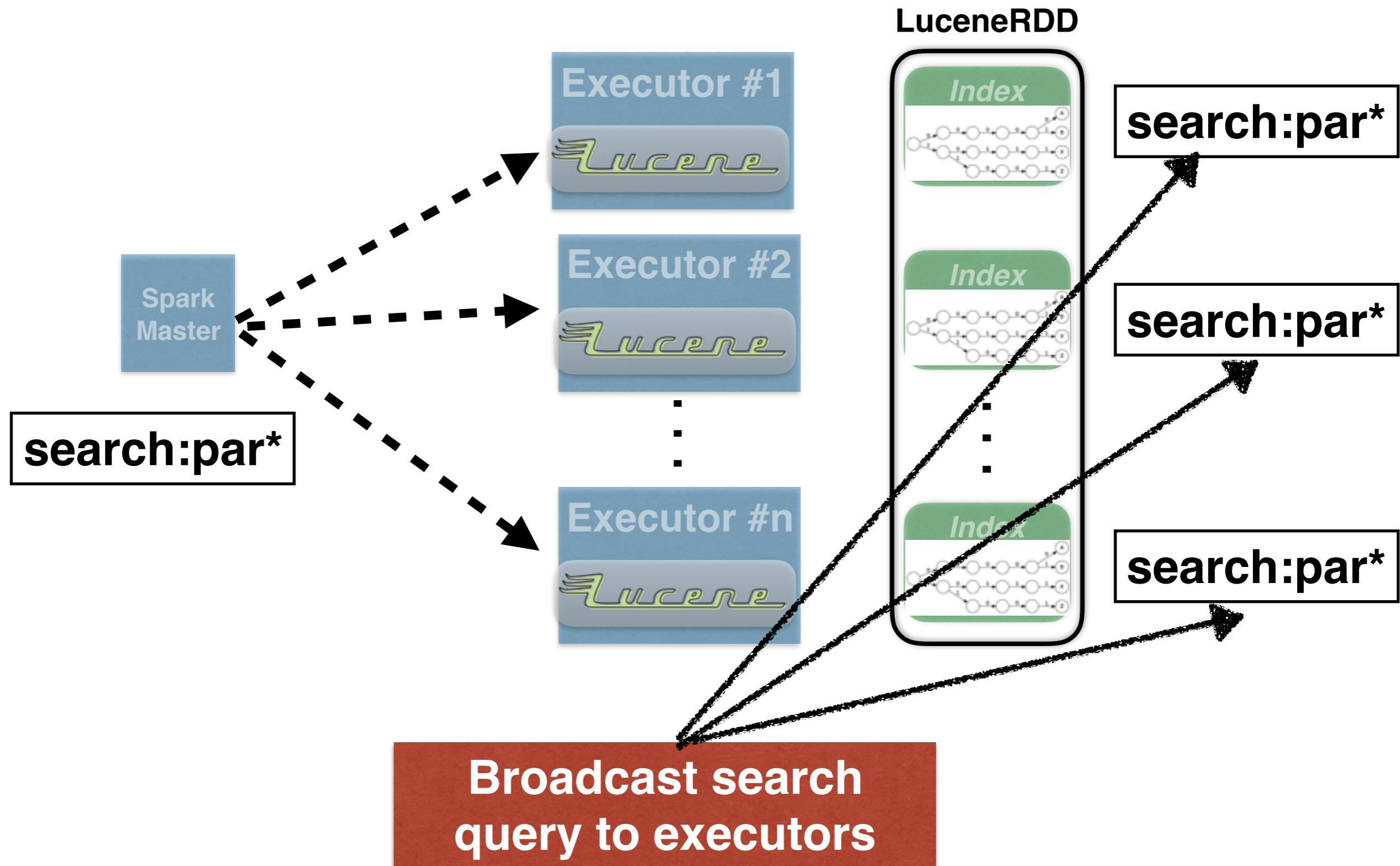
# Simplistic Example



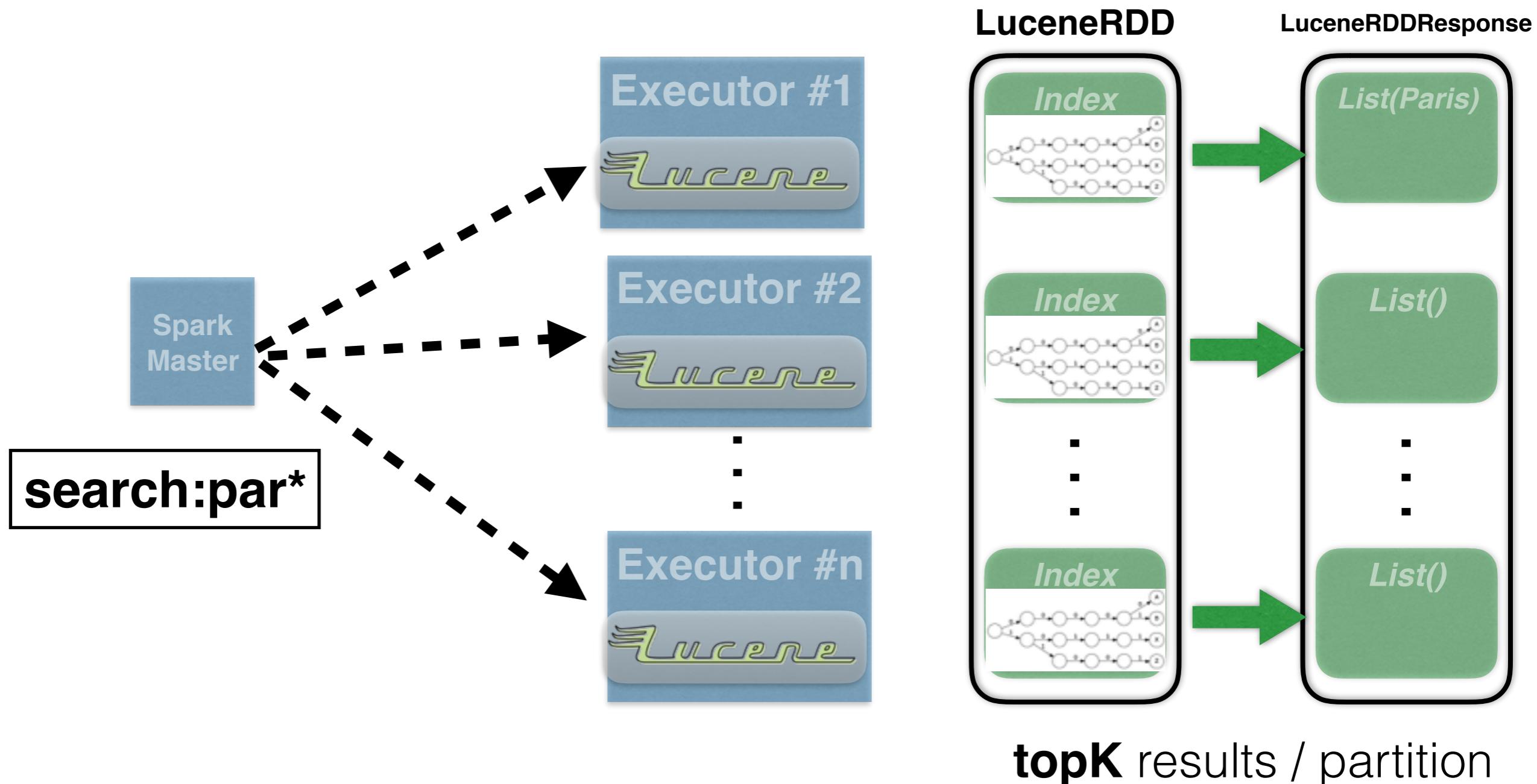
```

| val topK = 10
| val response = luceneRDD.prefixQuery("name", "par", topK)
  
```

# Simplistic Example



# Simplistic Example



**Lucene prefix Search  
on each partition**

