# High-Performance Big Data Processing Tools for Neuroscience and A Demo on Chameleon Cloud

Xiaoyi Lu

The Ohio State University
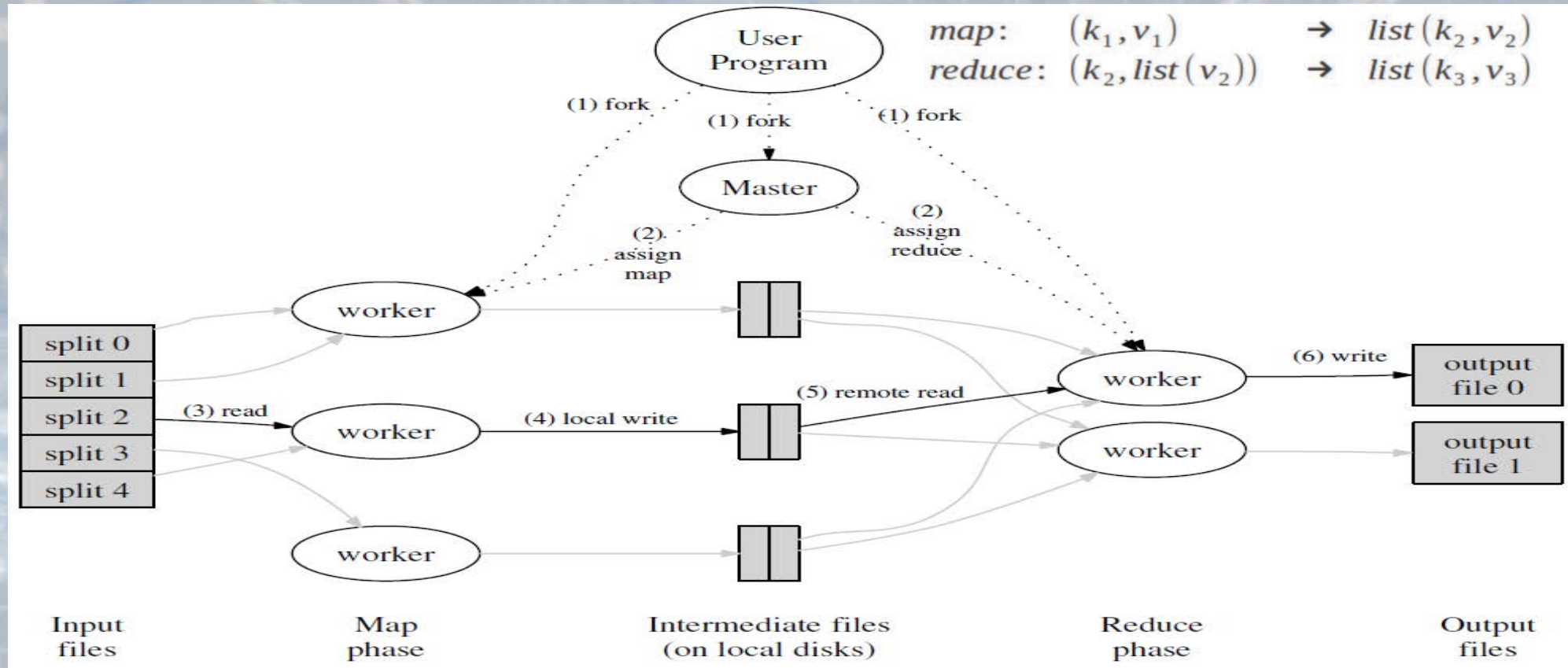
luxi@cse.ohio-state.edu

