
Big Data - Hadoop/MapReduce

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Logistics

- Quiz 1 next week. Sample released.
- Assignment 2 (tentative) draft uploaded.
- In-depth exploration topics
 - Kubernetes and DevOps: you will explore containers, Kubernetes and its application to cloud application devOps pipeline.
 - Kafka and Events/Logs: you will learn how it is used for handling events/logs including internals, use cases and how it is used for handling logs efficiently.
 - Elastic Search and data pipelines: you will learn how to effectively build data pipelines to process unstructured data and build data pipeline to process data efficiently.
 - Spark and cluster computing: learn how to effectively demonstrate compute in a cluster including scale, performance tuning, and use cases for performance tuning.
- Paper reviews

Assignment Project Exam Help

<https://eduassistpro.github.io/>

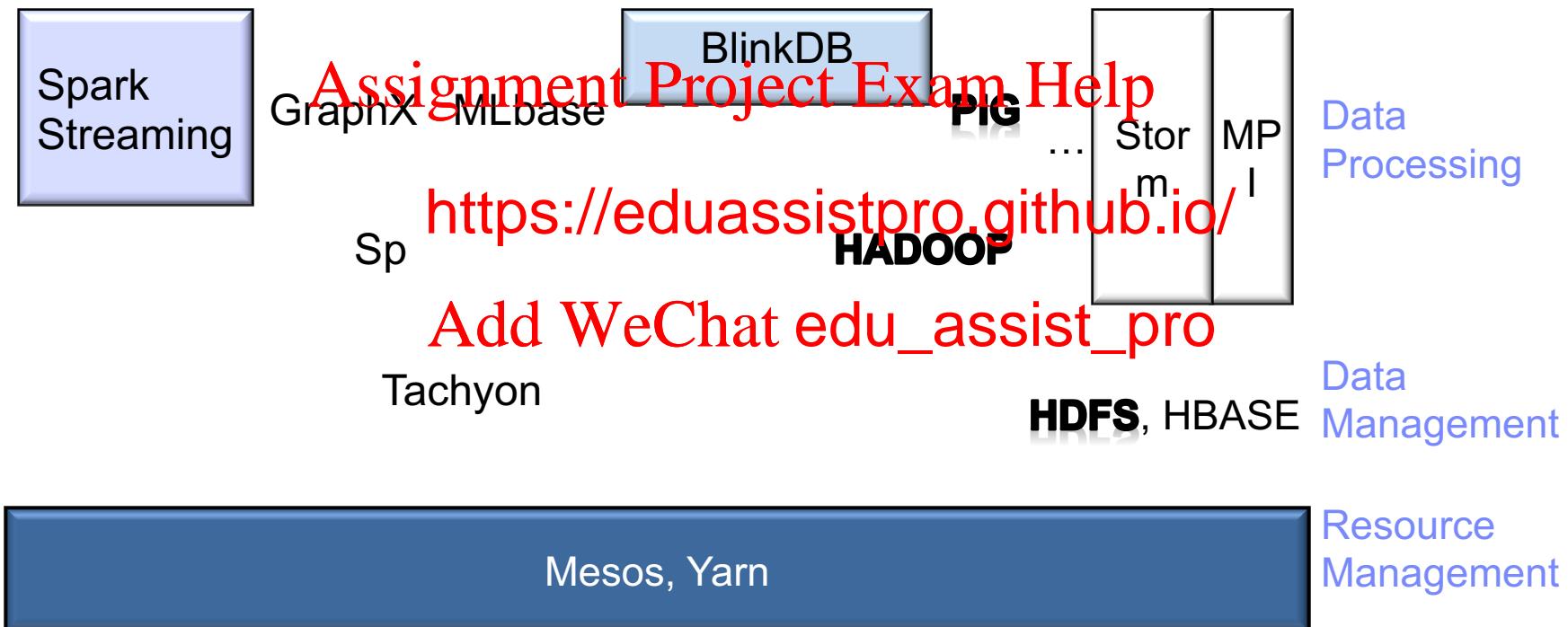
Add WeChat [edu_assist_pro](#)

Cluster Computing

- Designing a cluster
 - Automation of tasks in a cluster
 - Assignment Project Exam Help**
 - Programming in a cluster
 - Hadoop/MapReduce
 - Spark
- <https://eduassistpro.github.io/>**
- Add WeChat edu_assist_pro**

Big Data System Architecture

- Current system modules for each software layer



Agenda

- Why Big Data?

Assignment Project Exam Help

- Apache Hado <https://eduassistpro.github.io/>
 - Introduction Add WeChat edu_assist_pro
 - Architecture
 - Programming

Hypothetical Job

- You just got an awesome job at data-mining start-up ..
Congratulations !!

- Free Snacks, Soda and Coffee Yaaay
 - Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Your first day of work you are given
Add WeChat edu_assist_pro
 - The company has a new algorithm they want you to test.
 - Your boss gives you
 - The algorithm library
 - A test machine and
 - 1GB input data file

Java Document Scorer Program

Assignment Project Exam Help *Read Input*

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Throughput 1GB per hour.

What if we wanted to process 10GB data set? 10hours!!
How can we improve the performance?

Some Options

1. Faster CPU
 2. More Memory
 3. Increase the number of cores
 4. Increase the number of threads and cores
 5. Increase the number of threads and cores
- Assignment Project Exam Help**
- <https://eduassistpro.github.io/>**
- Add WeChat edu_assist_pro**

Java Document Scorer Program – Multi Threaded

Throughput 4GB per hour.

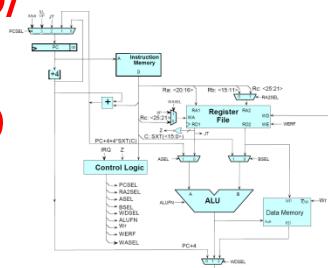
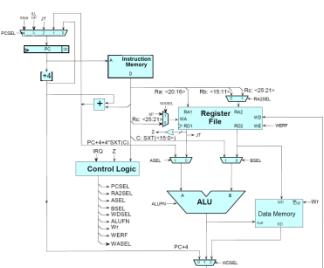
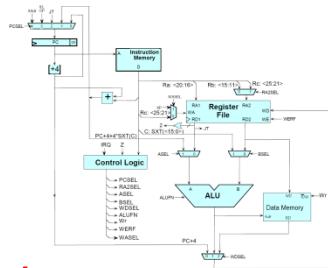
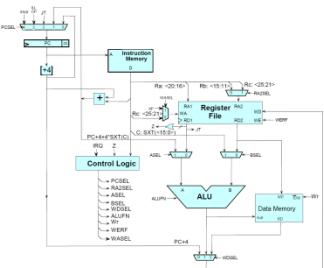
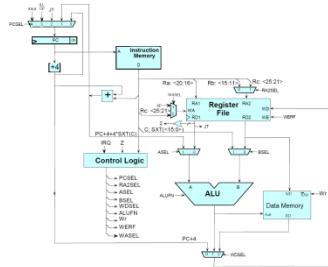
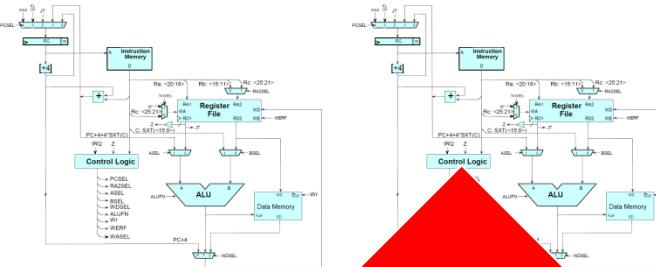
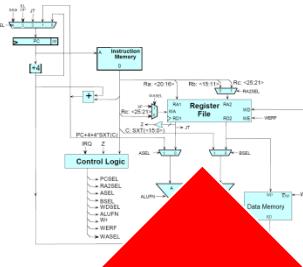
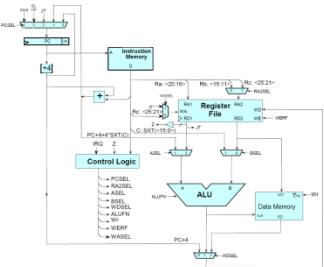
How long for 100GB?
What else can we do?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Get An Even Faster Machine with more Cores?



