



Lucidworks

# High-level Architecture

X

Solr

Dark SQL

Y



committer since April 2014, work for

Cloud features ... and bin/solr!

for Lucene / Solr 5.1

ction

nce working with Hadoop, Pig, Hive,  
started using Spark about 6 months

[k.apache.org/](http://k.apache.org/)

[www.cs.berkeley.edu/~matei/papers/2012/nsdi\\_spark](http://www.cs.berkeley.edu/~matei/papers/2012/nsdi_spark)

ed alternative to MapReduce

*sorted 100TB in 23 minutes (3x faster than Yahoo's previous)*

y Data

*algorithms (PageRank, K-Means, Logistic regression) & inter*

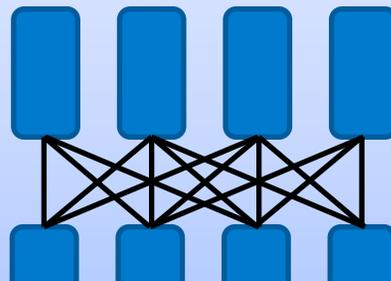
Spark  
Streaming

MLlib  
(machine  
learning)

Gr  
(E

## Spark Core

The Shuffle



Cachin



# Spark Master (daemon)

Keeps track of live workers

Web UI on port 8080

Task Scheduler

Restart failed tasks

## Spark Worker

Spark

## Spark Executors

Tasks

*losing a master prevents new applications from being executed*

*can achieve HA using ZooKeeper and multiple master nodes*

count  
sample  
reduce  
sortByKey  
distinct  
groupBy  
union  
take  
mapPartitions  
cogroup

```
val file = spark.textFile(  
val counts = file.flatMap(  
    flatMap
```

## Split lines into words

quick

brown

fox

quick

⋮

quick

⋮

## Map words into pairs with count of 1

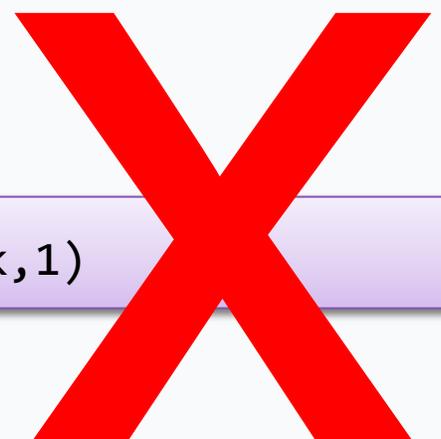
(quick,1)

(brown,1)

(fox,1)

(quick,1)

(quick,1)



system OR using a transformation of an

e of coarse-grained transformations (ma

DDs can be re-computed by re-playing th

ersist an RDD (for reusing during interact

ons

ETL jobs

job

matching against stored queries

tions (interactive data mining, machi

l query as Spark RDD (resilient distributed datas

ults from each shard in parallel

sequence of RDDs

you can use all the other Spark libs (MLlib, etc)

es

must be less than the batch interval time

:

by, join, etc)

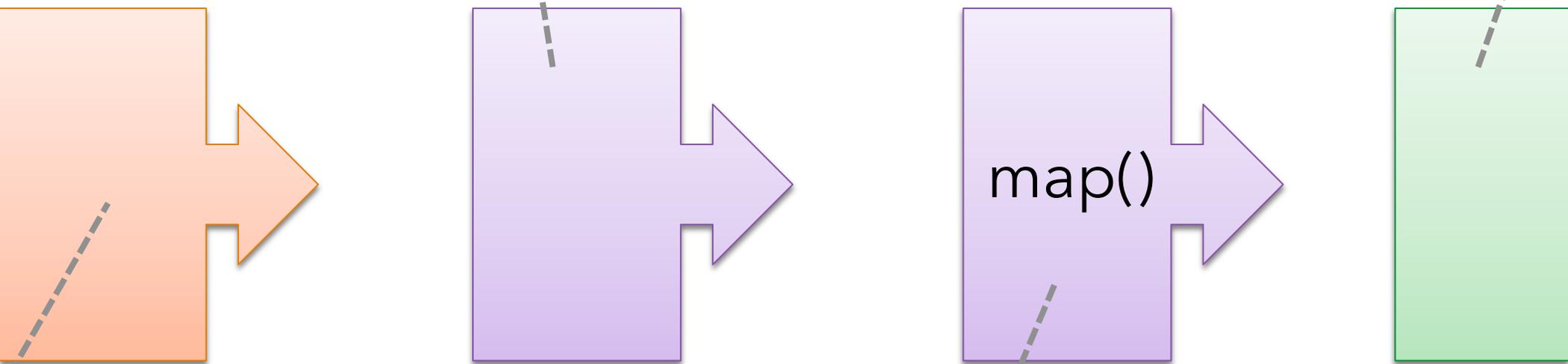
external sink, e.g. Solr)

e!

st localhost:2181 -collection social

Various transformations / enrichments  
on each tweet (e.g. sentiment analysis,  
language detection)

