# CS 744: RESILIENT DISTRIBUTED DATASETS

Shivaram Venkataraman

Fall 2019

# ADMINISTRIVIA

- Assignment 1: Due Sep 24
- Project details
  - Ideas posted on Piazza by Sat.
  - Come up with your own ideas!
  - Submit groups, topics by *9/30*
  - Meet? Office hours 9/23 or 9/30

# 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

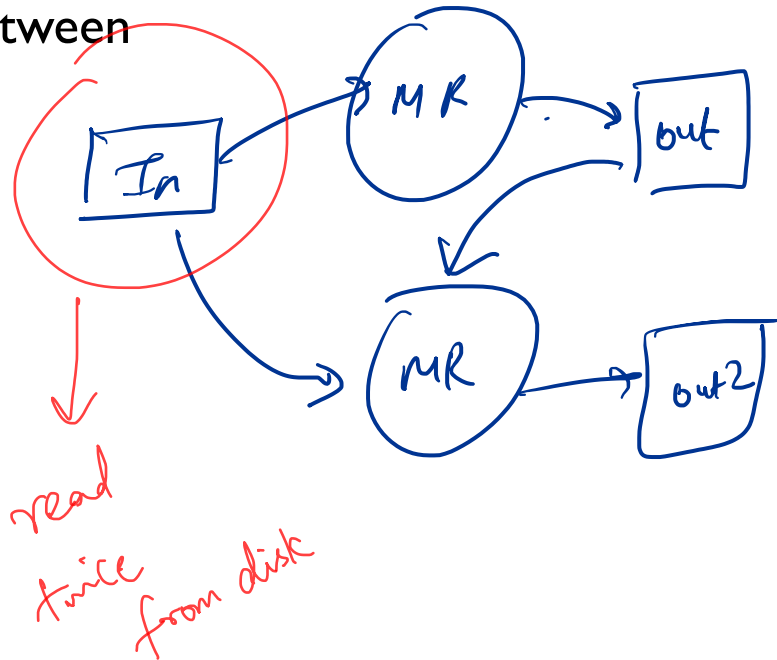- – 21 MR steps → 21 mapper and reducer classes

# 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

# PROGRAMMABILITY

## Google MapReduce WordCount:

```cpp
#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);
```

```cpp
// 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);
```

```cpp
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;
}
```
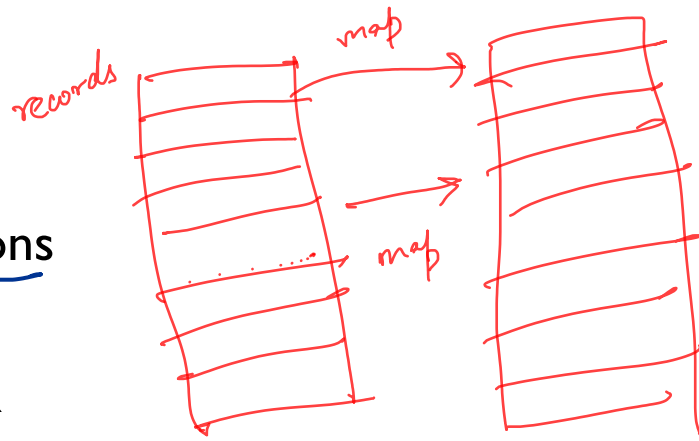
# APACHE SPARK PROGRAMMABILITY

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

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

# 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

Resilient distributed datasets (RDDs)
- – Immutable, partitioned collections of objects
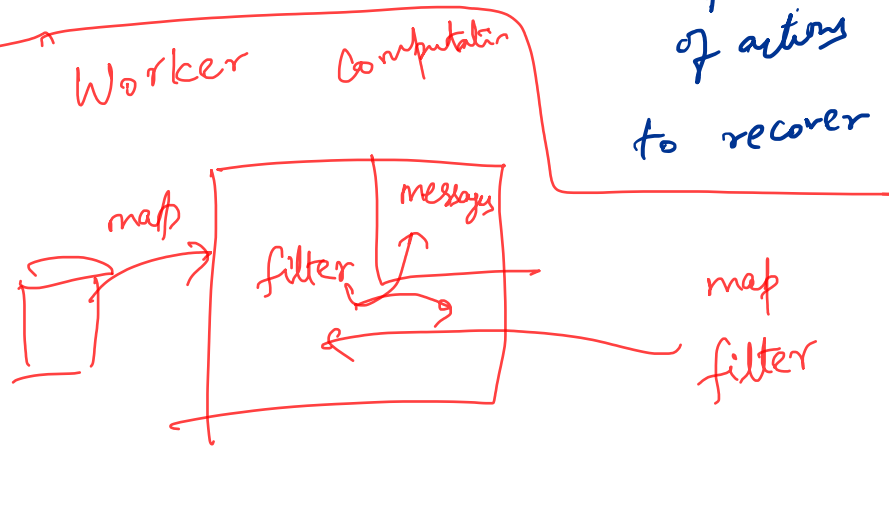- – May be cached in memory for fast reuse