# Aggregate Results

How to aggregate?  
TopK monoids in **action**

Spark  
Master

***response.take(k)***

List(Paris)

List()

List()

Executor #1

Executor #2

⋮

Executor #n

LuceneRDDResponse

List(Paris)

List()

List()

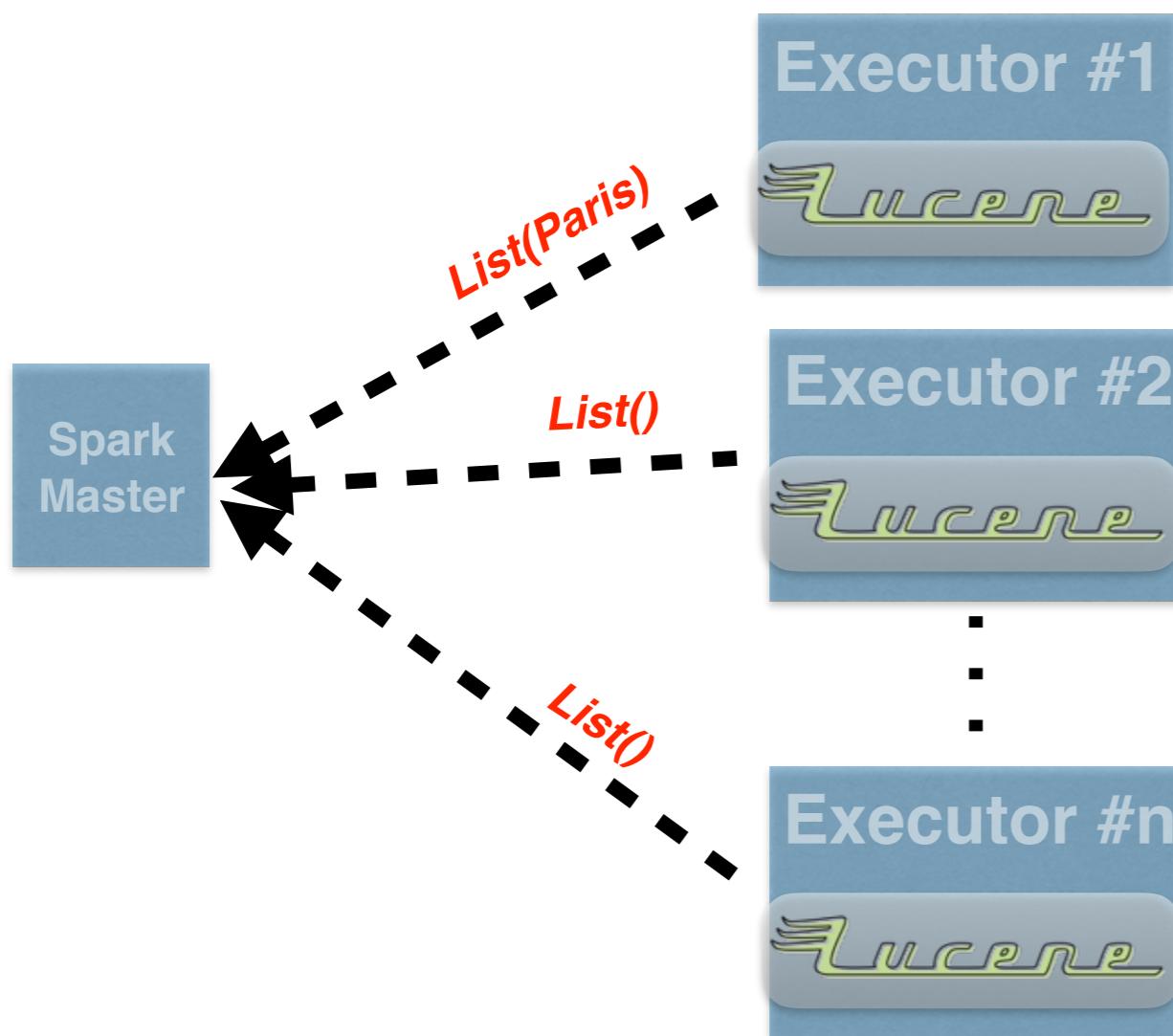
```
| response.take(k).foreach{println}
```

