

CS 744: BIG DATA SYSTEMS

Shivaram Venkataraman

Fall 2018

ADMINISTRIVIA

- Assignment 1: Due Oct 1
- Project
 - Ideas posted on Piazza
 - Meetings on Oct 4 / Oct 5

DISCUSSION: PROGRAMMABILITY

Most real applications require multiple MR steps

- Google indexing pipeline: 21 steps
- Analytics queries (e.g. sessions, top K): 2-5 steps
- Iterative algorithms (e.g. PageRank): 10's of steps

Multi-step jobs create spaghetti code

- 21 MR steps → 21 mapper and reducer classes

DISCUSSION: PERFORMANCE

MR only provides one pass of computation

- Must write out data to file system in-between

Expensive for apps that need to *reuse* data

- Multi-step algorithms (e.g. PageRank)
- Interactive data mining

SPARK

Programmability: clean, functional API

- Parallel transformations on collections
- 5-10x less code than MR
- Available in Scala, Java, Python and R

Performance

- In-memory computing primitives
- Optimization across operators



SPARK PROGRAMMABILITY

Google MapReduce WordCount:

```
• #include "mapreduce/mapreduce.h"
•
• // User's map function
• class Splitwords: public Mapper {
•     public:
•     virtual void Map(const MapInput& input)
•     {
•         const string& text = input.value();
•         const int n = text.size();
•         for (int i = 0; i < n; ) {
•             // Skip past leading whitespace
•             while (i < n && isspace(text[i]))
•                 i++;
•             // Find word end
•             int start = i;
•             while (i < n && !isspace(text[i]))
•                 i++;
•             if (start < i)
•                 Emit(text.substr(
•                     start,i-start),"1");
•         }
•     }
• };
•
• REGISTER_MAPPER(Splitwords);
•
• // User's reduce function
• class Sum: public Reducer {
•     public:
•     virtual void Reduce(ReduceInput* input)
•     {
•         // Iterate over all entries with the
•         // same key and add the values
•         int64 value = 0;
•         while (!input->done()) {
•             value += StringToInt(
•                 input->value());
•             input->NextValue();
•         }
•         // Emit sum for input->key()
•         Emit(IntToString(value));
•     }
• };
•
• REGISTER_REDUCER(Sum);
•
• int main(int argc, char** argv) {
•     ParseCommandLineFlags(argc, argv);
•     MapReduceSpecification spec;
•     for (int i = 1; i < argc; i++) {
•         MapReduceInput* in= spec.add_input();
•         in->set_format("text");
•         in->set_filepattern(argv[i]);
•         in->set_mapper_class("Splitwords");
•     }
•
•     // Specify the output files
•     MapReduceOutput* out = spec.output();
•     out->set_filebase("/gfs/test/freq");
•     out->set_num_tasks(100);
•     out->set_format("text");
•     out->set_reducer_class("Sum");
•
•     // Do partial sums within map
•     out->set_combiner_class("Sum");
•
•     // Tuning parameters
•     spec.set_machines(2000);
•     spec.set_map_megabytes(100);
•     spec.set_reduce_megabytes(100);
•
•     // Now run it
•     MapReduceResult result;
•     if (!MapReduce(spec, &result)) abort();
• }
```