http://web.cse.ohio-state.edu/~luxi/
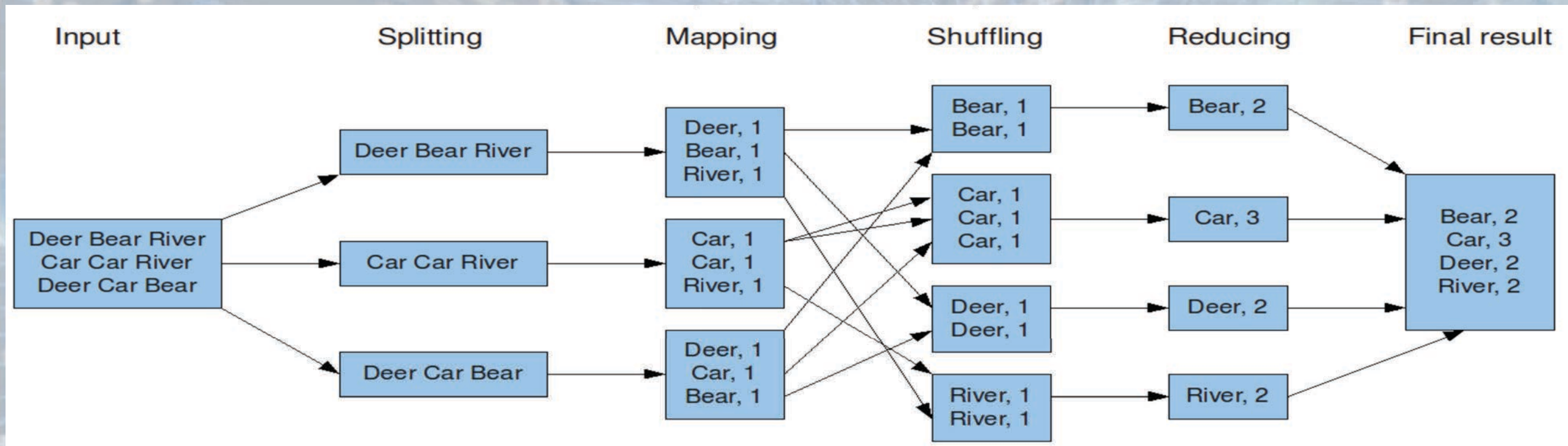
THE OHIO STATE
UNIVERSITY

# The MapReduce Model



J. Dean and S. Ghemawat. *MapReduce: Simplified Data Processing on Large Clusters*. In *Proceedings of the 6th Symposium on Operating Systems Design & Implementation (OSDI'04)*, 2004.
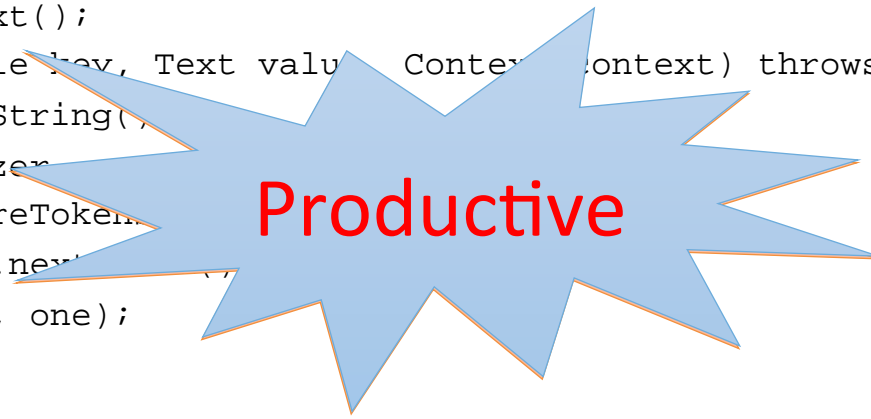
# WordCount Execution