Operations on RDDs
- – *Transformations* (build RDDs)
- – *Actions* (compute results)

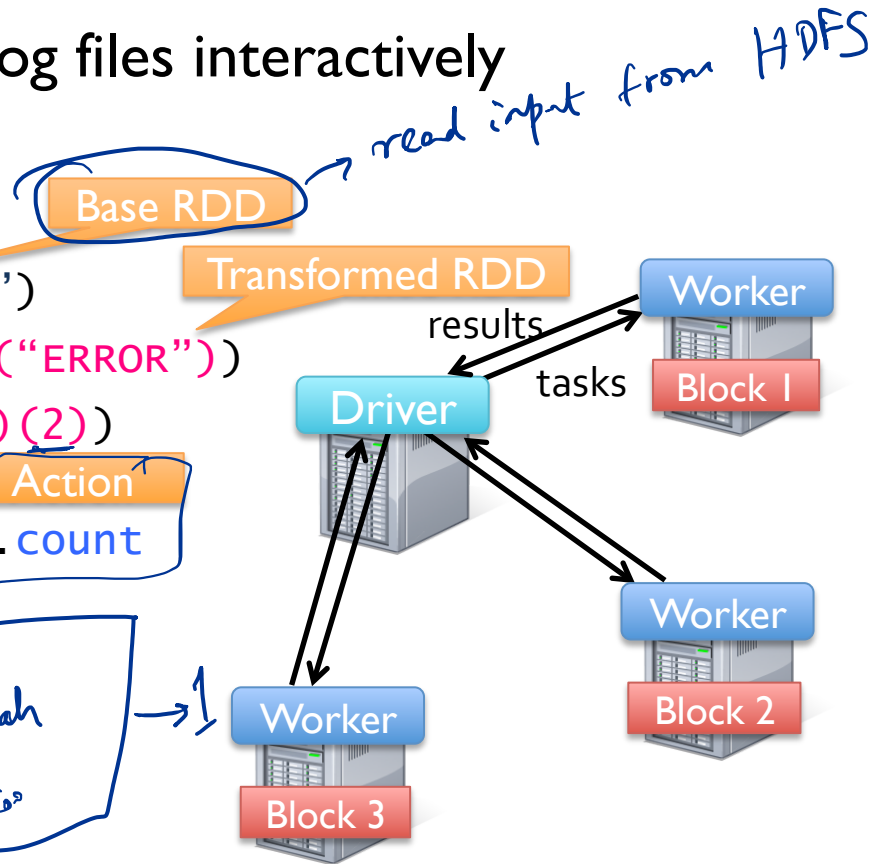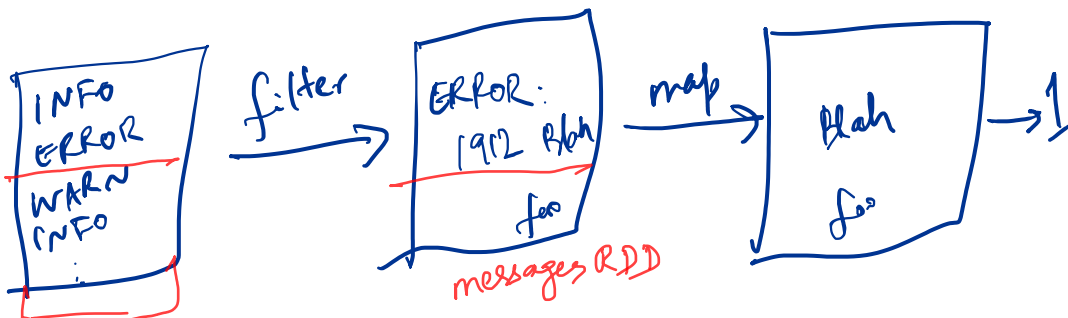Restricted shared variables
- – Broadcast, accumulators

*Lazy execution model*
*⇒ Optimizations*

# EXAMPLE: LOG MINING

Find error messages present in log files interactively

(Example: HTTP server logs)