SparkScoreDoc(1.0,9,0,Text fields:country:[France]name:[Paris])

SparkScoreDoc(1.0,43,0,Text fields:country:[Suriname]name:[Paramaribo])

# Aggregate Results

**LuceneRDDResponse**



*TopKMonoid((0.8, Paris), (0.7, Result2), (0.6, Result3))*

+

*TopKMonoid((0.71, XXX), (0.35, YYY), (0.25, ZZZ))*

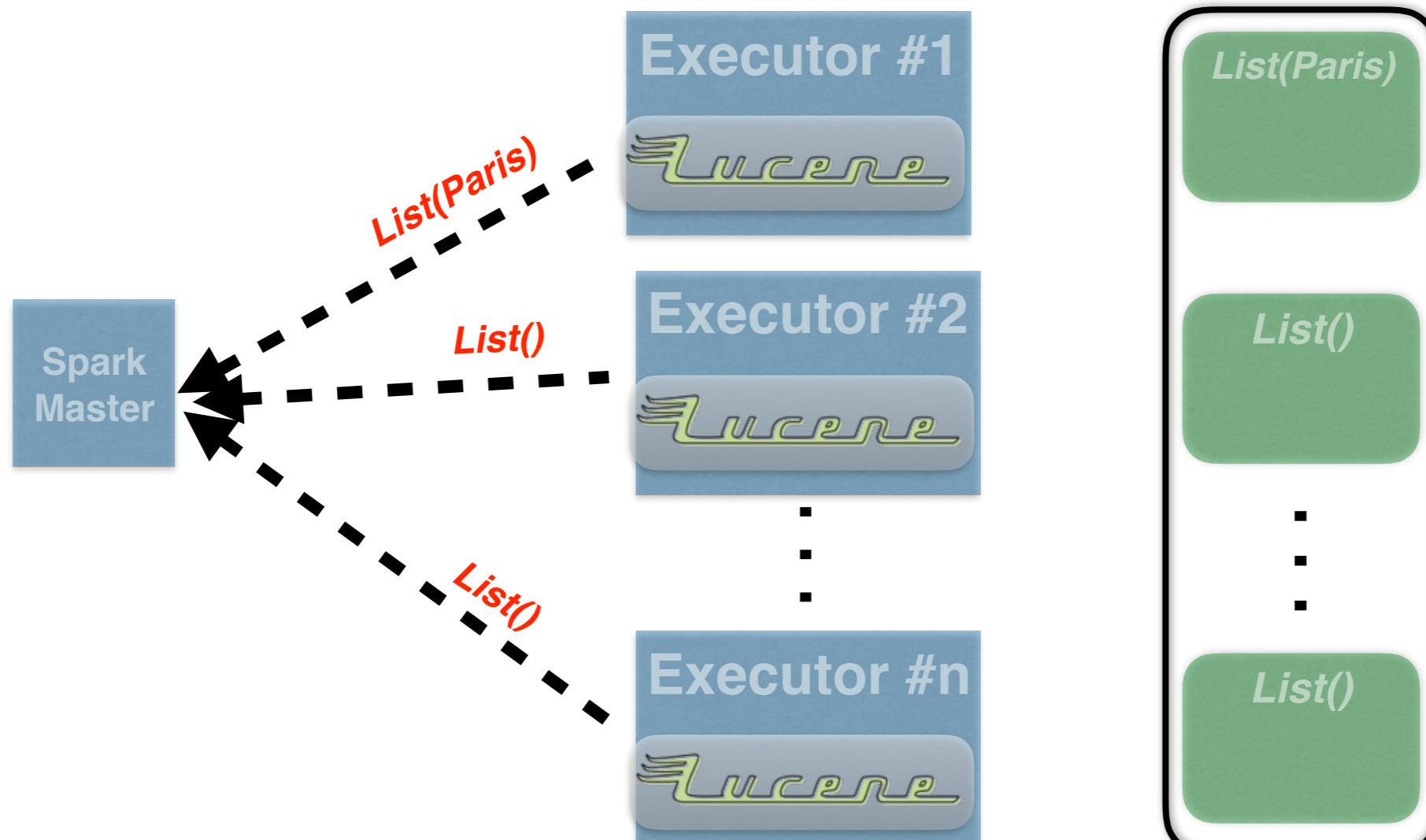
+

*TopKMonoid((0.12, XXX), (0.09, YYY), (0.05, ZZZ))*

***TopK Sorted  
Scored Results***

Twitter's Algebird TopKMonoid

## LuceneRDDResponse



## *LuceneRDDResponse.take(k)*

```
/*
 * Use [[TopKMonoid]] to take
 * @param num
 * @return
 */
override def take(num: Int): Array[SparkScoreDoc] = {
  val monoid = new TopKMonoid[SparkScoreDoc](num)(ordering)
  partitionsRDD.map(monoid.build(_))
    .reduce(monoid.plus).items.toArray
}
```