- The overall execution process of WordCount in MapReduce
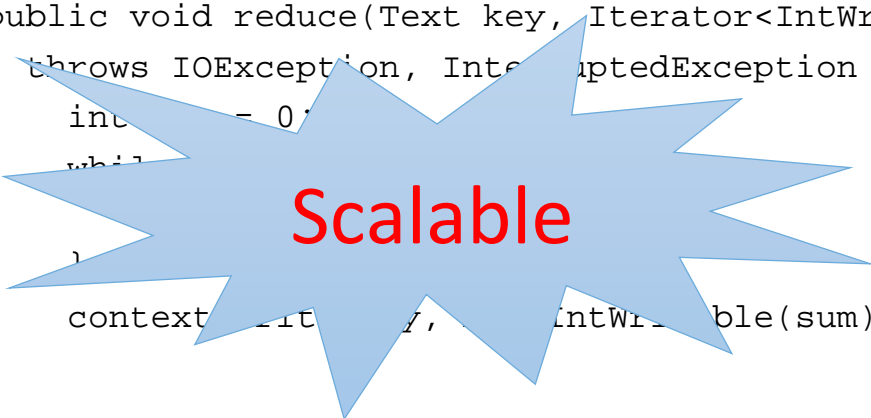
# Word Count in Hadoop!
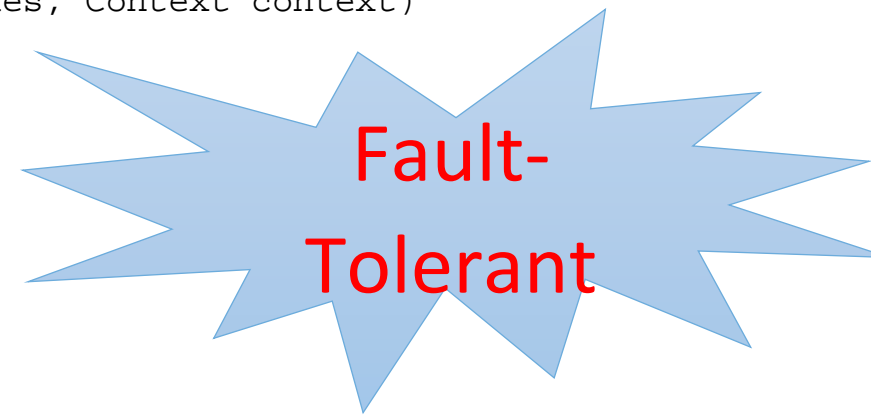
```
public class WordCount {
 public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
     private final static IntWritable one = new IntWritable(1);
     private Text word = new Text();
     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
         String line = value.toString();
         StringTokenizer tokenizer
         while (tokenizer.hasMoreToken