*read input from HDFS*

**Base RDD**

**Transformed RDD**

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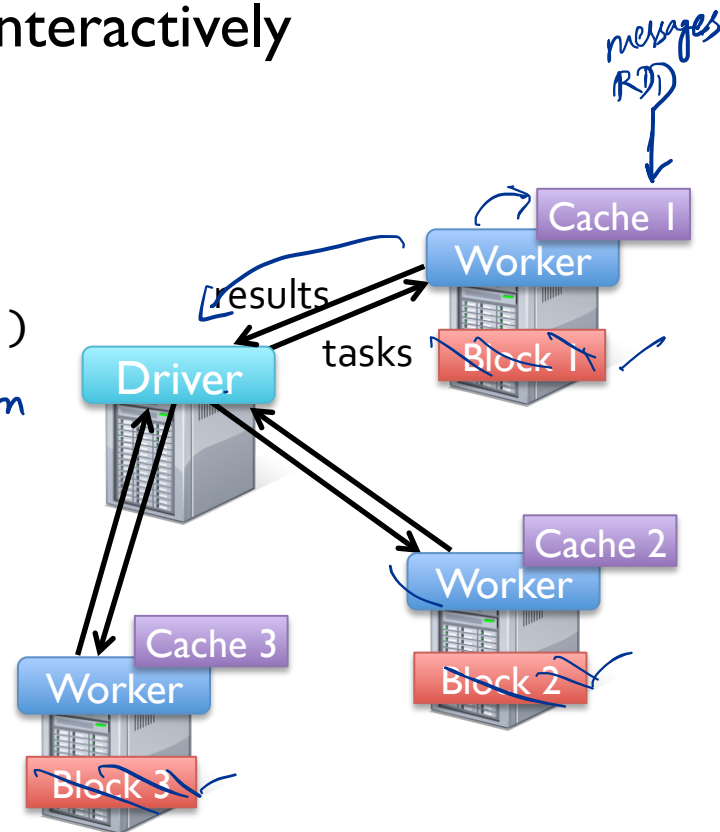
**Action**

# EXAMPLE: LOG MINING

*Locality*

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

*messages RDD*

```
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
. . . .
```
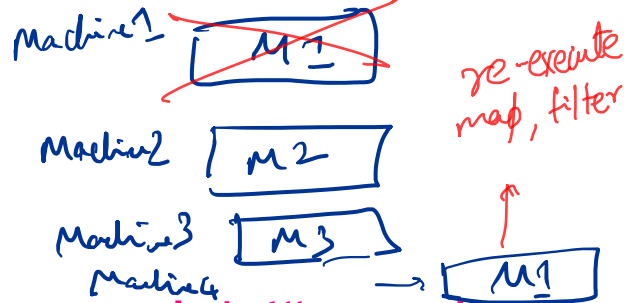
*Action*

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

Driver

results

tasks

Worker — Cache 1 — Block 1

Worker — Cache 2 — Block 2

Worker — Cache 3 — Block 3

# FAULT RECOVERY

lines =
lines. persist()
lines. take (10)

Machine1 ~~M1~~

Machine2  M2

Machine3  M3

Machine4  → M1

re-execute
map, filter

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

. filter ( ) . count



| HDFS File | → | Filtered RDD | → | Mapped RDD |

*filter*
(func = _.contains(...))

*map*
(func = _.split(...))

batch
records
~4K

filter → map → filter  → count

Iterator model

# SHARED VARIABLES

Program Driver

local variables

"Spark" variables

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

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

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)
```

input variables

shared variables

closure

local data

adds up results

worker closure.run()

(Task) Closure(M)

closure

# BROADCAST VARIABLES