Assignment 10: Pipelined System Help

<https://assist.mit.edu/~mitpro/github.io/>

Add your name to assist_pro

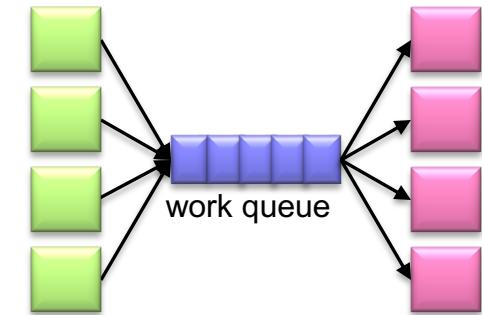
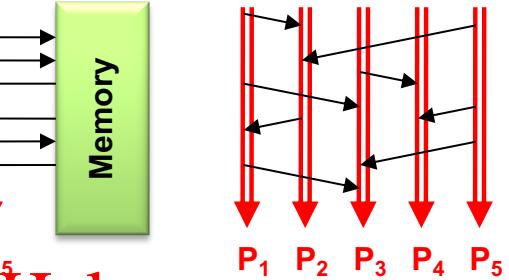
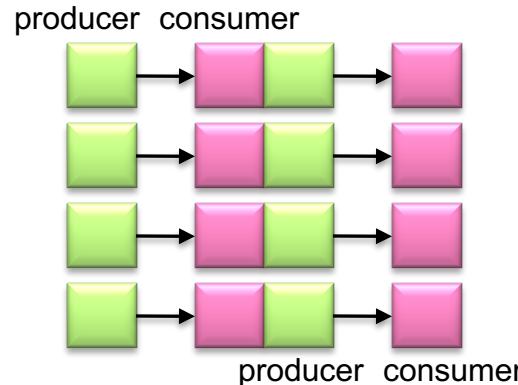
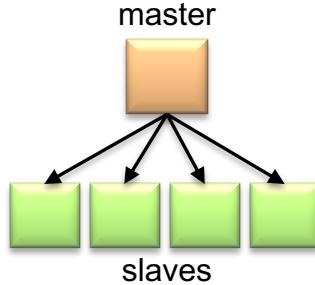
Current Tools

- Programming models
 - Shared memory (pthreads)
 - Message passing (MPI)
- Design Patterns

Assignment Project Exam Help

- Master-slaves
- Producer-cons <https://eduassistpro.github.io/>
- Shared work queues

Add WeChat edu_assist_pro



Where the rubber meets the road

- Concurrency is difficult to reason about
- Concurrency is even more difficult to reason about
 - At the scale of datacenters (even across datacenters)
 - In the presence of failures
 - In terms of mul
- Not to mention <https://eduassistpro.github.io/>
- The reality: [Add WeChat edu_assist_pro](#)
 - Lots of one-off solutions, custom code
 - Write you own dedicated library, then program with it
 - Burden on the programmer to explicitly manage everything

What's the common theme?

- To improve performance, you have to re-write the code
- The code has to adapt to the expected performance.
 - This doesn't work since you may not know the amount of data beforehand.
<https://eduassistpro.github.io/>
- The actual Intellectual Property
Add WeChat edu_assist_pro
e company is the analytic algorithm
 - However a lot of effort is spent on scaling the analytic

Big Data - Motivation

- Google processes 20 PB a day (2008)
- Wayback Machine has 3 PB + 100 TB/month (3/2009)
- Facebook has 2.5 PB of user data + 15 TB/day (4/2009)
- eBay has 6.5 PB o /2009)
- CERN's LHC will g

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
640K ought
to be enough
for anybody.



Enter .. Apache Hadoop

- Hadoop is a high-level Open Source project
 - Under Apache Software Foundation
 - Inspired by Google's MapReduce and GFS papers
- It contains several individual projects
 - HDFS
 - MapReduce
 - Yarn
- It also has a slew of related projects
 - PIG
 - HIVE
 - Hbase
- Has been implemented for the most part in Java.

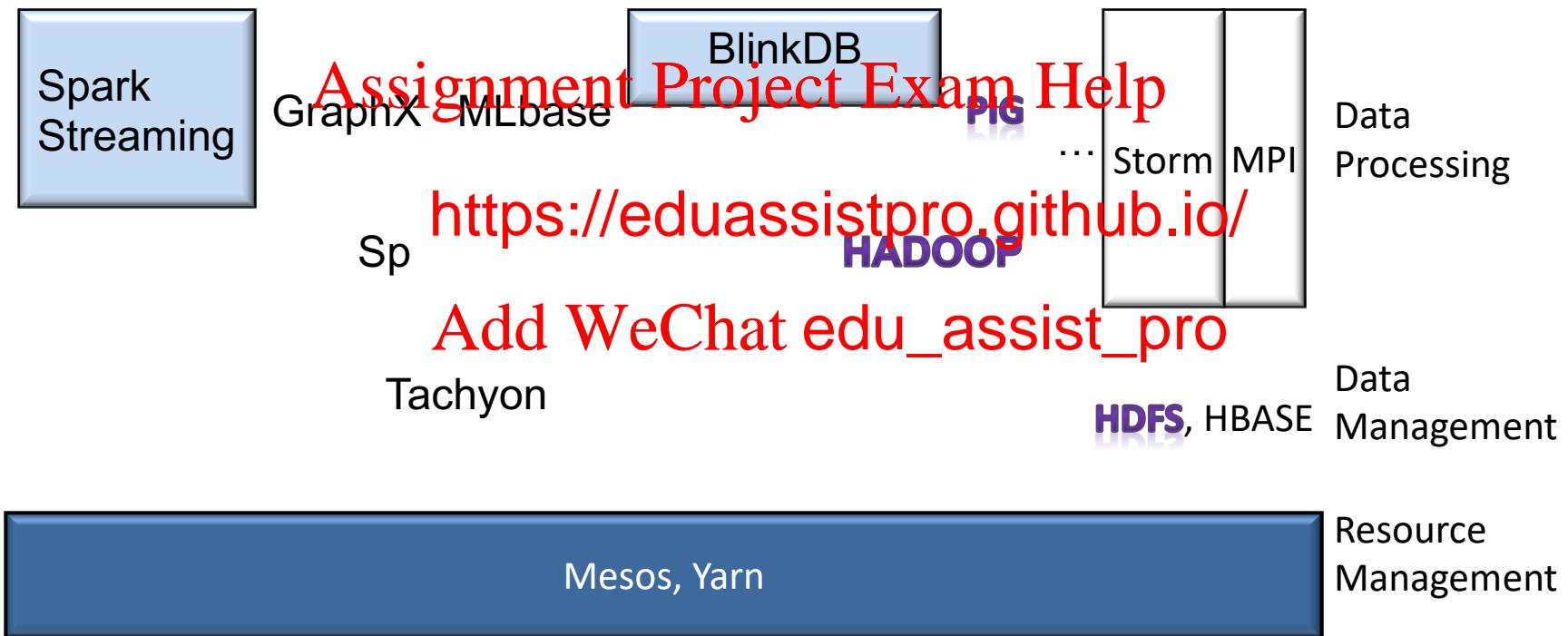
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Big Data System Architecture