             word.set(tokenizer.nex
             context.write(word, one);
         }
     }
 }
 public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
     public void reduce(Text key, Iterator<IntWritable> values, Context context)
        throws IOException, InterruptedException {
         int      = 0;
         whil
         }
         context          y,       intWr    ble(sum));
     }
 }
}
```
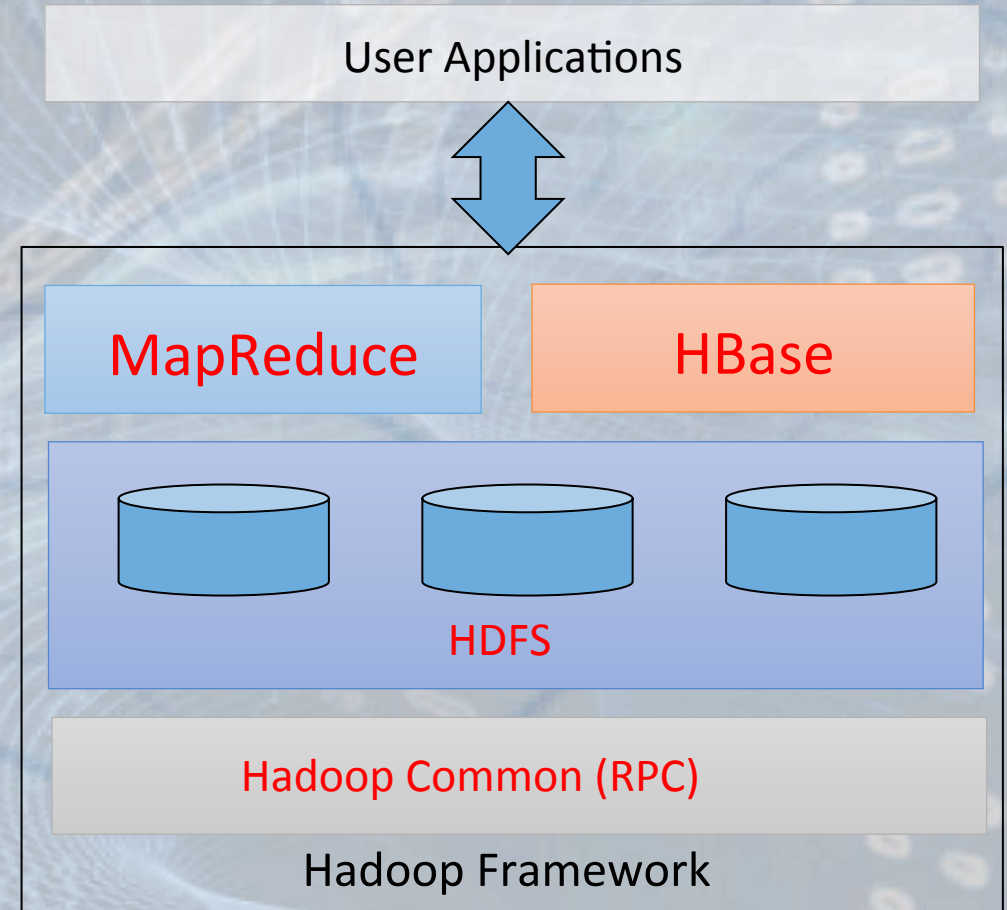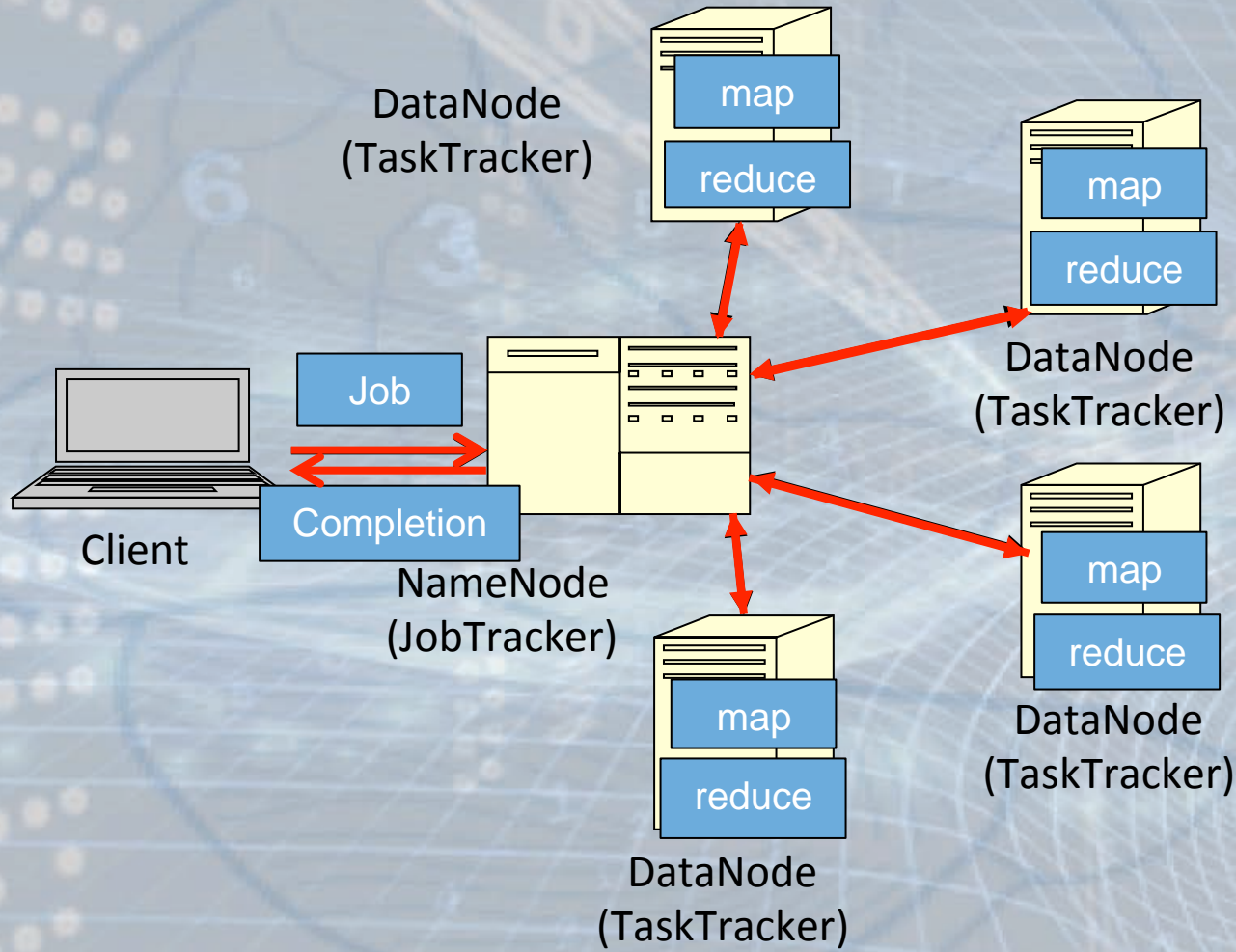
**Productive**