SolrSupport.index



is> **tweets** =

(null, filters):

javaDStream<SolrInputDocument> **docs** = **tweets**

```
status> tweets =  
m(jssc, null, filters);  
  
to SolrInputDocument objects for indexing in Solr  
ent> docs = tweets.map(  
  
InputDocument>() {  
ent call(Status status) {  
=  
ToSolrInputDoc("tweet-"+status.getId(), statu  
er_s", "twitter");  
_s", status.getUser().getScreenName());  
, status.isRetweet() ? "echo" : "post");
```

```
utDocument>, Void>() {
solrInputDocument> solrInputDocumentJavaRDD) throws E
.foreachPartition(
Iterator<SolrInputDocument>>() {
Iterator<SolrInputDocument> solrInputDocumentIterator)
solrServer = getSolrServer(zkHost);
ment> batch = new ArrayList<SolrInputDocument>();
cumentIterator.hasNext()) {
putDocumentIterator.next());
    >= batchSize)
r(solrServer, collection, batch);

())
solrServer, collection, batch);
```

determine which of a large set of stor

ernative flow paths through a stream,

into an embedded (in-memory) Solr

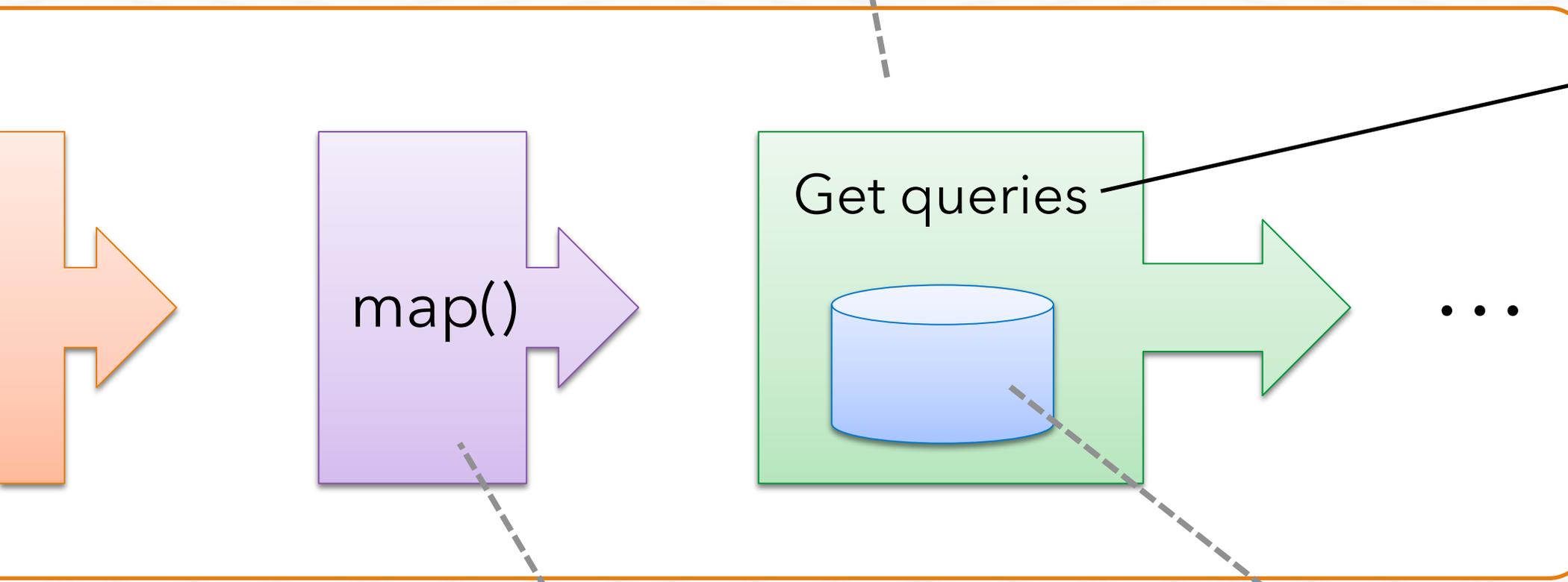
ch queries match

k; you have to decide where to load t

at to do when matches are found

... need to scale to many queries, ch

JavaDStream<SolrInputDocument> enriched =  
**SolrSupport.filterDocuments**(docFilterCon



Index docs into an  
EmbeddedSolrServ  
Initialized from conf  
stored in ZooKeepe

m<SolrInputDocument> docs = tweets.map(  
ction<Status SolrInputDocument>()

o each shard leader using metadata from  
nment, and ConcurrentUpdateSolrClient

```
Partitioner = new ShardPartitioner(zkHost, col  
itioner).foreachPartition(  
Tuple2<String, SolrInputDocument>>>() {  
r<Tuple2<String, SolrInputDocument>> tupleIter  
ient cuss = null;  
xt()) {  
oncurrentUpdateSolrClient once per partition
```

query and expose as an RDD

```
JavaRDD<SolrDocument>
```

```
needed (cursorMark)
```

It sets where global sort order doesn't  
execution by distributing requests across

```
nt> results =
```

struct RDD<Vector> which can then be

```
solrRDD(zkHost, collection);
```

```
rs =
```

```
solrVectors(jsc, solrQuery, field, numFeat
```

```
=
```

```
solrRDD().rdd(), numClusters, numIterations)
```

```
solrQuery(...);  
_t", "type_s");
```

```
solrRDD(zkHost, collection);  
solrJavaRDD = solrRDD.queryShards(jsc, solrQuery);  
new SQLContext(jsc);
```

```
Context, solrQuery, solrJavaRDD, zkHost, collection);
```

```
COUNT(type_s) FROM tweets WHERE type_s='echo'
```

```
s.javaRDD().map(new Function<Row, Long>() {
```

ntation of Solr and Spark on YARN

for reads and writes to Solr

- improving replication performance

to Solr Scale Toolkit

give visibility into performance

to me with questions:



