

# CS 744: RESILIENT DISTRIBUTED DATASETS

Shivaram Venkataraman

Fall 2020

# ADMINISTRIVIA

- Assignment 1: Due Sep 21, Monday at 10pm!
- Assignment 2: ML will be released Sep 22
- ~~Final project details: Next week~~  
*Course*

→ posted some  
submission notes to  
Piazza!

*Discussion*

# MOTIVATION: 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

*All of our inputs and outputs go to Disk*

# MOTIVATION: 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

↳ grep

Latency or Time to run  
MR  
has a lower bound  
in time to read input  
write output

to disk

# 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();
return 0;
}
```

# APACHE SPARK PROGRAMMABILITY

① fewer lines of code

② Is this because of Scala?  
lambda function

Map Y → functions  
Reduce Y → functions

inline functions

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

counts.save("out.txt")

④ Type checking allows for chaining safely

③ Scala to operate on collections (local)

same API to compute in a distributed setting as local

Scala |  
Python

# APACHE 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 CONCEPTS

once we create this  
we cannot change  
its contents

track changes using lineage

## 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)

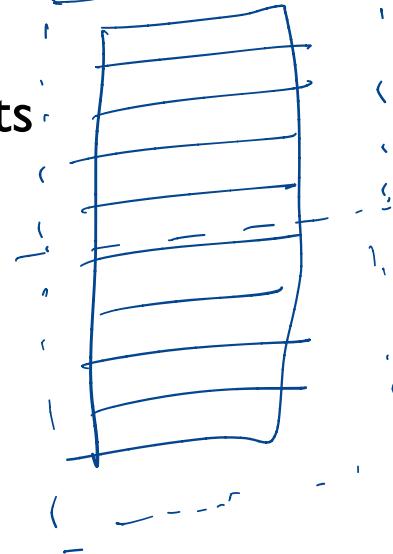
## Restricted shared variables

- Broadcast, accumulators

Int, Type  
File

List class <Type>  
Collection of records

split or  
partitioned



RDD → RDD

list < Int > → list < String >

```

for ( ) {
    read a line
    ↳ filter, map, map, inc count
    ↳ Put this in cache
}

```

# EXAMPLE: LOG MINING



Find error messages present in log files interactively

(Example: HTTP server logs)

*Create an RDD from a file*