# LuceneRDD Operations

| Operation   | Syntax  | Description   |
|---|---|---|
| <b>Flexible multi-fields query:</b><br>term/Fuzzy/Prefix/etc sub-queries<br>+ boolean queries | <code>LuceneRDD.termQuery(field, query, topK)</code><br><code>LuceneRDD.fuzzyQuery(field, query, maxEdits, topK)</code> | Exact term search   |
| Phrase Query  | <code>LuceneRDD.phraseQuery(field, query, topK)</code>  | Phrase search   |
| Prefix Query  | <code>LuceneRDD.prefixSearch(field, prefix, topK)</code>  | Prefix search   |
| Query Parser  | <code>LuceneRDD.query(queryString, topK)</code>   | Query parser search   |
| Faceted Search  | <code>FacetedLuceneRDD.facetQuery(queryString, field, topK)</code>  | Faceted Search  |
| Record Linkage  | <code>LuceneRDD.link(otherEntity: RDD[T], linkageFct: T =&gt; searchQuery, topK)</code>                                 | Record linkage via Lucene queries   |
| Circle Search   | <code>ShapeLuceneRDD.circleSearch((x,y), radius, topK)</code>   | Search within radius  |
| Bbox Search   | <code>ShapeLuceneRDD.bboxSearch(lowerLeft, upperLeft, topK)</code>  | Bounding box  |
| Spatial Linkage   | <code>ShapeLuceneRDD.linkByRadius(RDD[T], linkage: T =&gt; (x,y), radius, topK)</code>                                  | <b>ShapeLuceneRDD</b><br>Spatial radius linkage<br>LuceneRDD + Spatial search |

# *Entity Linkage*

## *a.k.a. approximate join*

# Entity Linkage (approx. join)

## Left Dataset

|                                 |  |
|---------------------------------|--|
| James Stephenson                | One Ford Way, Dearborn, MI               |
| Ford Petaprod                   |  |
| Alicia Thomson                  | 4789 Woodward Avenue Detroit             |
| Detroitics                      | MegaProd,TeraProd                        |
| Jack Jones                      | One Microsoft Way Redmond                |
| Microsoft USA                   | PetaProd                                 |
| James Jones                     | 1234 Woodward Ave. Detroit "Home", MI    |
| Moonlighting                    | MegaProd                                 |
| James Jones                     | 4789 Woodward Ave. Detroit "Work", MI    |
| Detroitics                      | GigaProd,MegaProd                        |
| Al Shepard                      | New York, NY                             |
| Newyorkonics                    | SuperProd,MegaProd                       |
| Mary Barry                      | GM Renaissance Center, Detroit, MI 48243 |
| General Motors                  | TeraProd,SuperProd                       |
| Ram Kumar                       | New Orchard Road, Armonk, NY             |
| International Business Machines | PetaProduct                              |
| Joel Smith                      | New York, NY                             |
| New York-onics                  | MegaProduct,SuperProd                    |
| Mike Taylor                     | Unknown                                  |
| Lockheed-M                      | Tera Product                             |

?

## Right Dataset

|                        |  |
|------------------------|--|
| James (Jim) Stephenson | One Ford Way, Dearborn, MI               |
| 48126                  |  |
| Ford Petaprod          |  |
| Alice Thompson         | 4789 Woodward Ave, Detroit, MI           |
| Detroitics             | MegaProduct,TeraProd                     |
| Jackob Jones           | One Microsoft Way Redmond, WA 98052-7329 |
| Microsoft              | PetaProduct                              |
| Jim Jones              | 1234 Woodward Ave. Detroit, MI           |
| Moonlighting, Inc      | MegaProd                                 |
| James Jones            | 4789 Woodward Ave. Detroit, MI           |
| Detroitics             | GigaProd                                 |
| Joe Smith              | 1234 56 St, New York, NY                 |
| Newyorkonics           | SuperProd,MegaProd                       |
| Mary Barry             | Renaissance Center, Detroit              |
| GM                     | Tera product                             |
| Ramkrishna (Ram) Kumar | New Orchard Road, Armonk, New York 10504 |
| IBM                    | GigaProduct, PetaProduct                 |
| Alicia Shepard         | 1234 56 St, New York, NY                 |
| Newyorkonics           | ???                                      |
| Michael Taylor         | Classified, Bethesda, MD                 |
| Lockheed Martin        | TeraProd                                 |

Join left & right dataset  
according to a “similarity”  
(NO ids available)

Simplicity:  
one-to-one linkage

# Entity Linkage (approx. join)

## Left Dataset

|                                 |  |
|---------------------------------|--|
| James Stephenson                | One Ford Way, Dearborn, MI               |
| Ford Petaprod                   |  |
| Alicia Thomson                  | 4789 Woodward Avenue Detroit             |
| Detroitics                      | MegaProd,TeraProd                        |
| Jack Jones                      | One Microsoft Way Redmond                |
| Microsoft USA                   | PetaProd                                 |
| James Jones                     | 1234 Woodward Ave. Detroit "Home", MI    |
| Moonlighting                    | MegaProd                                 |
| James Jones                     | 4789 Woodward Ave. Detroit "Work", MI    |
| Detroitics                      | GigaProd,MegaProd                        |
| Al Shepard                      | New York, NY                             |
| Newyorkonics                    | SuperProd,MegaProd                       |
| Mary Barry                      | GM Renaissance Center, Detroit, MI 48243 |
| General Motors                  | TeraProd,SuperProd                       |
| Ram Kumar                       | New Orchard Road, Armonk, NY             |
| International Business Machines | PetaProduct                              |
| Joel Smith                      | New York, NY                             |
| New Yorkonics                   | MegaProduct,SuperProd                    |
| Mike Taylor                     | Unknown                                  |
| Lockheed-M                      | TeraProduct                              |

?