SPARK PROGRAMMABILITY

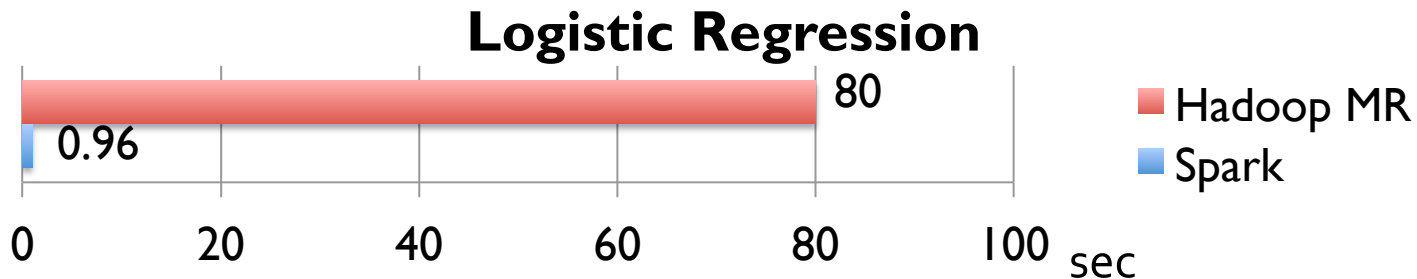
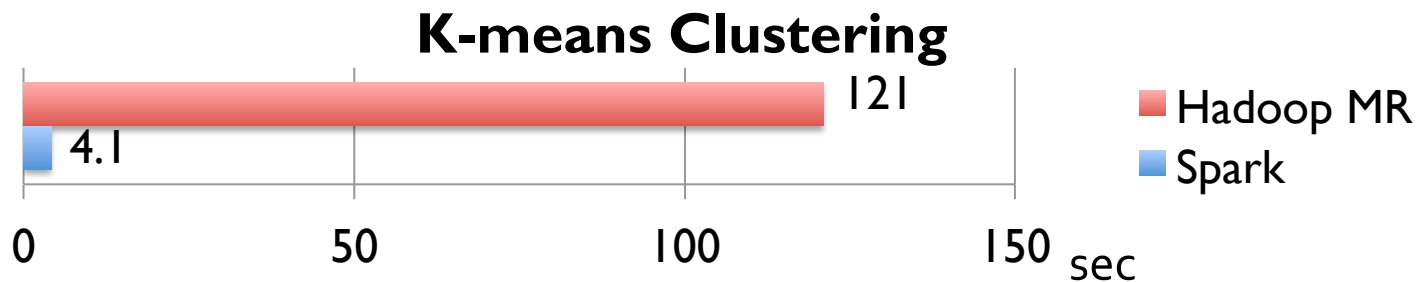
Spark WordCount:

```
val file = spark.textFile("hdfs://...")
val counts = file.flatMap(line => line.split(" "))
                  .map(word => (word, 1))
                  .reduceByKey(_ + _)

counts.save("out.txt")
```

SPARK PERFORMANCE

Iterative algorithms:



SPARK CONCEPTS

Resilient distributed datasets (RDDs)

- Immutable, partitioned collections of objects
- May be cached in memory for fast reuse

Operations on RDDs

- *Transformations* (build RDDs)
- *Actions* (compute results)

Restricted shared variables

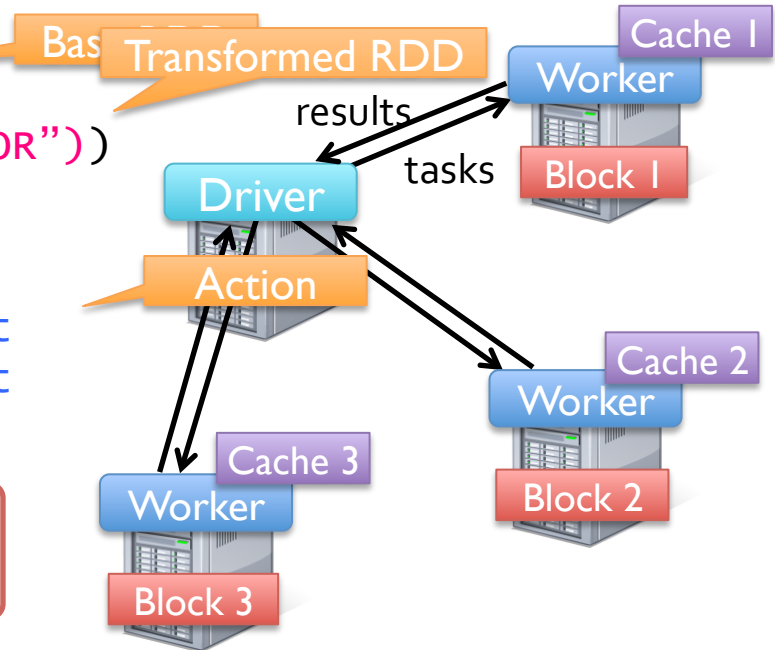
- Broadcast, accumulators

EXAMPLE: LOG MINING

Find error messages present in log files interactively
(Example: HTTP server logs)