```

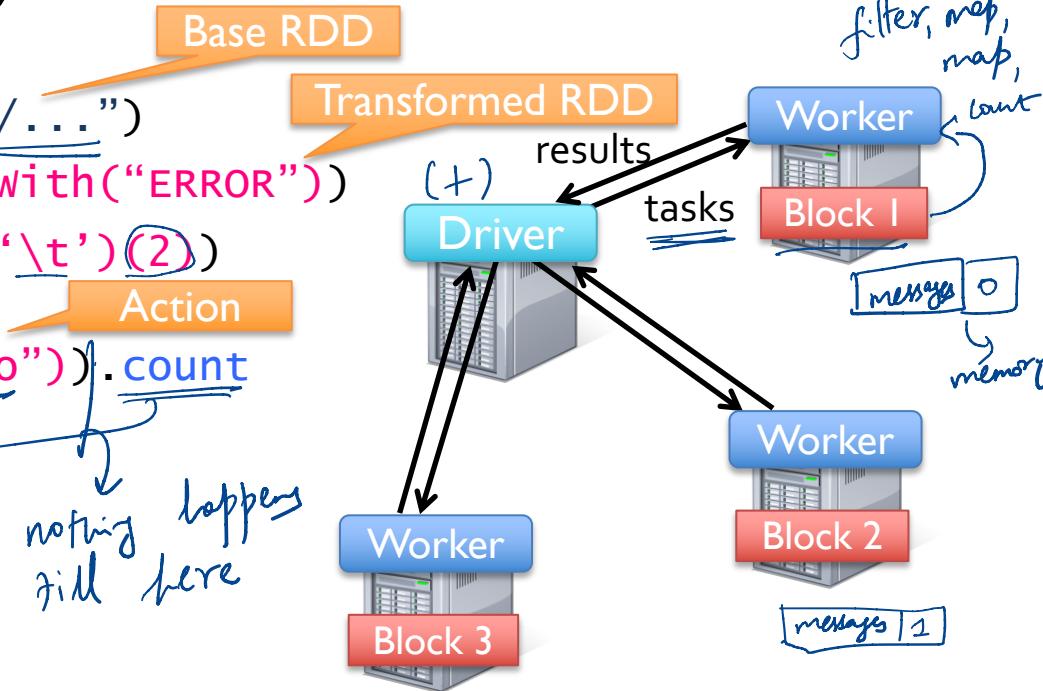
lines = spark.textFile("hdfs://...")          Base RDD
errors = lines.filter(_.startsWith("ERROR"))  Transformed RDD
messages = errors.map(_.split('\t')(2))        Action
messages.cache()                            ???
messages.filter(_.contains("foo")).count()     count

```

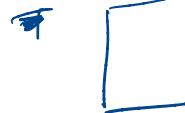
*are handled by Driver*

*memory*

*nothing happens here*



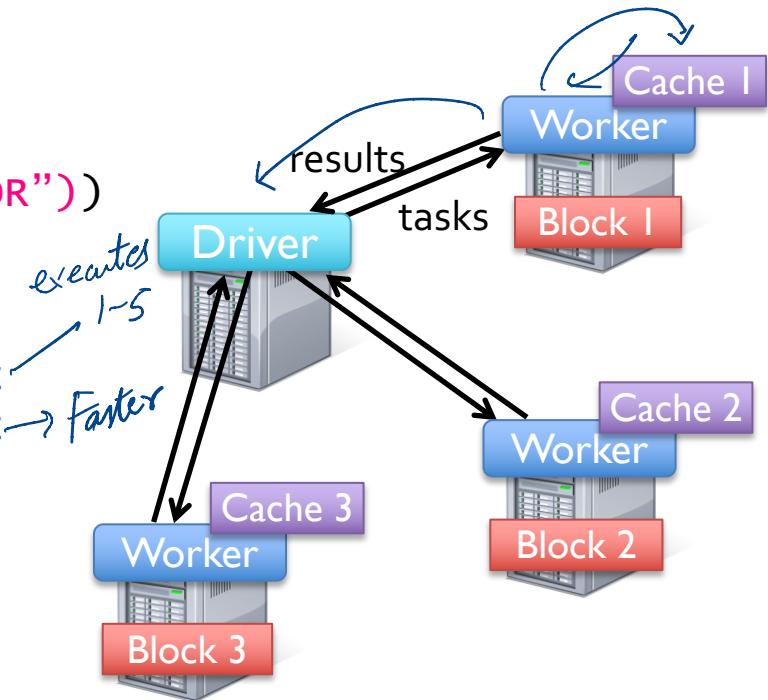
# EXAMPLE: LOG MINING



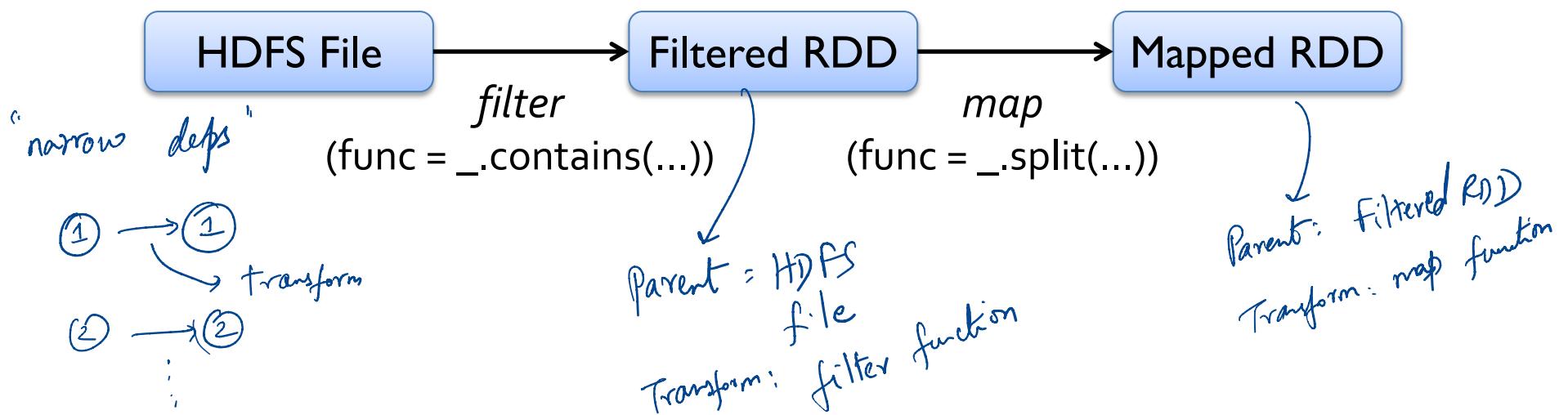
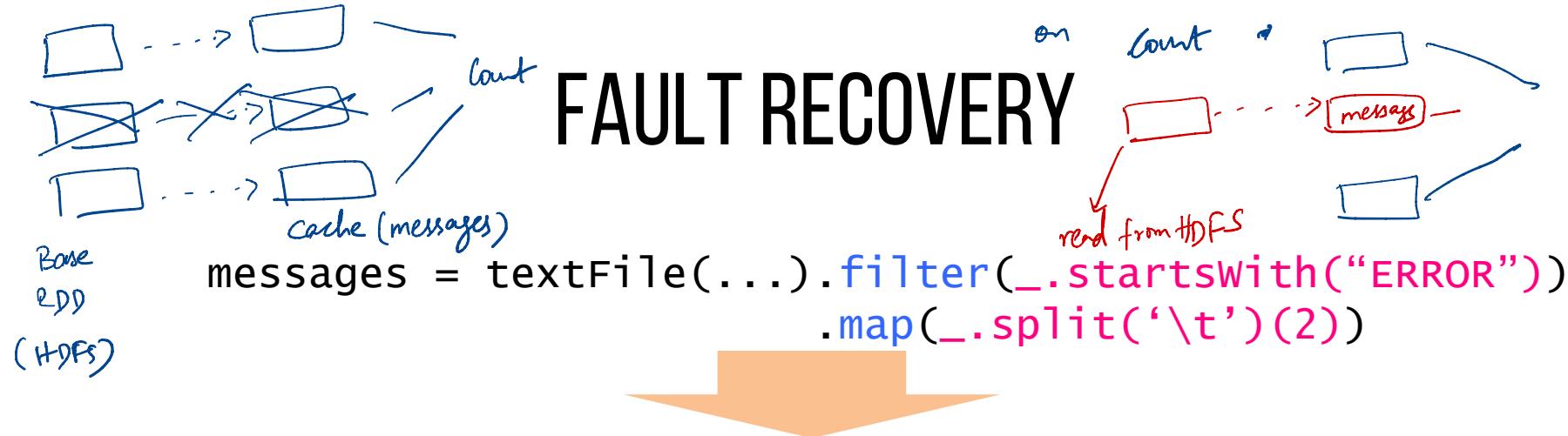
Find error messages present in log files interactively  
(Example: HTTP server logs)