## Right Dataset

|                        |  |
|------------------------|--|
| James (Jim) Stephenson | One Ford Way, Dearborn, MI               |
| 48126                  |  |
| Ford Petaprod          |  |
| Alice Thompson         | 4789 Woodward Ave, Detroit, MI           |
| Detroitics             | MegaProduct,TeraProd                     |
| Jackob Jones           | One Microsoft Way Redmond, WA 98052-7329 |
| Microsoft              | PetaProduct                              |
| Jim Jones              | 1234 Woodward Ave. Detroit, MI           |
| Moonlighting, Inc      | MegaProd                                 |
| James Jones            | 4789 Woodward Ave. Detroit, MI           |
| Detroitics             | GigaProd                                 |
| Joe Smith              | 1234 56 St, New York, NY                 |
| Newyorkonics           | SuperProd,MegaProd                       |
| Mary Barry             | Renaissance Center, Detroit              |
| GM                     | TeraProduct                              |
| Ramkrishna (Ram) Kumar | New Orchard Road, Armonk, New York 10504 |
| IBM                    | GigaProduct, PetaProduct                 |
| Alicia Shepard         | 1234 56 St, New York, NY                 |
| Newyorkonics           | ????                                     |
| Michael Taylor         | Classified, Bethesda, MD                 |
| Lockheed Martin        | TeraProd                                 |

### Naive approach

Compute pairwise scores does not scale:  $O(n^2)$

# Entity Linkage API

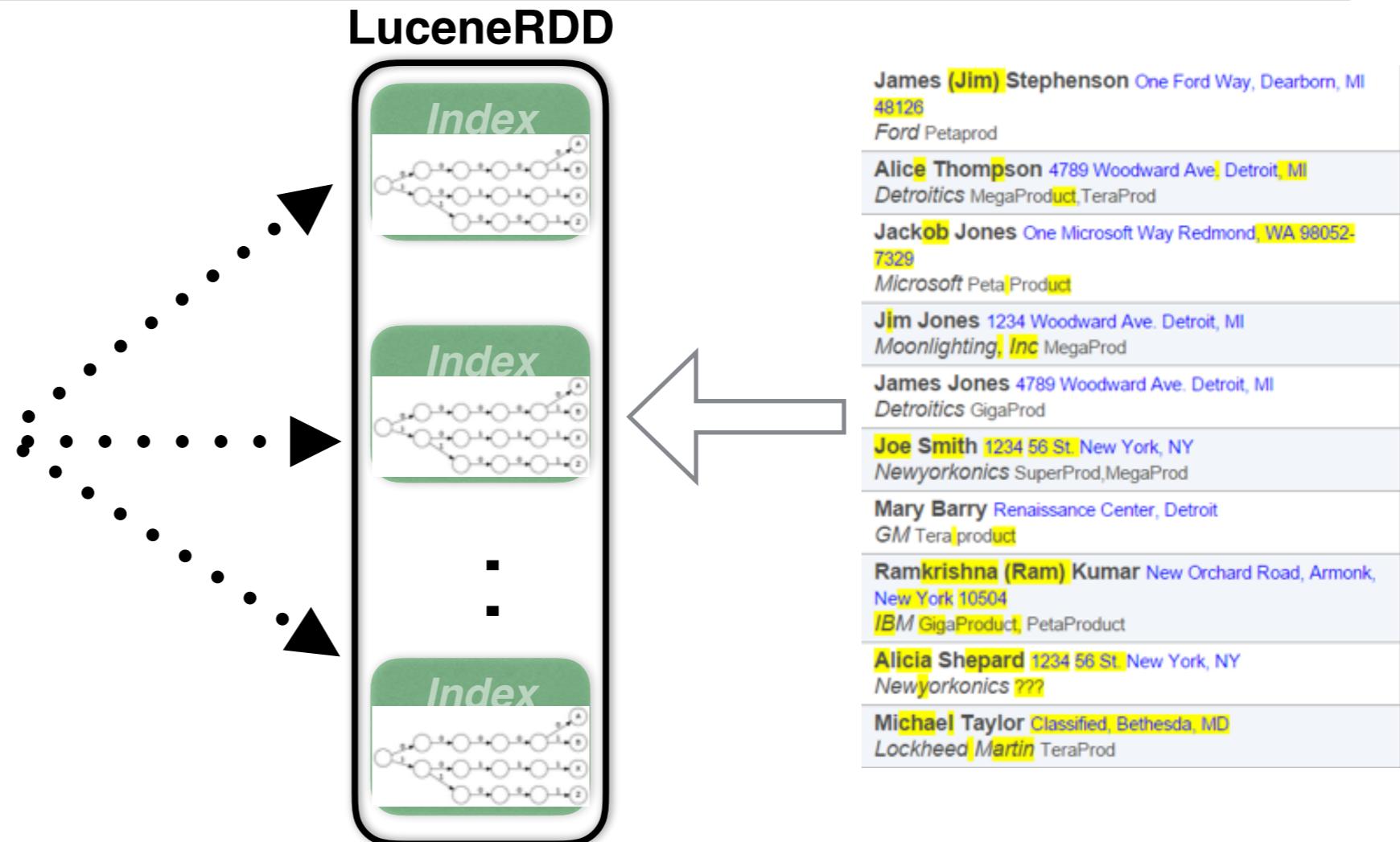
| Operation       | Syntax  | Description                       |
|-----------------|---|-----------------------------------|
| Term Query      | <code>LuceneRDD.termQuery(field, query, topK)</code>                                    | Exact term search                 |
| Fuzzy Query     | <code>LuceneRDD.fuzzyQuery(field, query, maxEdits, topK)</code>                         | Fuzzy term search                 |
| Phrase Query    | <code>LuceneRDD.phraseQuery(field, query, topK)</code>                                  | Phrase search                     |
| Prefix Query    | <code>LuceneRDD.prefixSearch(field, prefix, topK)</code>                                | Prefix search                     |
| Query Parser    | <code>LuceneRDD.query(queryString, topK)</code>   | Query parser search               |
| Faceted Search  | <code>FacetedLuceneRDD.facetQuery(queryString, field, topK)</code>                      | Faceted Search                    |
| Record Linkage  | <code>LuceneRDD.link(otherEntity: RDD[T], linkageFct: T =&gt; searchQuery, topK)</code> | Record linkage via Lucene queries |
| Circle Search   | <code>ShapeLuceneRDD.circleSearch((x,y), radius, topK)</code>                           | Search within radius              |
| Bbox Search     | <code>ShapeLuceneRDD.bboxSearch(lowerLeft, upperLeft, topK)</code>                      | Bounding box                      |
| Spatial Linkage | <code>ShapeLuceneRDD.linkByRadius(RDD[T], linkage: T =&gt; (x,y), radius, topK)</code>  | Spatial radius linkage            |