- Current system modules for each software layer



A closer look

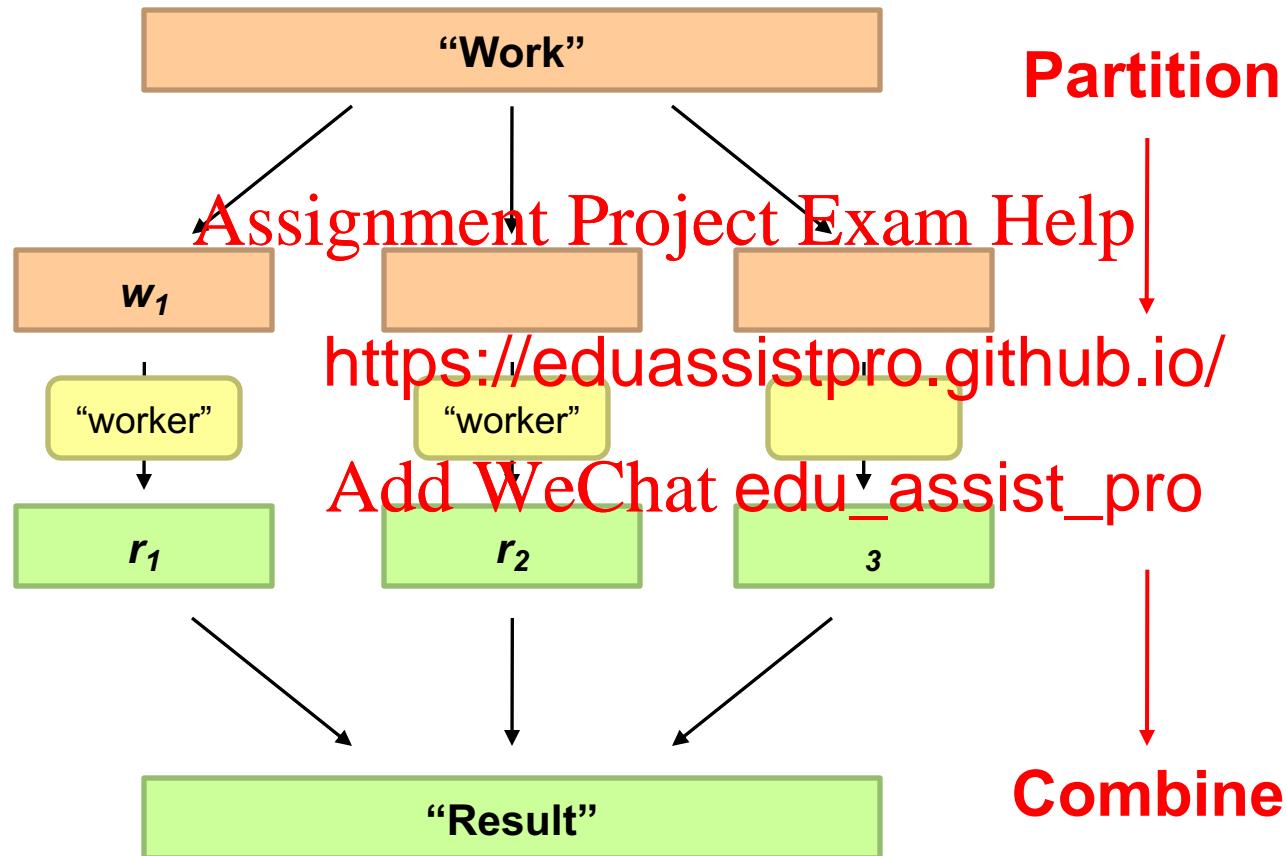
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assisttion_pro` Work

Combine Results

Divide and Conquer



Parallelization Challenges

- How do we assign work units to workers?
- What if we have more work units than workers?
- What if workers need to share partial results?
Assignment Project Exam Help
- How do we agg
- How do we know finished?
<https://eduassistpro.github.io/>
- What if workers die?
Add WeChat edu_assist_pro

What is the common theme of all of these problems?

What's the point?

- It's all about the right level of abstraction
 - The von Neumann architecture has served us well, but is no longer appropriate for the multi-core/cluster environment
- Hide system-level details from the developers
 - No more race conditions to consider
- Separating the <https://eduassistpro.github.io/>
 - Developer specifies the computation to be performed
 - Execution framework (“runtime”) for parallel execution

The datacenter *is* the computer!

“Big Ideas”

- Scale “out”, not “up”
 - Limits of SMP and large shared-memory machines
- Move processing to the data
 - Cluster have limited bandwidth
- Process data s
 - Seek access
- Seamless scalability
 - From the mythical man-month to the tradable machine-hour

Hadoop

- Platform for distributed **storage** and **computation**
 - HDFS
 - MapReduce
 - Ecosystem

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

What are we missing here?

Sequential File Read

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro Work

Combine Results

Hadoop

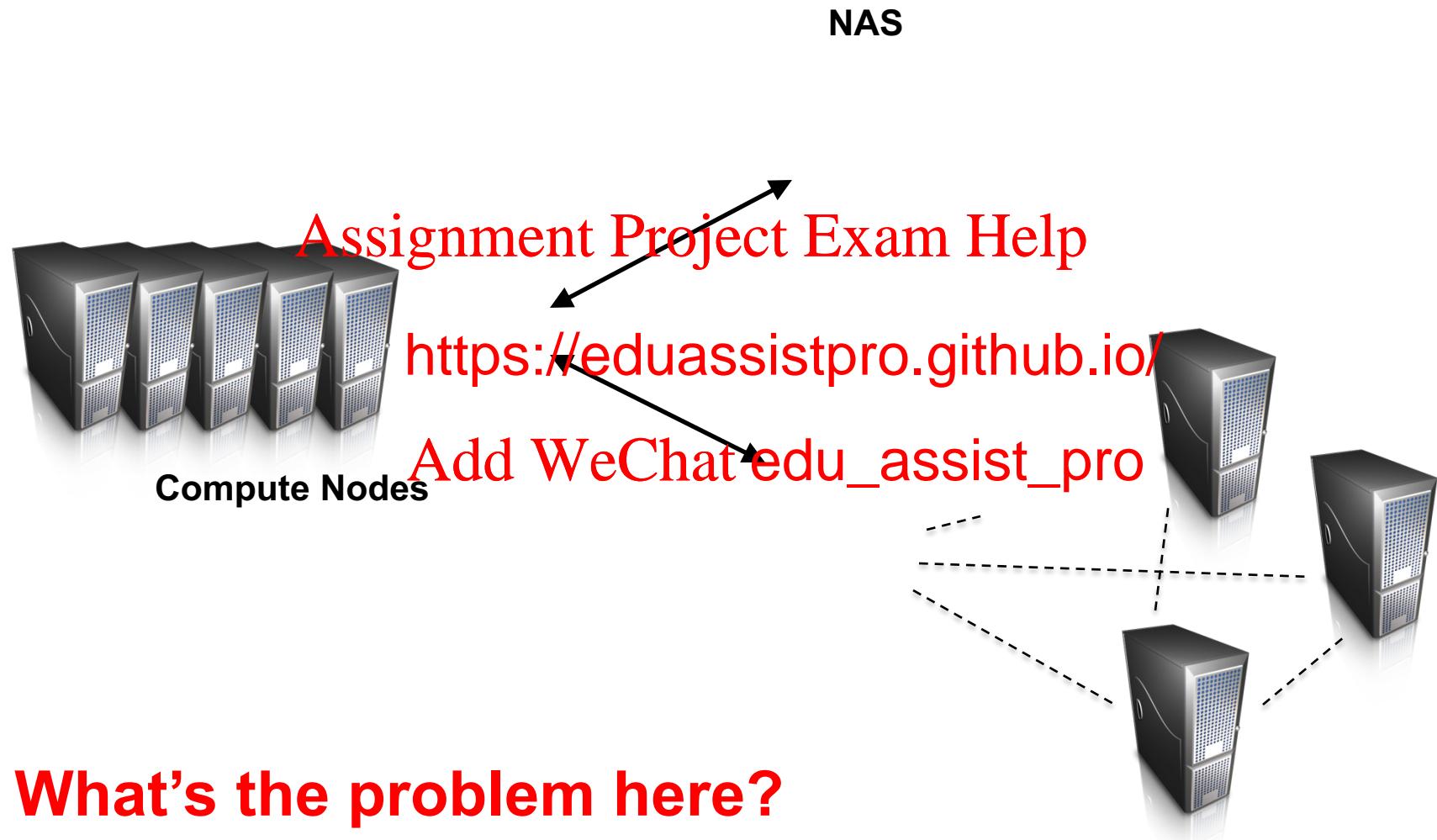
- Platform for distributed **storage** and **computation**
 - **HDFS**
 - **MapReduce**
 - **Ecosystem**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