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

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



# FAULT RECOVERY

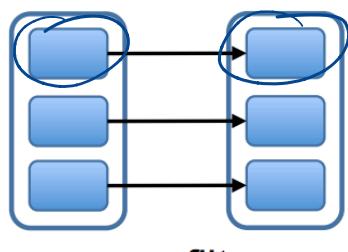


# OTHER RDD OPERATIONS

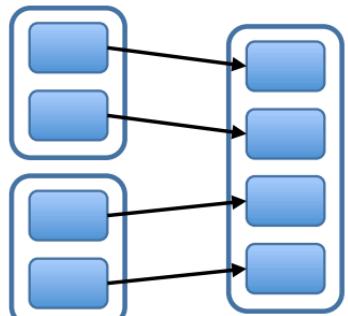
|  |   |   |
|--|---|---|
| <b>Transformations</b><br>(define a new RDD) | map<br>filter<br>sample<br>groupByKey<br>reduceByKey<br>cogroup | flatMap<br>union<br>join<br>cross<br>mapValues<br>... |
| <b>Actions</b><br>(output a result)          | collect<br>reduce<br>take<br>fold                     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

# DEPENDENCIES

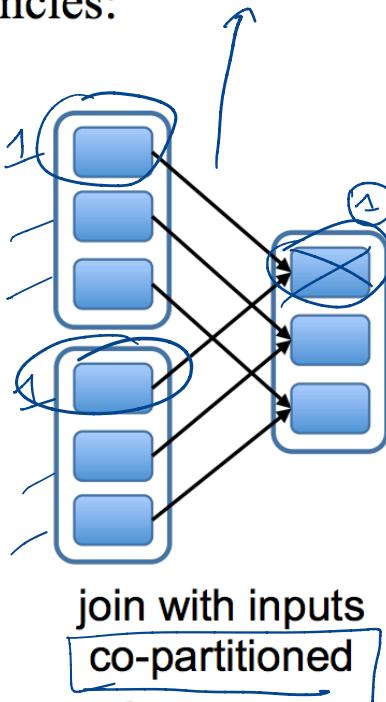
Narrow Dependencies:



map, filter