```scala
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)
```
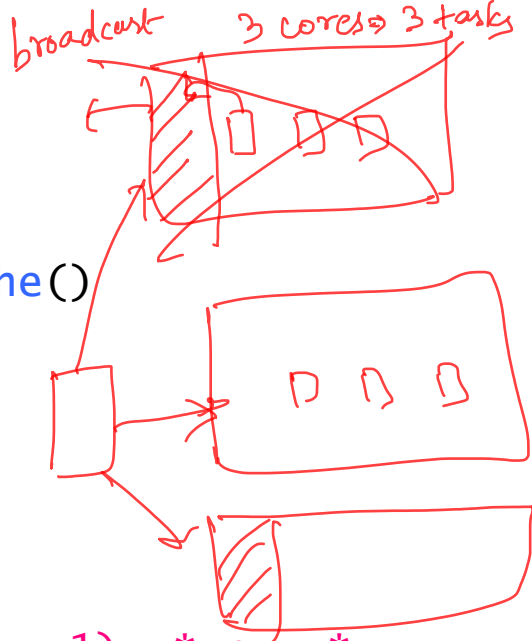
*(handwritten annotations:)* broadcast · 3 cores ⇒ 3 tasks · M. destroy()

# OTHER RDD OPERATIONS

*flatMap*
*reduceByKey*
*save As Hadop*

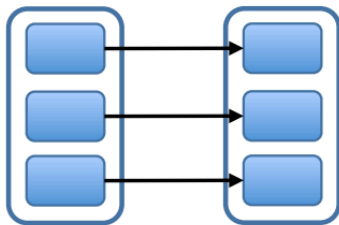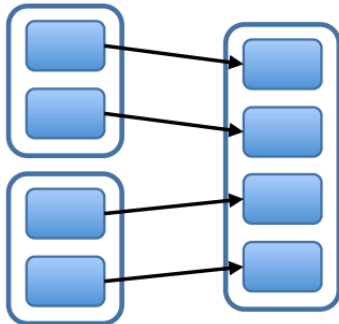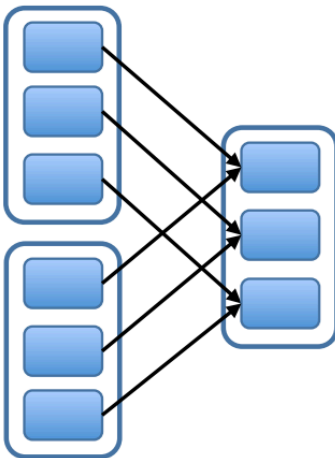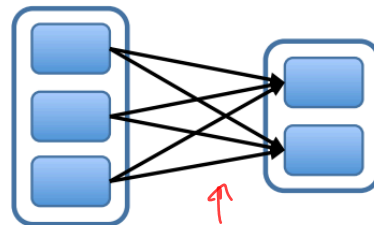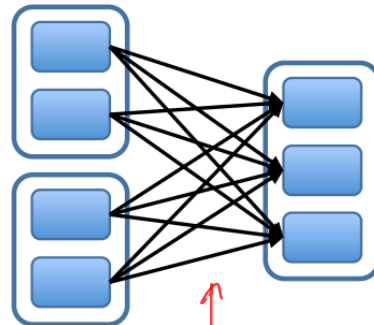|  |  |  |
|---|---|---|
| **Transformations**<br>(define a new RDD) | map<br>filter<br>sample<br>groupByKey<br>reduceByKey<br>cogroup | flatMap<br>union<br>join<br>cross<br>mapValues<br>... |
| **Actions**<br>(output a result) | collect<br>reduce<br>take<br>fold | count<br>saveAsTextFile<br>saveAsHadoopFile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

# DEPENDENCIES

Narrow Dependencies:



map, filter

union

join with inputs
co-partitioned

Wide Dependencies:
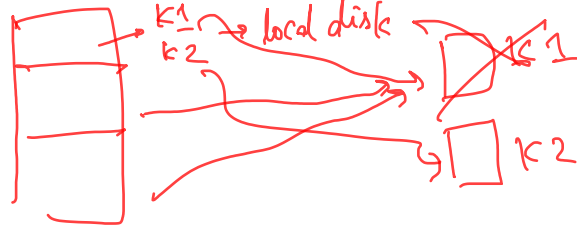
groupByKey

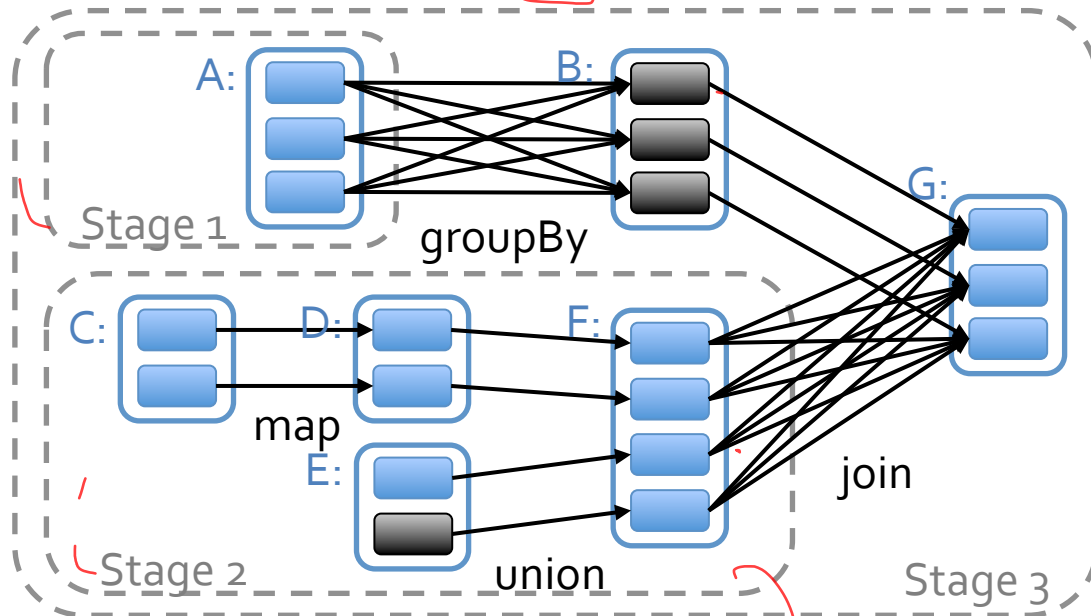join with inputs not