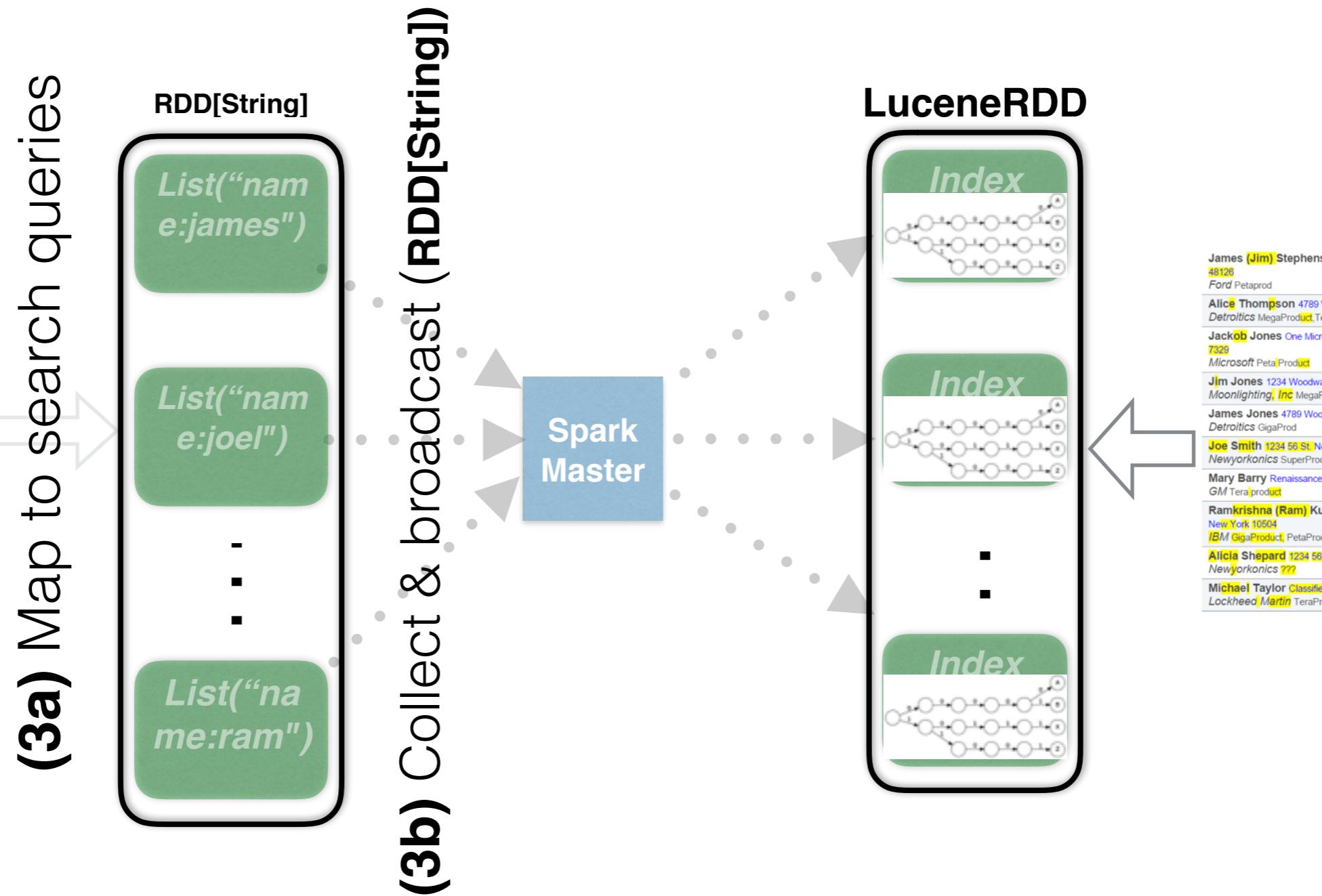
## LuceneRDD Approximate Join

- 1) Index **right** dataset using **LuceneRDD**
- 2) Define “search similarity function” between **each left entity & right** dataset
- 3) For every **left** entity, **search** & **zip** using **topK** results

|                |   |                                   |
|----------------|---|-----------------------------------|
| Record Linkage | <code>LuceneRDD.link(otherEntity: RDD[T], linkageFct: T =&gt; searchQuery, topK)</code> | Record linkage via Lucene queries |
|----------------|---|-----------------------------------|

|  |
|--|
| James Stephenson One Ford Way, Dearborn, MI<br>Ford Petaprod                             |
| Alicia Thomson 4789 Woodward Avenue Detroit<br>Detroitics MegaProd,TeraProd              |
| Jack Jones One Microsoft Way Redmond<br>Microsoft USA PetaProd                           |
| James Jones 1234 Woodward Ave. Detroit "Home", MI<br>Moonlighting Mega Prod              |
| James Jones 4789 Woodward Ave. Detroit "Work", MI<br>Detroitics GigaProd,MegaProd        |
| AI Shepard New York, NY<br>Newyorkonics SuperProd,MegaProd                               |
| Mary Barry GM Renaissance Center, Detroit, MI 48243<br>General Motors TeraProd,SuperProd |
| Ram Kumar New Orchard Road, Armonk, NY<br>International Business Machines PetaProduct    |
| Joel Smith New York, NY<br>New York-onics MegaProduct,SuperProd                          |
| Mike Taylor Unknown<br>Lockheed-M Tera Product   |





## LuceneRDD Entity Linkage

- 1) Index **right** dataset using **LuceneRDD**
- 2) Define “search similarity function” between **left** & **right** dataset
- 3) For every left entity, **search** & **zip** by search similarly