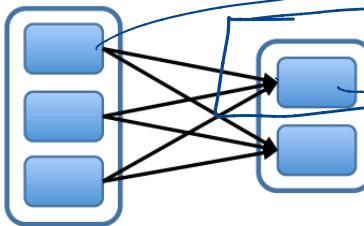


union

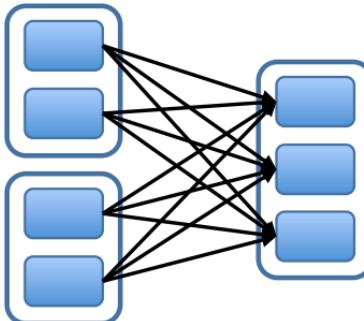


join with inputs  
co-partitioned

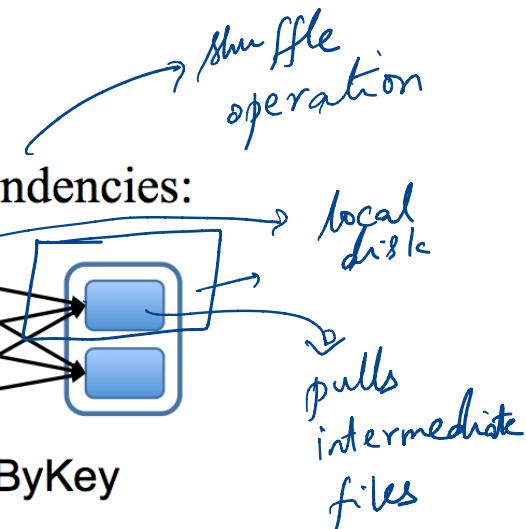
Wide Dependencies:



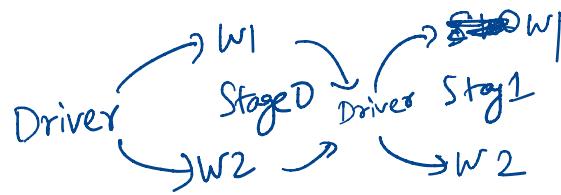
groupByKey



join with inputs not  
co-partitioned



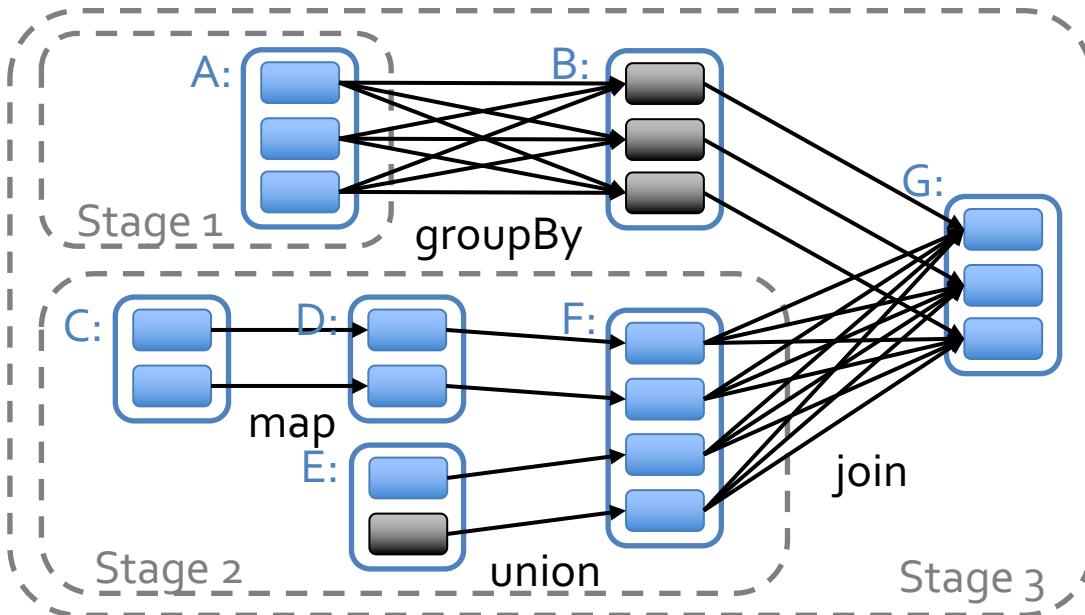
# JOB SCHEDULER (1)



Captures RDD dependency graph

Pipelines functions into “stages”

filter, map, map → one Task  
narrow



# JOB SCHEDULER (2)

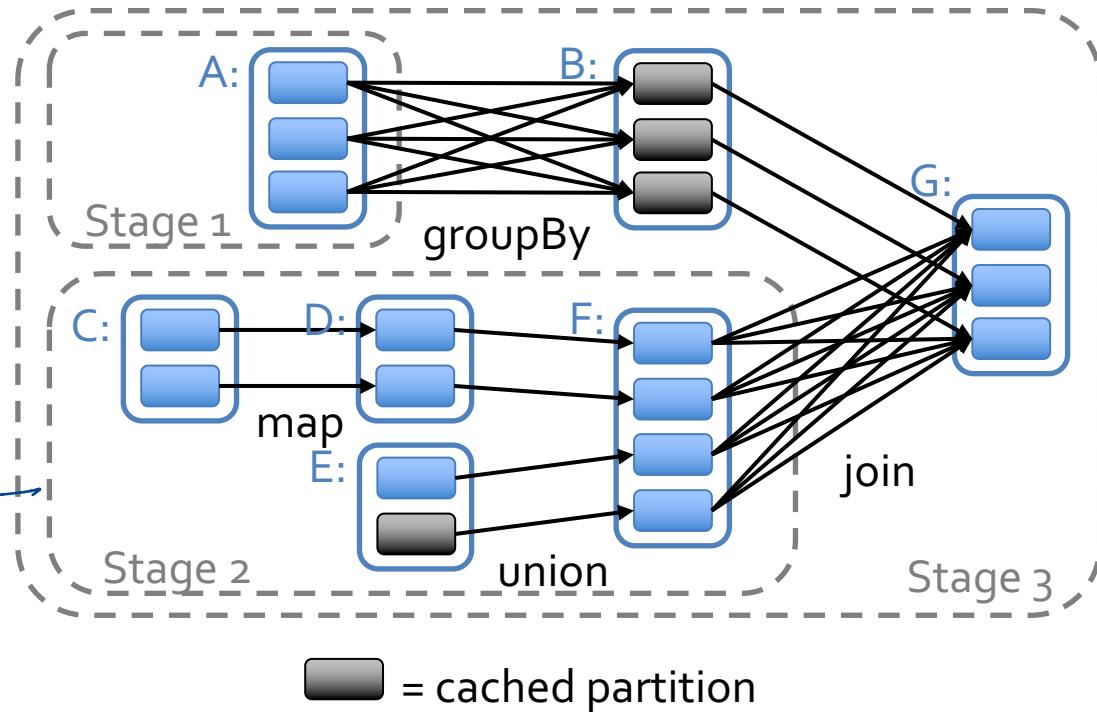
Cache-aware for data reuse, locality

also in MR

Partitioning-aware to avoid shuffles

↓  
co-partition

Scheduler fails?? \*  
Job aborts



# CHECKPOINTING