**Scalable**

**Fault-Tolerant**

# Big Data Processing with Hadoop Components

- Major components included in this tutorial:
  - MapReduce (Batch)
  - HBase (Query)
  - HDFS (Storage)
  - RPC (Inter-process communication)
- Underlying Hadoop Distributed File System (HDFS) used by both MapReduce and HBase
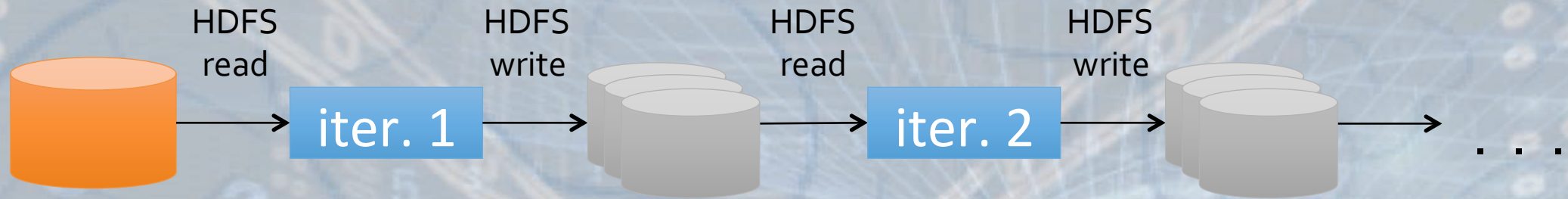- Model scales but high amount of communication during intermediate phases can be further optimized

# MapReduce Job Execution on Hadoop



- Main Features
  - Replication (e.g. 3)
  - Data locality for Maps
  - HTTP-based Shuffle
  - Speculative execution
  - Independence among tasks
  - ...

- Goals
  - Fault Tolerance
  - Scalability

# Data Sharing Problems in MapReduce



Input → HDFS read → iter. 1 → HDFS write → HDFS read → iter. 2 → HDFS write → . . .

Slow due to replication, serialization, and disk IO

In-Memory?

Input → iter. 1 → iter. 2 → . . .