# partitionsRDD: linkage code

Query each **partition** and **zip** query with **topK** results

**Map** each **left** entity to search query & **broadcast** to each partition

```
def link[T1: ClassTag](other: RDD[T1], searchQueryGen: T1 => String, topK: Int = DefaultTopK): RDD[(T1, List[SparkScoreDoc])] = {
  val monoid = new TopKMonoid[SparkScoreDoc](topK)(SparkScoreDoc.descending)

  val queries = other.map(searchQueryGen).collect()
  val queriesB = partitionsRDD.context.broadcast(queries)

  val resultsByPart: RDD[(Long, TopK[SparkScoreDoc])] = partitionsRDD.flatMap { partition =>
    queriesB.value.zipWithIndex.map { case (qr, index) =>
      val results = partition.query(qr, topK).map(x => monoid.build(x))
      (index.toLong, results.reduceOption(monoid.plus).getOrElse(monoid.zero))
    }
  }

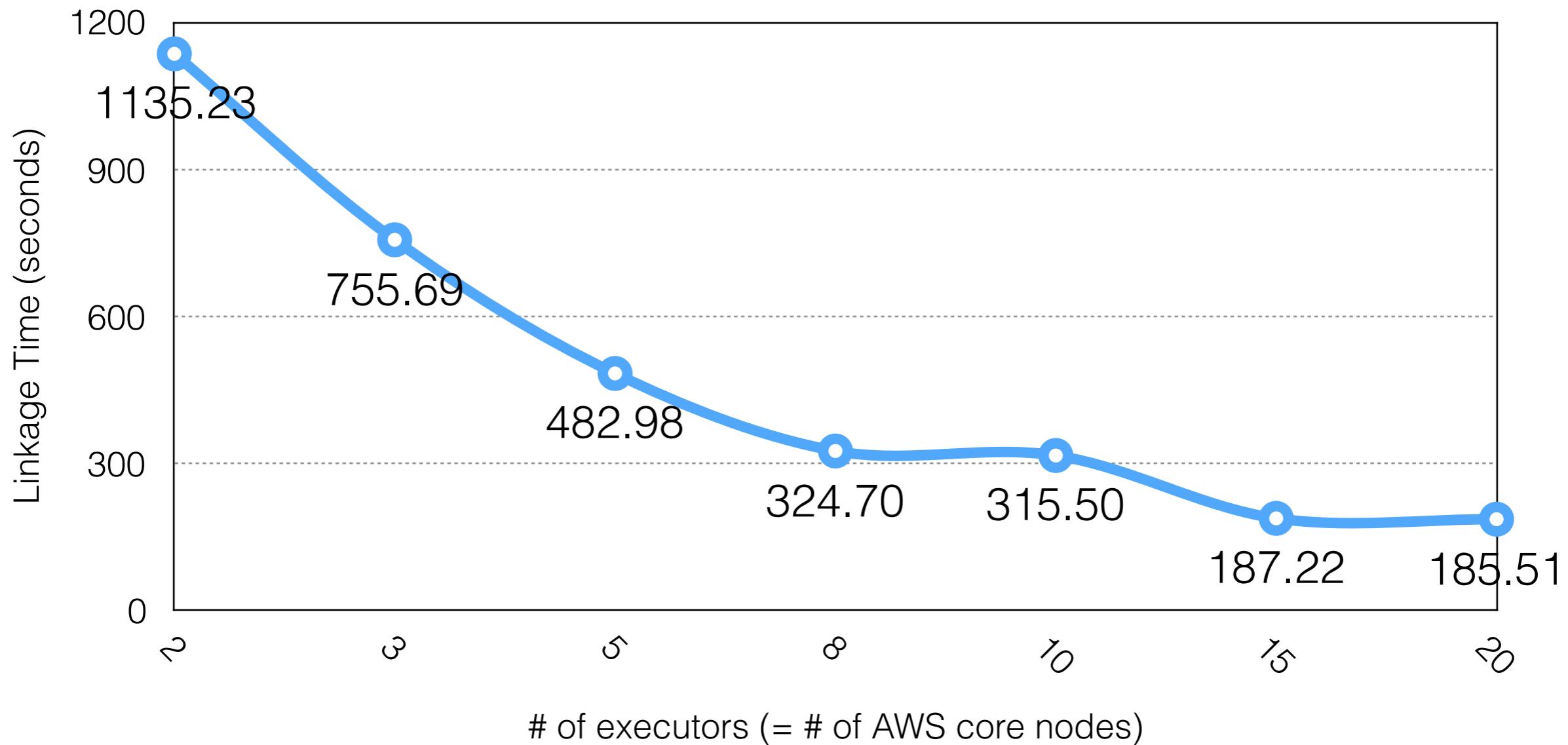
  val results = resultsByPart.reduceByKey(monoid.plus)
  other.zipWithIndex.map(_.swap).join(results).values
    .map(joined => (joined._1, joined._2.items.take(topK)))
}
```

Reduce **topK** monoids per query

*Does LuceneRDD approx. join scale?*

<https://github.com/zouzias/spark-lucenerdd-aws>

## LuceneRDD Entity Linkage Scalability (EMR 4.8.0 / Spark 1.6.2)



## Joined Datasets

- \* H1B US Visa Applications (2.6 million records)
- \* Geonames US Cities (4.4 million records)

AWS: master: m3.xlarge, Core: r3.xlarge, Exec. mem: 9G

*Demo*

*LuceneRDD Linkage API*

*(ACM vs DBLP research articles)*

Köpcke, H.; Thor, A.; Rahm, E.