```
rdd = sc.parallelize(1 to 100, 2).map(x → 2*x)  
rdd.checkpoint()
```

# SUMMARY

Spark: Generalize MR programming model

Support in-memory computations with RDDs

Job Scheduler: Pipelining, locality-aware

reduce → why at Driver?  
Map Reduce  
↳ map  
reduceByKey  
join in Spark  
partitioned

mimic Scala API  
reduce → Spark  
behaves like  
this

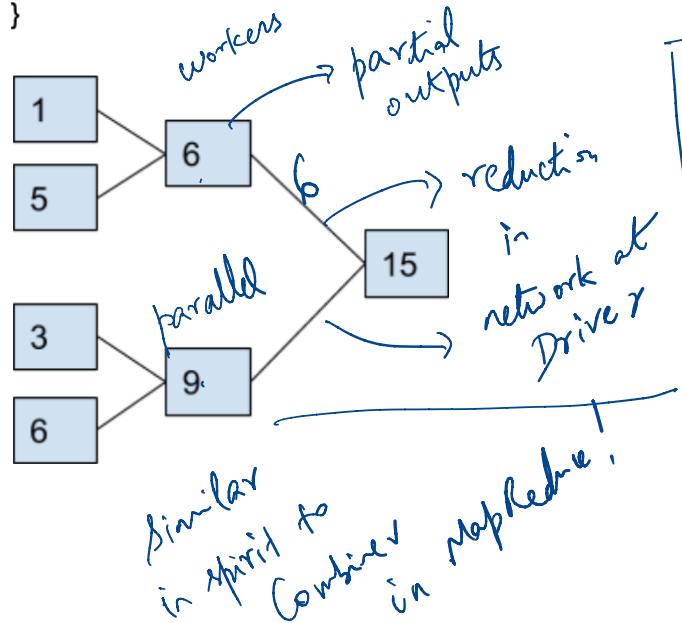
# DISCUSSION

<https://forms.gle/4JDXfpRuVaXmQHxD8>

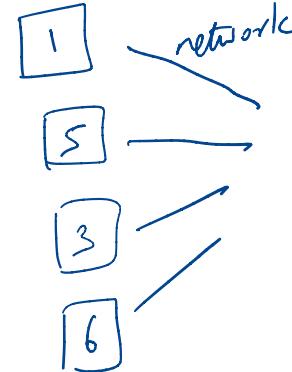


action (like, count)