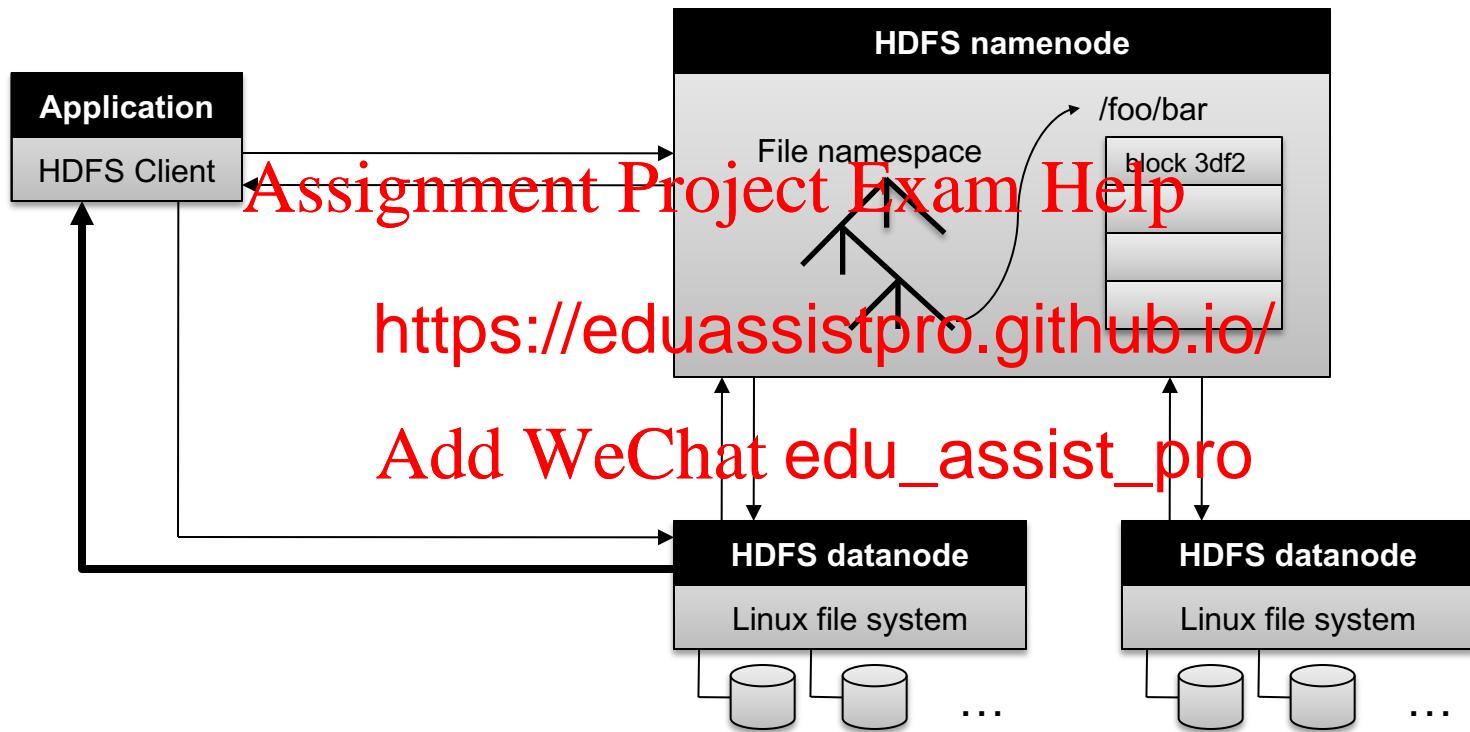
How do we get data to the workers?



HDFS: Assumptions

- Commodity hardware over “exotic” hardware
 - Scale “out”, not “up”
 - High component failure rates
 - Inexpensive commodity components fail all the time
 - “Modest” number of huge files
 - Multi-gigabyte files are
 - Files are write-once, read-many
 - Perhaps concurrently
 - Large streaming reads over random access
 - High sustained throughput over low latency
- Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

HDFS Architecture



How HDFS works

- When an input file is added to HDFS
 - File is split into smaller blocks of fixed size
 - Each block is replicated
 - Each replicated block is stored on a different host
- Block size is configurable. Default is 128/256MB.
- Replication level is c
 - Replication is nec
 - Scaling
 - High Availability
- In case a host crashes or is removed
 - All blocks on that host are automatically replicated to other hosts
- In case a host is added
 - Blocks will be rebalanced so that some blocks from other hosts will be placed on the new host

HDFS Component Responsibilities

- Name Node
 - Managing the file system namespace:
 - Holds file/directory structure, metadata, file-to-block mapping, access permissions, etc.
 - Coordinating file operations:
 - Directs clients to datanodes for reads and writes
 - No data is moved through the namenode
 - Maintaining overall health
 - Periodic communication
 - Block re-replication and repair
 - Garbage collection
- Data Node
 - Actual storage and management of data block on a single host
 - Provides clients with access to data

Assignment Project Exam Help

<https://eduassistpro.github.io/>

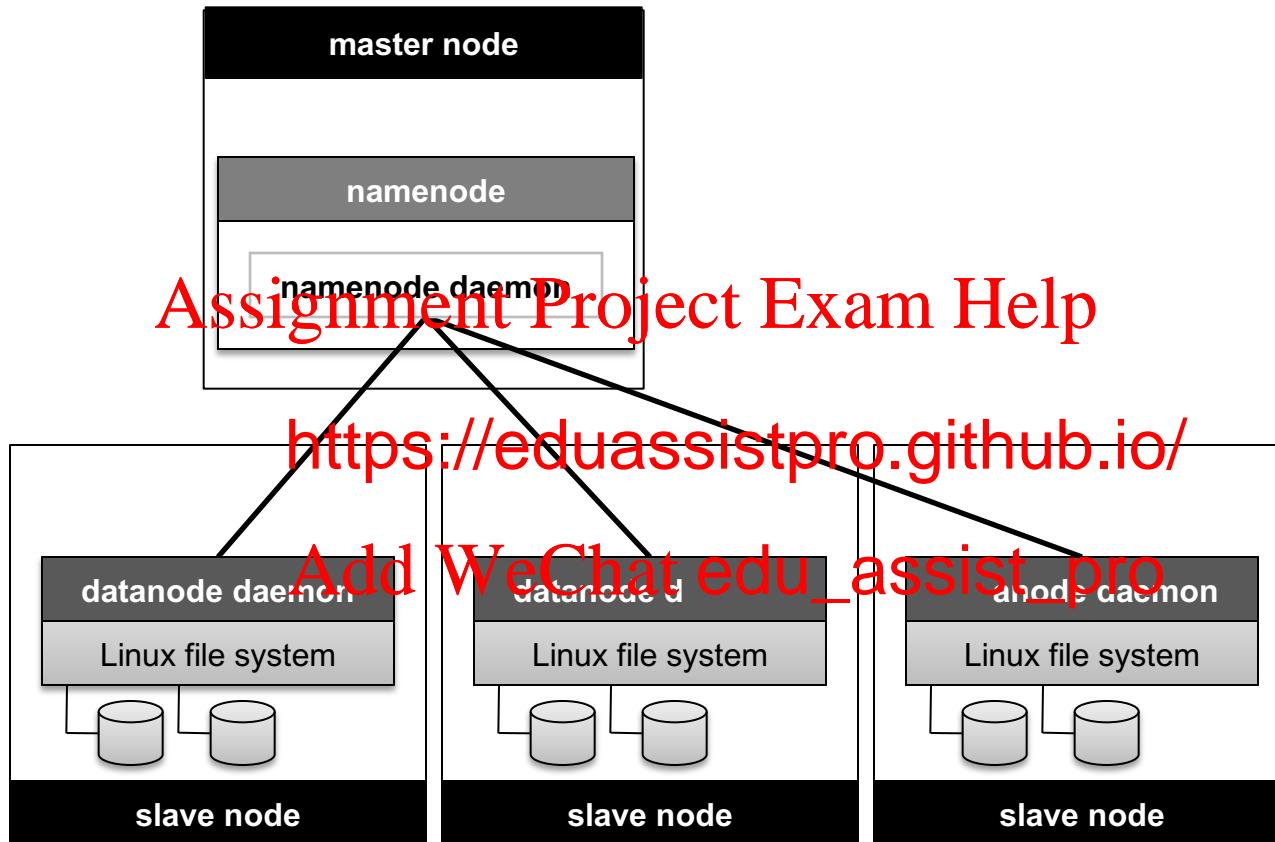
Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

HDFS Components in Cluster



Hadoop

- Platform for distributed **storage** and **computation**
 - HDFS
 - **MapReduce**
 - Ecosystem

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

MapReduce (MR) can refer to...