10-100× faster than network and disk

# RDD Programming Model in Spark

- Key idea: *Resilient Distributed Datasets (**RDDs**)*
  - Immutable distributed collections of objects that can be cached in memory across cluster nodes
  - Created by transforming data in stable storage using data flow operators (map, filter, groupBy, …)
  - Manipulated through various parallel operators
  - Automatically rebuilt on failure
    - rebuilt if a partition is lost

- Interface
  - Clean language-integrated API in Scala (Python & Java)
  - Support R, distributed machine learning using MLlib
  - Can be used *interactively* from Scala console

# RDD Operations

| Transformations (define a new RDD) | Actions (return a result to driver) |
| --- | --- |
| map | reduce |
| filter | collect |
| sample | count |
| union | first |
| groupByKey | Take |
| reduceByKey | countByKey |
| sortByKey | saveAsTextFile |
| join | saveAsSequenceFile |
| ... | ... |

More Information:
- https://spark.apache.org/docs/latest/programming-guide.html#transformations
- https://spark.apache.org/docs/latest/programming-guide.html#actions

# Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

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

cachedMsgs.filter(_.contains("foo")).count
cachedMsgs.filter(_.contains("bar")).count
. . .
```

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

Base Transformed RDD

results

tasks

Driver

Action

Worker — Cache 1 — Block 1

Worker — Cache 2 — Block 2

Worker — Cache 3 — Block 3

Courtesy: https://spark.apache.org/

# Example: Word Count in Spark!

```scala
val fi        extFile("hdfs.        )
val cou    e.flatMap(line =>       )
           .map(word => (word,  ))
           .reduceByKey(_ + _)