```
for (i <- 1 to numIters) {
    val modelBC = sc.broadcast(model)
    val grad = data.mapPartitions(iter => gradient(iter, modelBC.value))
    val aggGrad = grad.reduce(case(x, y) => add(x, y))
    model = computeUpdate(aggGrad, model)
}
```



Workers

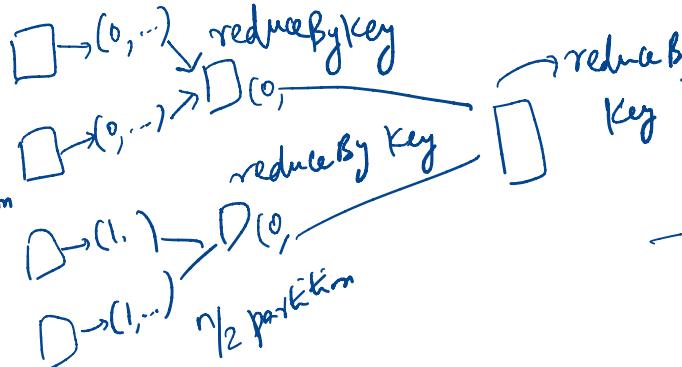


Driver  $(1+5+3+6)$   
 $= 15$

bottleneck

$\hookrightarrow p$  partitions  
 $d$  bytes  
 $p \times d$  bytes = Driver

How can we do this?



tree reduce  
RDD, Scala

When would reduction trees be better than using `reduce` in Spark?

When would they not be ?

Overhead with doing work in stages

→ scheduling, task creation

shuffle overheads

→ Compute & data transmitted is small

→ tree Reduce might be slow.

Device is full error

↳ Disk is full

↳ Shuffle data → local disk

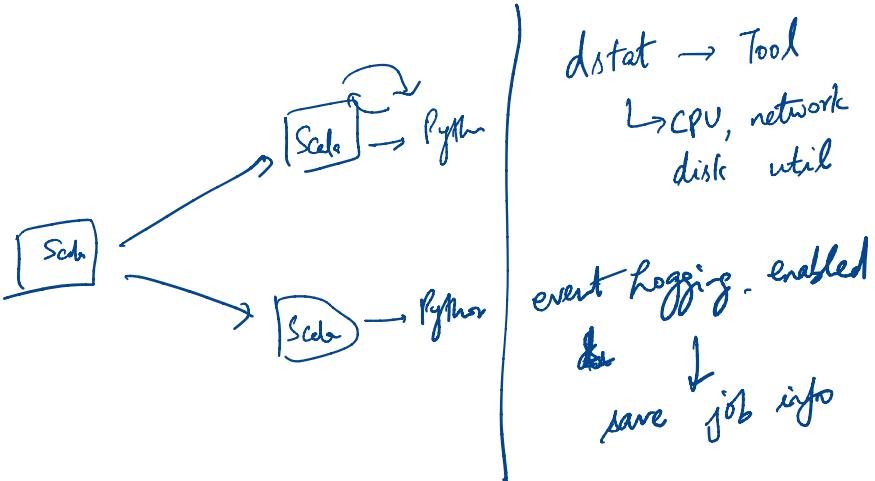
“df - h” worker

/ partition

/mnt → big!

Python

Python



# NEXT STEPS

## Next week: Resource Management

- Mesos, YARN
- DRF

Assignment I is due soon!

FT → Memory replication  
vs  
lineage  
↓  
Frequency of failure

Trade-offs  
network speed vs memory speed

Review form  
↳ When is MR better than Spark  
vice versa  
→ Spark is better when multiple passes over data & data fits in memory  
  
MR is better when → single pass  
→ scheduling overhead?