- The execution framework (aka “runtime”)
- The programming model
- The specific implementation

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Usage is usually clear from context!

MR Framework Components

- Job Tracker
 - Central component responsible for managing job lifecycles
 - One Job Tracker per MR framework instance
 - Accepts job submissions, queries etc. from clients
 - Enqueues jobs and schedules individual tasks.
 - Communicates with Task Trackers to assign and run tasks
 - Attempts to assign tasks to Task Trackers based on availability.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Task Tracker
 - One Task Tracker per host
 - Runs and manages individual tasks
 - Communicates progress of tasks back to Job Tracker.

MR Programming Model

- Programmers specify two functions:
 $\text{map } (k, v) \rightarrow \langle k', v' \rangle^*$
 $\text{reduce } (k', v') \rightarrow \langle k', v' \rangle^*$
 - All values with the same key are sent to the same reducer
- The MR Execution framework handles everything else...

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

MapReduce

- **Everything Else**

- Handles scheduling
 - Assigns workers to map and reduce tasks
- Handles “data distribution”
 - Moves processes to dat
- Handles synchronization
 - Gathers, sorts, and shuffles intermediate data

Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

- Handles errors and faults
 - Detects worker failures and restarts
- Everything happens on top of a distributed FS (HDFS)

Our Scoring Algorithm as a Map Reduce Program

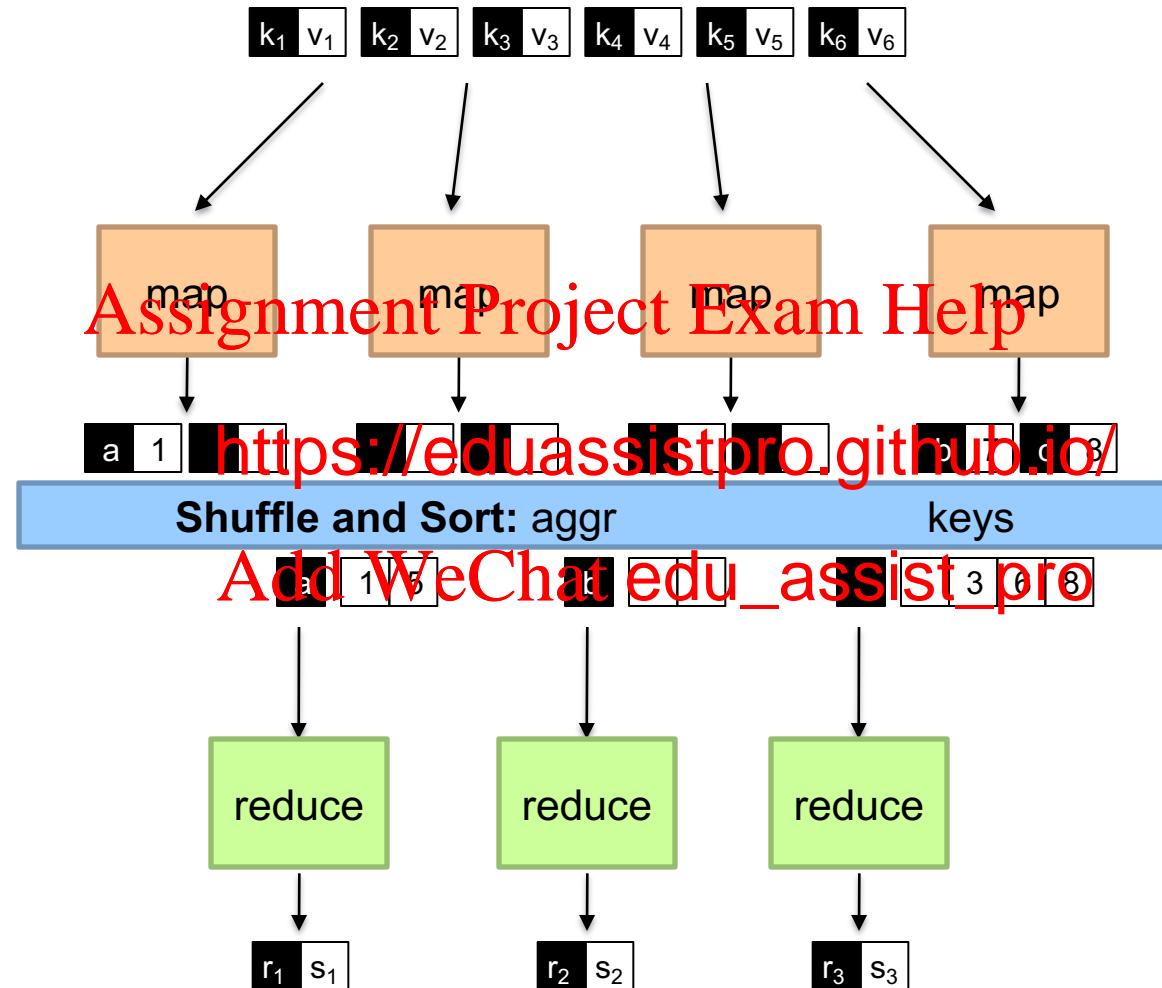
Assignment Project Exam Help

<https://eduassistpro.github.io/>

*Our Analytic
Add WeChat edu_assist_pro*

Basic Hadoop API*

- Mapper
 - void map(K1 key, V1 value, OutputCollector<K2, V2> output, Reporter reporter)
 - void configure(JobConf job)
 - void close() throws IOException
- Reducer/Comb <https://eduassistpro.github.io/>
 - void reduce(K2 key, Iterator<V2 OutputCollector<K3,V3> output, Reporter reporter)
 - void configure(JobConf job)
 - void close() throws IOException
- Partitioner
 - void getPartition(K2 key, V2 value, int numPartitions)



Lets Talk Numbers

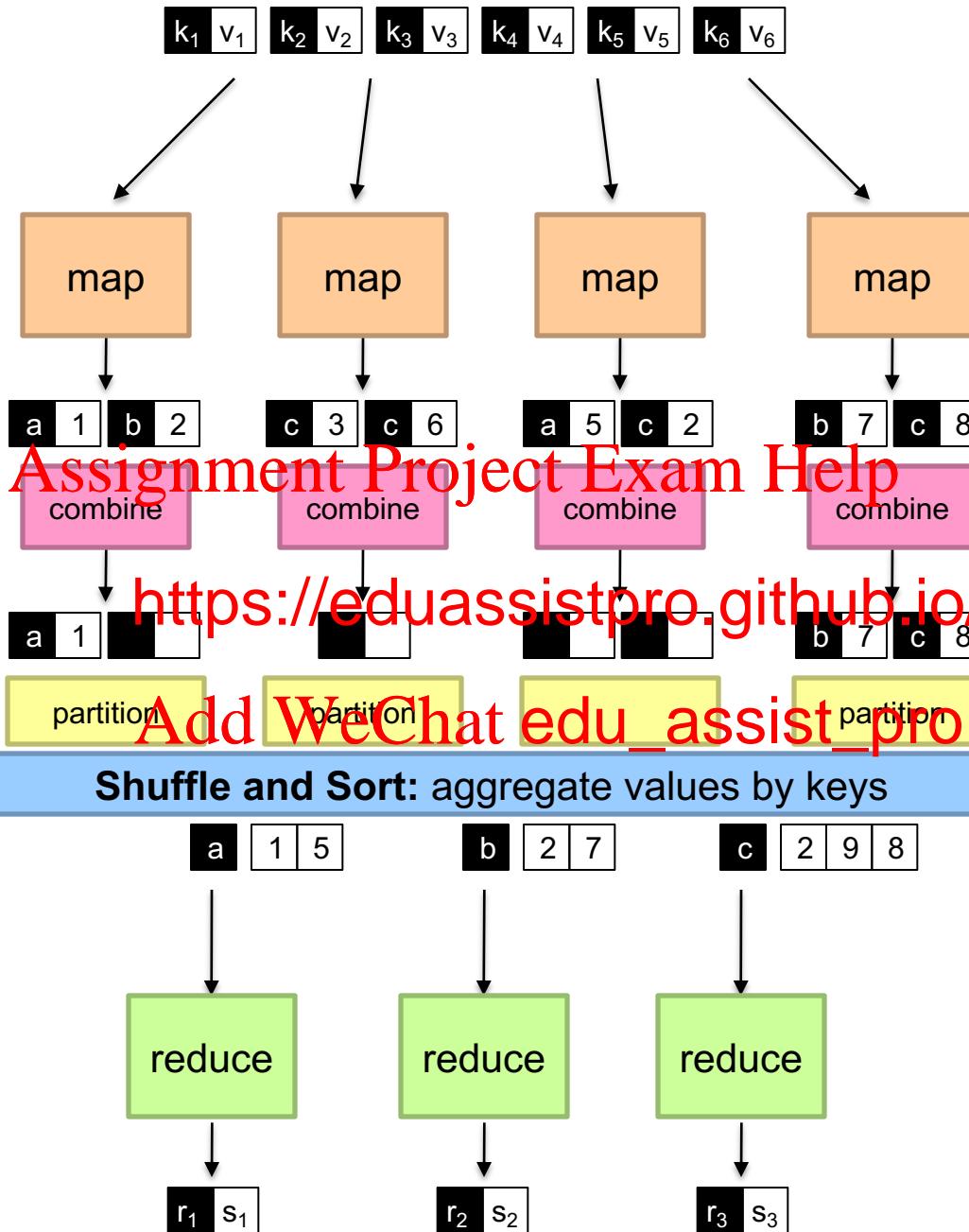
- How many mappers?
 - Depends on the size of input data
 - Typically 1 mapper per data block
 - So 1 GB input data will have around 8 Mappers
 - Assuming 128MB block size
- How many reducers?
 - Depends on cluster reducer capacity
 - Can be set depending on the number of keys
 - For large data sets, set it to cluster reducer capacity