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
messages.cache()
messages.filter(_.contains("foo")).count
messages.filter(_.contains("bar")).count
. . .
```

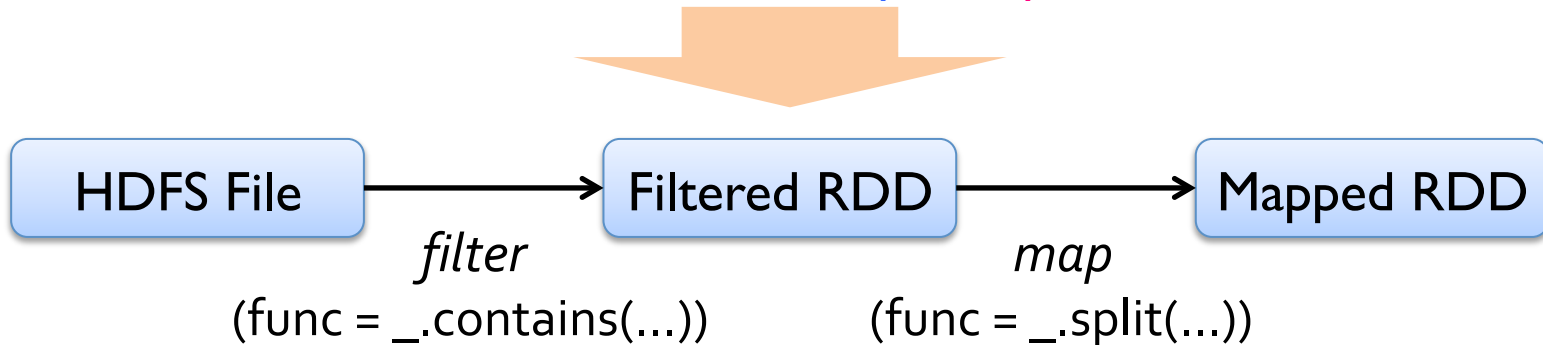
Result: search 1 TB data in 5-7 sec
(vs 170 sec for on-disk data)



FAULT RECOVERY

RDDs track *lineage* information that can be used to efficiently reconstruct lost partitions

Ex: `messages = textFile(...).filter(_.startsWith("ERROR")).map(_.split('\t')(2))`



SHARED VARIABLES

RDD operations: use local variables from scope

Two other kinds of shared variables:

- Broadcast Variables

- Accumulators

BROADCAST VARIABLES

```
val data = spark.textFile(...).map(readPoint).cache()

// Random Projection
val M = Matrix.random(N)

var w = Vector.random(D)

for (i <- 1 to ITERATIONS) {
  val gradient = data.map(p =>
    (1 / (1 + exp(-p.y*(w.dot(p.x.dot(M)))))) - 1)
    * p.y * p.x
  ).reduce(_ + _)
  w -= gradient
}

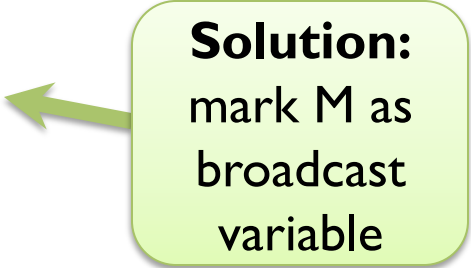
println("Final w: " + w)
```

Large Matrix

Problem:
M re-sent to all
nodes in each
iteration

BROADCAST VARIABLES