**Evaluation of entity resolution approaches on real-world match problems**

Proc. 36th Intl. Conference on Very Large Databases (VLDB), 2010

<https://github.com/zouzias/spark-lucenerdd-examples>

# Spatial/Entity Linkage API

| Operation       | Syntax  | Description                       |
|-----------------|---|-----------------------------------|
| Term Query      | <code>LuceneRDD.termQuery(field, query, topK)</code>                                    | Exact term search                 |
| Fuzzy Query     | <code>LuceneRDD.fuzzyQuery(field, query, maxEdits, topK)</code>                         | Fuzzy term search                 |
| Phrase Query    | <code>LuceneRDD.phraseQuery(field, query, topK)</code>                                  | Phrase search                     |
| Prefix Query    | <code>LuceneRDD.prefixSearch(field, prefix, topK)</code>                                | Prefix search                     |
| Query Parser    | <code>LuceneRDD.query(queryString, topK)</code>   | Query parser search               |
| Faceted Search  | <code>FacetedLuceneRDD.facetQuery(queryString, field, topK)</code>                      | Faceted Search                    |
| Record Linkage  | <code>LuceneRDD.link(otherEntity: RDD[T], linkageFct: T =&gt; searchQuery, topK)</code> | Record linkage via Lucene queries |
| Circle Search   | <code>ShapeLuceneRDD.circleSearch((x,y), radius, topK)</code>                           | Search within radius              |
| Bbox Search     | <code>ShapeLuceneRDD.bboxSearch(lowerLeft, upperLeft, topK)</code>                      | Bounding box                      |
| Spatial Linkage | <code>ShapeLuceneRDD.linkByRadius(RDD[T], linkage: T =&gt; (x,y), radius, topK)</code>  | Spatial radius linkage            |

*Demo*

*ShapeLuceneRDD Linkage*

*(Country polygons vs Capital points)*

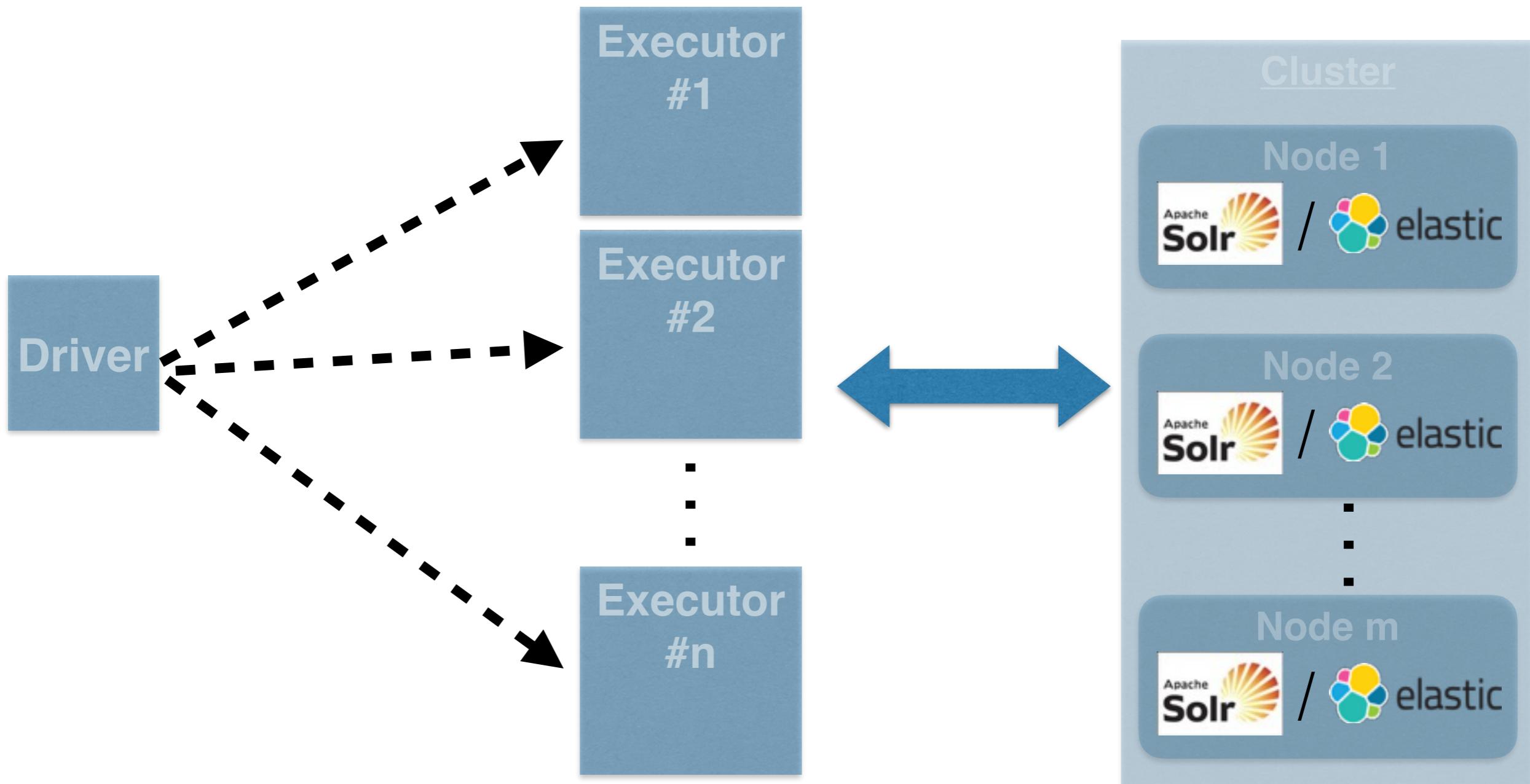
<https://github.com/zouzias/spark-lucenerdd-examples>

# Future Work

- \* Contributors / use cases are welcome!
- \* More performance tests with approx. join
- \* Spatial-linkage performance tests using open Street Map data
- \* Bypass the collect-broadcast Spark pattern?

*Merci pour votre attention!*  
Questions?

# Elastic/SolrCloud Connector



## Limitations

- **m** is **fixed** & (usually) small compared to **n**
- additional network communication: execs & elastic/solr