MapReduce

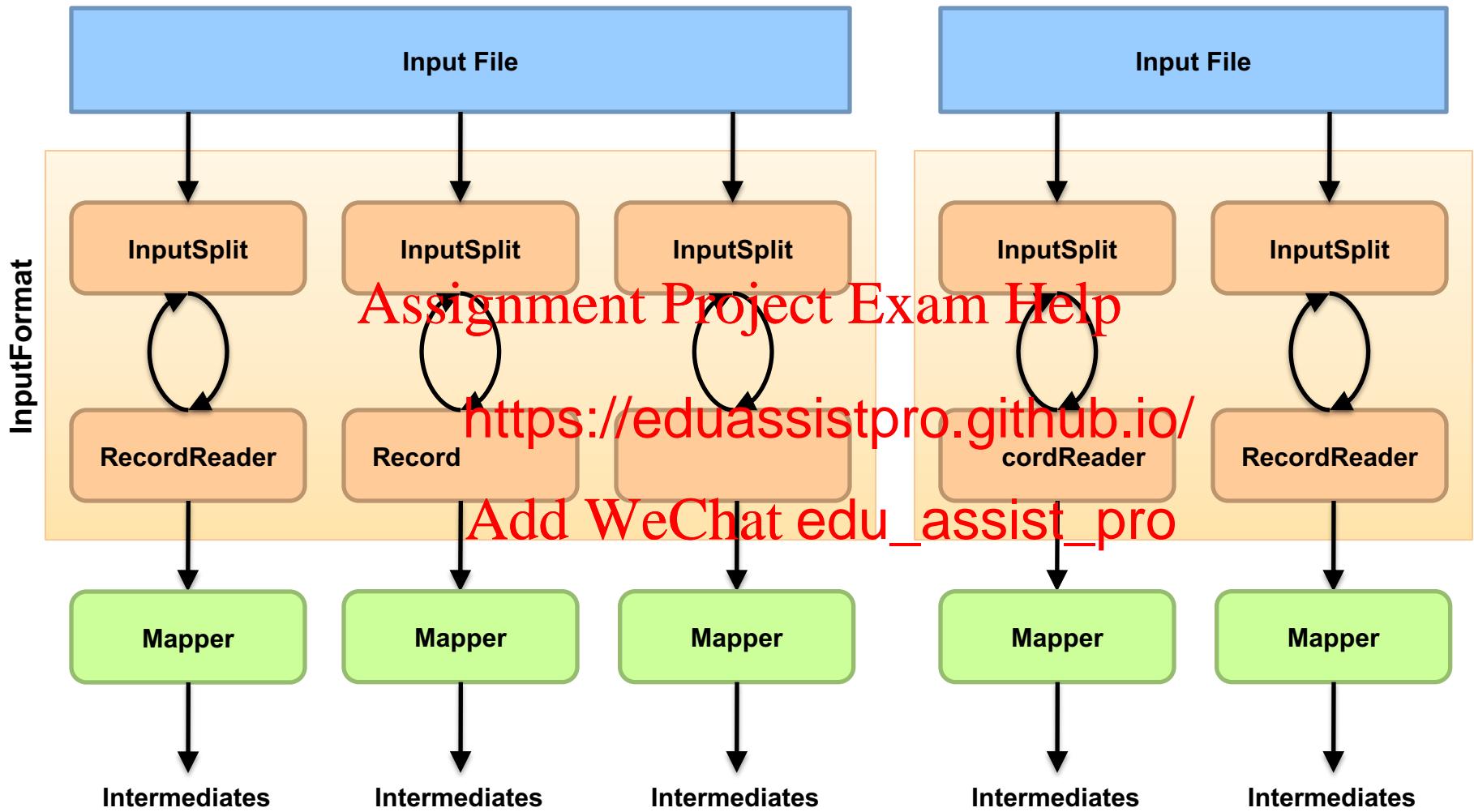
- Programmers specify two functions:
 - map** (k, v) $\rightarrow \langle k', v' \rangle^*$
 - reduce** (k', v') $\rightarrow \langle k', v' \rangle^*$
 - All values with the same key are reduced together
- The execution framework handles everything else...
- Not quite...usu <https://eduassistpro.github.io/>
 - combine** (k', v') $\rightarrow \langle k', v' \rangle^*$
 - Mini-reducers that run in memory during the reduce phase
 - Used as an optimization to reduce network traffic
 - partition** (k' , number of partitions) \rightarrow partition for k'
 - Often a simple hash of the key, e.g., $\text{hash}(k') \bmod n$
 - Divides up key space for parallel reduce operations

Two more details...