```
val data = spark.textFile(...).map(readPoint).cache()
// Random Projection
val M = spark.broadcast(Matrix.random(N))
var w = Vector.random(D)
for (i <- 1 to ITERATIONS) {
  val gradient = data.map(p =>
    (1 / (1 + exp(-p.y*(w.dot(p.x.dot(M.value)))) - 1) * p.y *
    p.x
  ).reduce(_ + _)
  w -= gradient
}
println("Final w: " + w)
```



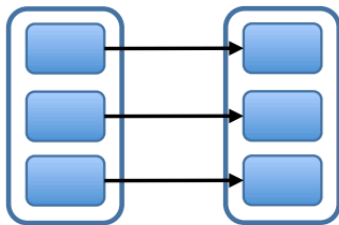
Solution:
mark M as
broadcast
variable

OTHER RDD OPERATIONS

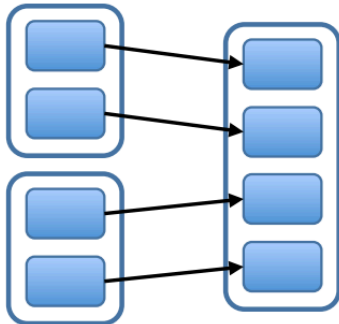
Transformations (define a new RDD)	map filter sample groupByKey reduceByKey cogroup	flatMap union join cross mapValues ...
Actions (output a result)	collect reduce take fold	count saveAsTextFile saveAsHadoopFile ...

DEPENDENCIES

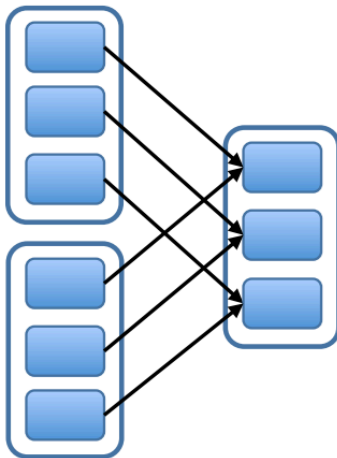
Narrow Dependencies:



map, filter

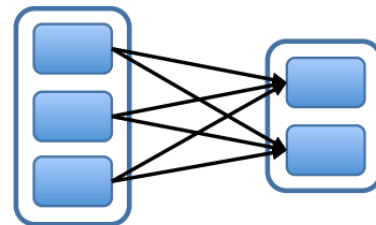


union

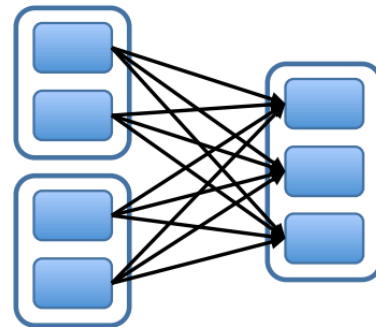


join with inputs
co-partitioned

Wide Dependencies:



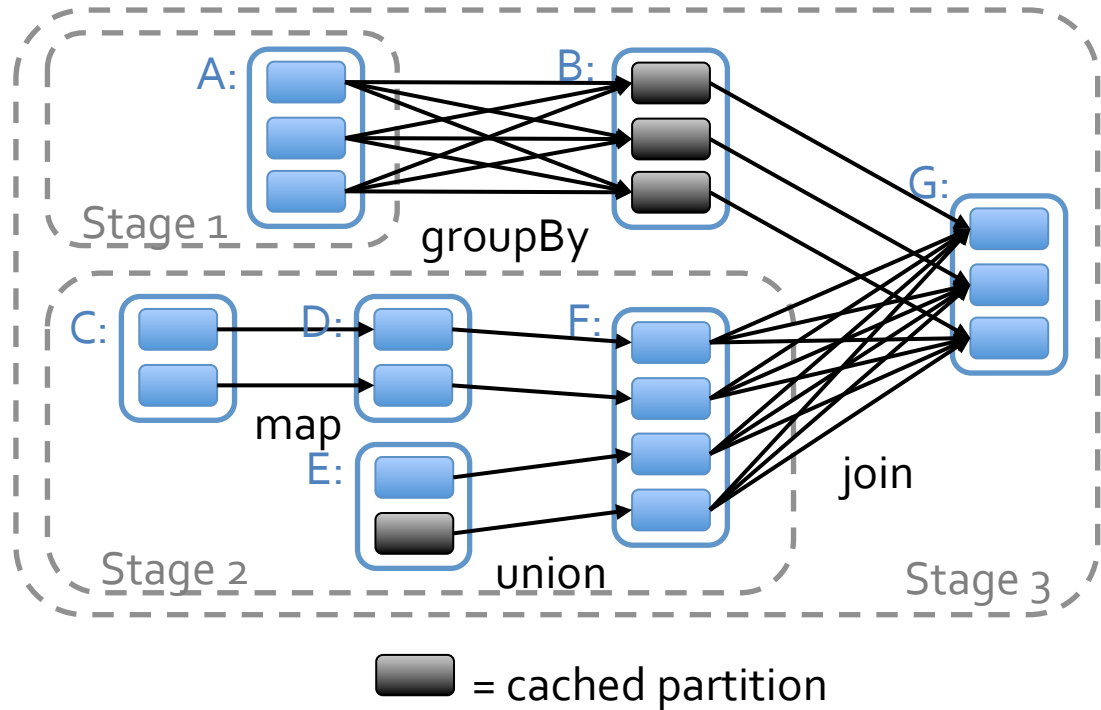
groupByKey



join with inputs not
co-partitioned

JOB SCHEDULER

Captures RDD
dependency graph
Partitions functions
into “stages”
Cache-aware for
data reuse & locality
Partitioning-aware
to avoid shuffles



SPARK ADOPTION

Open source Apache Project, > 1000 contributors

Extensions to SQL, Streaming, Graph processing

Unified Platform for Big Data Applications

QUESTIONS / DISCUSSION ?