counts.saveAsTextFile("hdfs://...")
```

Productive

High-Performance

Scalable

Fault-Tolerant

# Data Movement in MapReduce



- Map and Reduce Tasks carry out the total job execution
  - Map tasks read from HDFS, operate on it, and write the intermediate data to local disk
  - Reduce tasks get these data by shuffle from TaskTrackers, operate on it and write to HDFS
- Communication in shuffle phase uses HTTP over Java Sockets
- Similar bottlenecks exist in HDFS, RPC, Spark, HBase, Memcached, ec.

# The High-Performance Big Data (HiBD) Project

- **http://hibd.cse.ohio-state.edu**

- RDMA for Apache Spark

- RDMA for Apache Hadoop 2.x (RDMA-Hadoop-2.x)

  - Plugins for Apache, Hortonworks (HDP) and Cloudera (CDH) Hadoop distributions

- RDMA for Apache HBase

- RDMA for Memcached (RDMA-Memcached)

- RDMA for Apache Hadoop 1.x (RDMA-Hadoop)

- OSU HiBD-Benchmarks (OHB)

  - HDFS, Memcached, and HBase Micro-benchmarks

- Users Base: 190 organizations from 26 countries

- More than 17,800 downloads from the project site

- RDMA for Impala (upcoming)

**Available for InfiniBand and RoCE**

HiBD
High-Performance Big Data

Network Based Computing
Laboratory

THE OHIO STATE
UNIVERSITY

# Performance Benefits of HiBD on PageRank, Sort, and TeraSort



**SDSC Comet 32 Worker Nodes, 768 cores**
**PageRank Total Time**

**SDSC Comet 64 Worker Nodes, 1536 cores**
**PageRank Total Time**

**TACC Stampede 32 Worker Nodes,**
**Sort Total Time**

**TACC Stampede 32 Worker Nodes,**
**TeraSort Total Time**

# A Case Study with NeuroPigPen

- **NeuroPigPen Toolkit[1]:** a data management toolkit using Hadoop Pig for processing electrophysiological signals in Neuroscience applications
  - Dr. Sahoo, etc., Case Western Reserve University



The data processing workflow supported by the NeuroPigPen toolkit modules consists of multiple steps with EDF files as input and CSF files as output. The Load functions in the toolkit extend the Hadoop FileInputFormat and FileInputRecordReader classes to support signal data. The Map functions in the toolkit are automatically compiled into MapReduce tasks by the Apache Pig compiler. The intermediate and final results are stored in Hadoop Distributed File System (HDFS), which provides a reliable and scalable storage platform for signal data.

1. Sahoo SS, Wei A, Valdez J, Wang L, Zonjy B, Tatsuoka C, Loparo KA, Lhatoo SD. NeuroPigPen: a Data Management Toolkit using Hadoop Pig for Processing Electrophysiological Signals in Neuroscience Applications, http://www.ncbi.nlm.nih.gov/pubmed/27375472

# Demo on NSF Supported Chameleon Cloud

- Five InfiniBand bare-metal nodes (with SR-IOV support) on Chameleon Cloud
- 25 VMs in total, 6 VMs per node
- RDMA-Hadoop appliance available on Chameleon Cloud
  - Pre-installed OS, drivers, and software packages
  - Automatically configuration
- Deployment of virtual machines
- RDMA-Hadoop examples
  - RandomWriter
  - Sort

# NSF Chameleon Cloud: A Powerful and Flexible Experimental Instrument

- Large-scale instrument
  - Targeting Big Data, Big Compute, Big Instrument research
  - ~650 nodes (~14,500 cores), 5 PB disk over two sites, 2 sites connected with 100G network

- Reconfigurable instrument
  - Bare metal reconfiguration, operated as single instrument, graduated approach for ease-of-use

- Connected instrument
  - Workload and Trace Archive
  - Partnerships with production clouds: CERN, OSDC, Rackspace, Google, and others
  - Partnerships with users

- Complementary instrument
  - Complementing GENI, Grid'5000, and other testbeds

- Sustainable instrument
  - Industry connections

http://www.chameleoncloud.org/

# Demo - Login to Chameleon

# Demo – Overview of Your Chameleon Resources

# Demo – Create Your Lease

# Demo – Create Your Lease

# Demo – Create Your Lease

# Demo – Create Your Lease

# Demo – Launch Your Instances

# Demo – Launch Your Instances

# Demo – Launch Your Instances

# Demo – Associate Floating IP

# Demo – Login to Your Instances

# Demo – Play with Your Hadoop on Chameleon Cloud

- ssh cc@129.114.108.229
- cd rdma-hadoop-2.x-0.9.8/
- ./sbin/start-all.sh
- bin/hdfs dfsadmin –report
- bin/mapred job -list-active-trackers
- ./bin/hdfs dfs -ls /
- bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-*.jar randomwriter -Dmapreduce.randomwriter.mapsperhost=4 -Dmapreduce.randomwriter.bytespermap=1000 /rw-output
- ./bin/hdfs dfs -ls /
- bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-*.jar sort /rw-output /sort-output
- ./bin/hdfs dfs -ls /
- ./sbin/stop-all.sh

# Conclusion

- Overview of MapReduce Programming Model and Hadoop Architecture

- Overview of RDD Programming Model and Spark Architecture

- Overview of HiBD project: accelerating Big Data processing middleware (e.g., Spark, Hadoop, Memcached) by taking advantage of HPC technologies, such as InfiniBand/RDMA, SSD, etc.

- A case study with NeuroPigPen

- Demo of RDMA-Hadoop on NSF supported Chameleon Cloud

- Soliciting collaboration with other groups
  - Interested in exploring BigData processing requirements from Neuroscience researchers
  - We can explore how to collaborate further

# Thank You!

luxi@cse.ohio-state.edu

Network-Based Computing Laboratory
http://nowlab.cse.ohio-state.edu/

The MVAPICH2 Project
http://mvapich.cse.ohio-state.edu/

The High-Performance Big Data Project
http://hibd.cse.ohio-state.edu/

**Computational Neuroscience Network (ACNN)**

http://www.NeuroscienceNetwork.org/ACNN_Workshop_2016.html