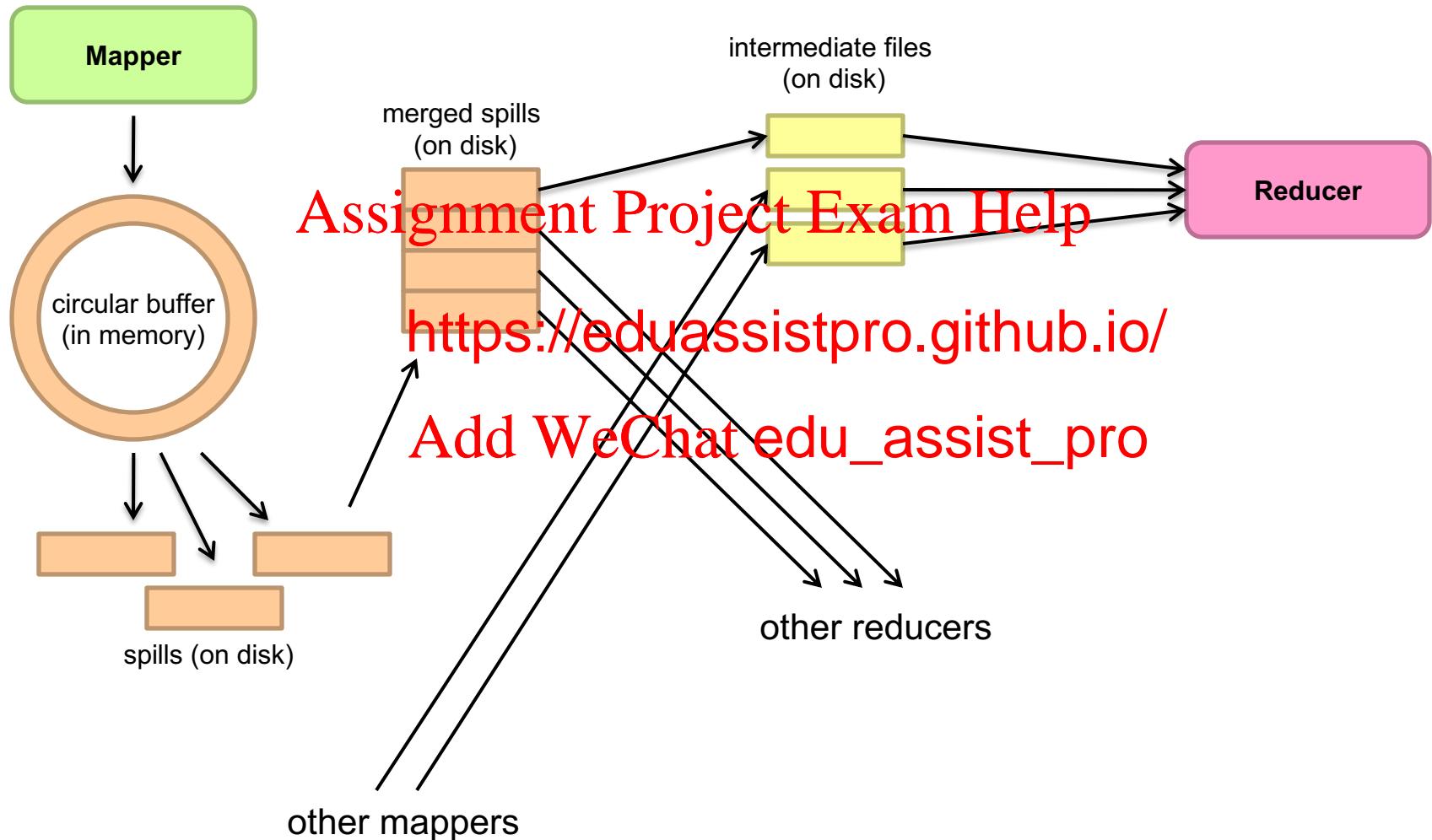
- Barrier between map and reduce phases
 - But we can begin copying intermediate data earlier
- Keys arrive at each reducer in sorted order
 - No enforced or
 - <https://eduassistpro.github.io/>
 - Add WeChat edu_assist_pro



Input To Mappers

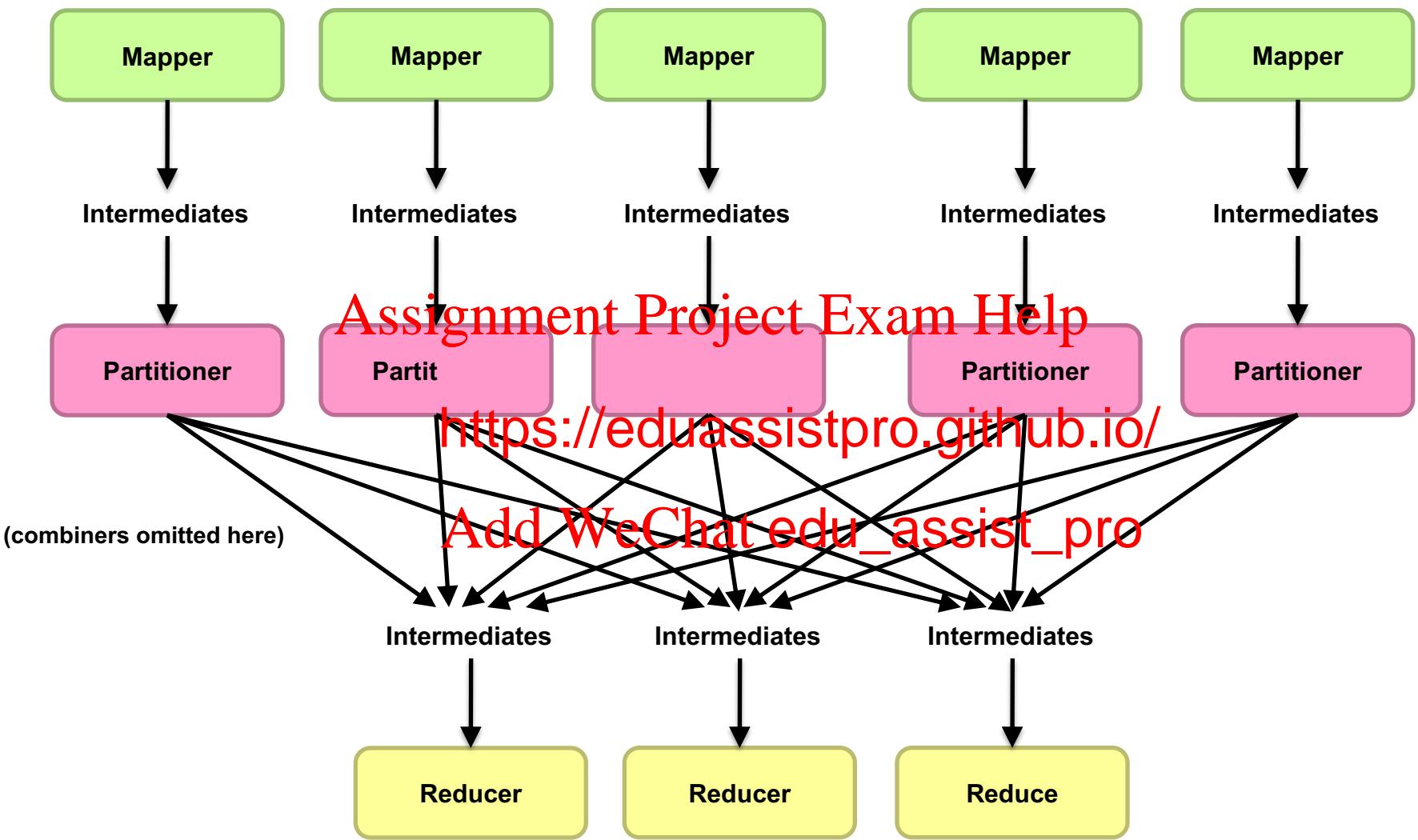


Shuffle and Sort

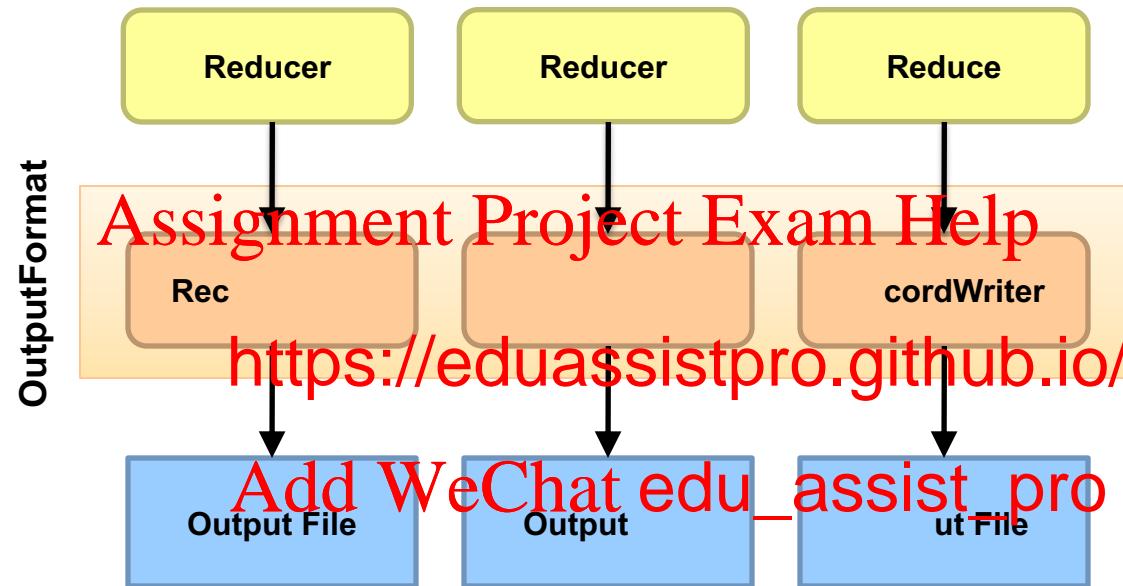


Shuffle and Sort in Hadoop

- Probably the most complex aspect of MapReduce!
- Map side
 - Map outputs are buffered in memory in a circular buffer
 - When buffer reaches threshold contents are “spilled” to disk
 - Spills merged into sorted within each partition): com <https://eduassistpro.github.io/>
- Reduce side
 - First, map outputs are copied over machine
 - “Sort” is a multi-pass merge of map outputs (happens in memory and on disk): combiner runs here
 - Final merge pass goes directly into reducer



Reducer to Output



Input and Output

- InputFormat:

- TextInputFormat
- KeyValueTextInputFormat
- SequenceFileInputFormat
- ...