co-partitioned

# JOB SCHEDULER

B = A.groupByKey()

Captures RDD dependency graph

Pipelines functions into "stages"

Cache-aware for data reuse, locality
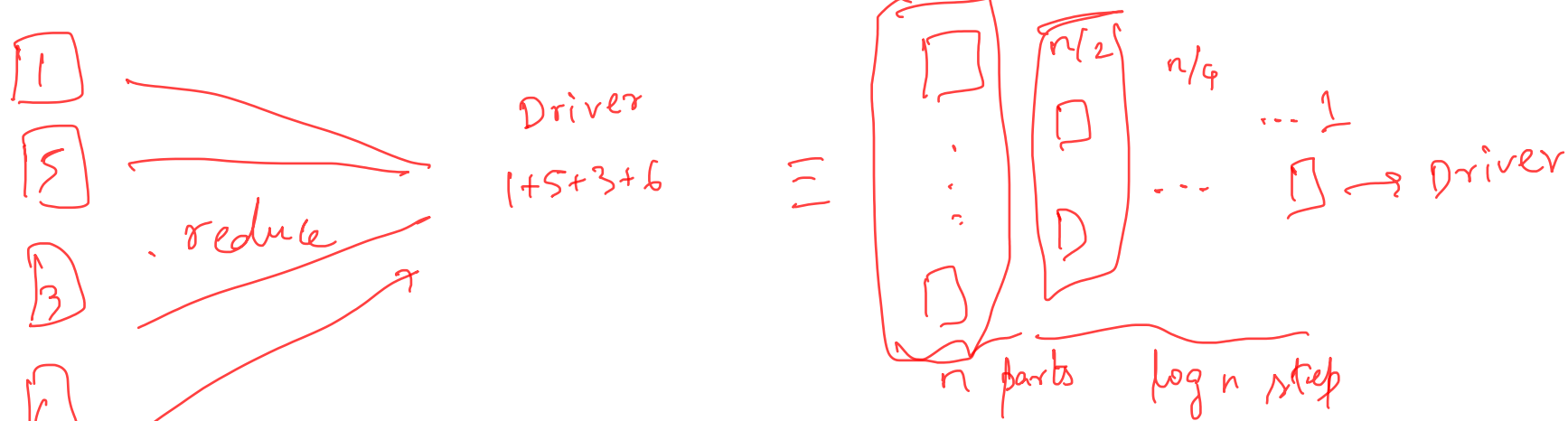
Partitioning-aware to avoid shuffles

A:

B:

Stage 1

groupBy

C:

D:

map

E:

F:

union

G:

join

Stage 2

Stage 3

= cached partition

narrow deps

# CHECKPOINTING

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

# DISCUSSION

https://forms.gle/Gg2K1hsGFJpFbmSj9

1
5
3
6

. reduce

Driver
1+5+3+6

$\equiv$

n/2
n/4
... 1

$\rightarrow$ Driver

n parts        log n step

Co group        2 RDDs

↳ reduce By Key        n/2 keys

↳ Partitioner to
Control
keyspace

while (
rdd = rdd. mapPartitions { (r, idx) =>
    (idx/2, r)
}. reduce By Key ( n/2 )
n=n/2
}

# SPARK ADOPTION

Open source Apache Project, > 1000 contributors

Extensions to SQL, Streaming, Graph processing

Unified Platform for Big Data Applications

# NEXT STEPS

- Next week: Resource Management
    - Mesos, YARN
    - DRF
- Assignment 1 is due soon!