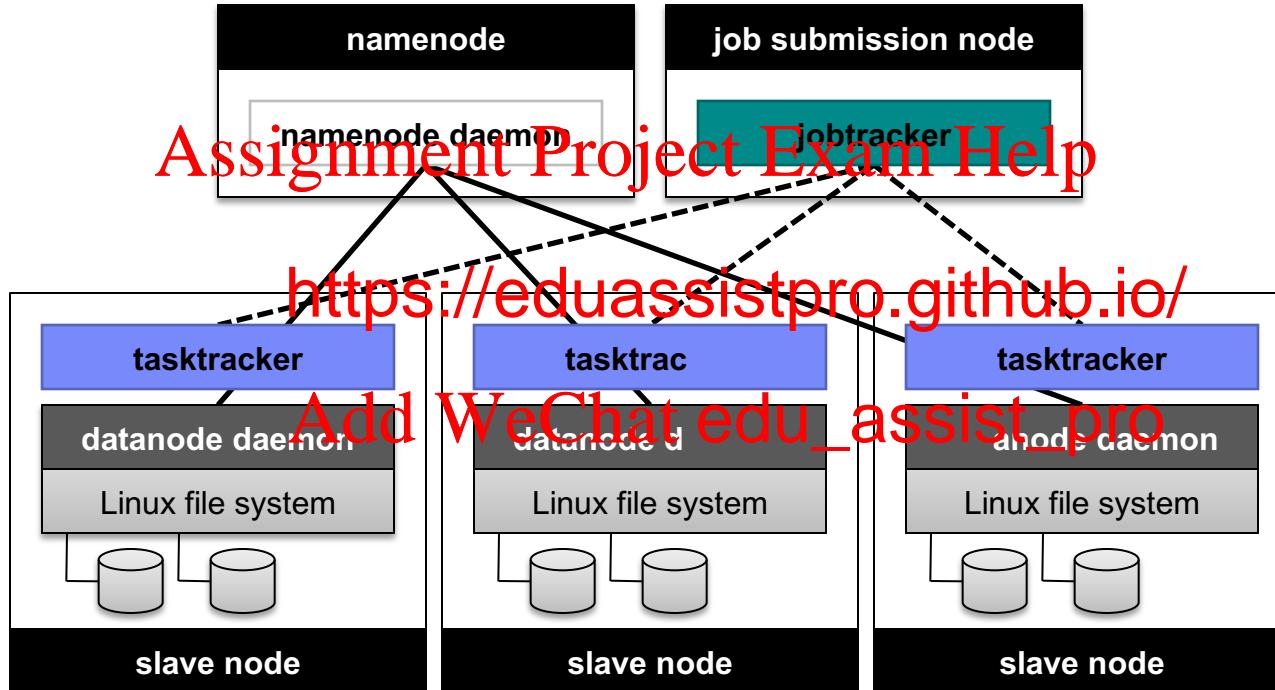
Assignment Project Exam Help

- OutputFormat: <https://eduassistpro.github.io/>

- TextOutputFormat
- SequenceFileOutputFormat
- ...

Add WeChat edu_assist_pro

Putting everything together...



HADOOP Architecture

- Master
 - NameNode
 - JobTracker
- Slaves
 - Data Node
 - Compute Node
 - Why together?
 - Data Locality

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

One More Thing

- Distributed Cache
 - Usually used for files of small size
 - Provides a convenient way to propagate applications and configuration files
 - HDFS is not used handle such files due to their small size
 - Shared across all <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dizzy Yet?

- OK, we went through a lot of details
- Whatever happened to the simplicity of programming??
Assignment Project Exam Help
- Do I really have to write **https://eduassistpro.github.io/** me I want to run a new analytic?

Add WeChat edu_assist_pro

We went from..

Multi-Threaded

Map-Reduce

Assignment Project Exam Help

<https://eduassistpro.github.io/>

 Add WeChat edu_assist_pro

Enter PIG ... Oink!

- High Level Languages for Map-Reduce
 - PIG
 - Developed by Yahoo
 - HIVE
 - Developed by Facebook
 - JAQL
 - Developed by IBM
- All of these languages p <https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat edu_assist_pro

- All of them allow users to plug in their own functions (UDFs)

Lets get Practical – From Setup to Results

Setting up a Hadoop Cluster

- Minimum recommended configuration (4 Hosts)
 - 1 Host Dedicated for Management Services (Job Tracker, Name Node etc)
 - 3 Hosts as Slave nodes (Data Node , Task Trackers)

Assignment Project Exam Help

- Data nodes should have
 - This is where all your data is stored
- How much total disk space?
 - Depends on input data to be processed
 - Effective Storage Space Recommended: Typically 3 times the size of your input data
 - Actual Storage Space: Effective Storage Space * 3 (replication level)
- Single node installation is fine for development/testing on very small data
 - Perhaps not the best for testing performance
- Installation instructions vary from provider to provider

Some cluster configuration parameters

- HDFS configuration parameters
 - Stored in hdfs-site.xml
 - Block size
 - Default replication count

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- MapReduce configuration
 - Stored In “mapred-site.xml”
 - Java heap size for mappers/reducers
 - Number of mappers/reducers per host
 - See <http://wiki.apache.org/hadoop/HowManyMapsAndReduces>

Add WeChat edu_assist_pro

- **IMPORTANT**
 - Job Tracker URL: http://<masterhost>:50030
 - Name Node URL: http://<masterhost>:50070

Job Tracker Web Page (port 50030)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Working with data

- Lets say you have 1 GB of data in your local filesystem (mydata.txt)
- Load into HDFS
 - hadoop fs –mkdir /path/mydirectory
 - hadoop fs –put mydata.txt /path/mydirectory
 - where /path/mydirectory is in HDFS
- List the file you just uplo <https://eduassistpro.github.io/>
 - hadoop fs –ls /path/mydirectory
- “hadoop fs” works similar to linux filesystem commands
 - However HDFS is not POSIX compliant.
 - It cannot be mounted as a regular filesystem

Writing your program .. see the simplicity!!

- JAQL program for running our scorer

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- PIG program for running our scorer

Add WeChat edu_assist_pro

All languages provide similar functionality

- LOAD (various data formats)
- JOIN
- FOR-EACH
- GROUP
- SORT
- FILTER
- Pluggable UDFs

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Hadoop Programming Tips

- Thinking at scale
 - Filter unwanted data earlier in the flow
 - Store intermediate data
 - Use “sequence” format for storing data.

Assignment Project Exam Help

- These are not iterative
 - i.e. No *for* or *while* loops
- Watch out for obvious bottlenecks
 - Single key for all mapper output will send all data to one reducer
 - Too much data sent to a UDF will result in OOM errors

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

Submitting a Job

- Create and save your PIG script (myscript.pig)
- To deploy (pig command will be in your installation)
 - pig –f myscript.pig
 - Command will complete once your job completes

Assignment Project Exam Help

- To check the status of y
 - Use the Job Tracker U <https://eduassistpro.github.io/>
 - hadoop job –list (will print all job ids)
 - hadoop job –status <jobid> (will print the job)
- To get the results
 - hadoop fs –get /path/results.txt .

Add WeChat edu_assist_pro

Anatomy of a Job

- MapReduce program in Hadoop = Hadoop job
 - Jobs are divided into map and reduce tasks
 - An instance of running a task is called a task attempt
 - Multiple jobs can be composed into a workflow
- Job submission
 - Client (i.e., driver) configures it, and submits it to job tracker
 - JobClient computes input splits ()
 - Job data (jar, configuration XML) are sent to JobTracker
 - JobTracker puts job data in shared location, enqueues tasks
 - TaskTrackers poll for tasks
 - Off to the races...

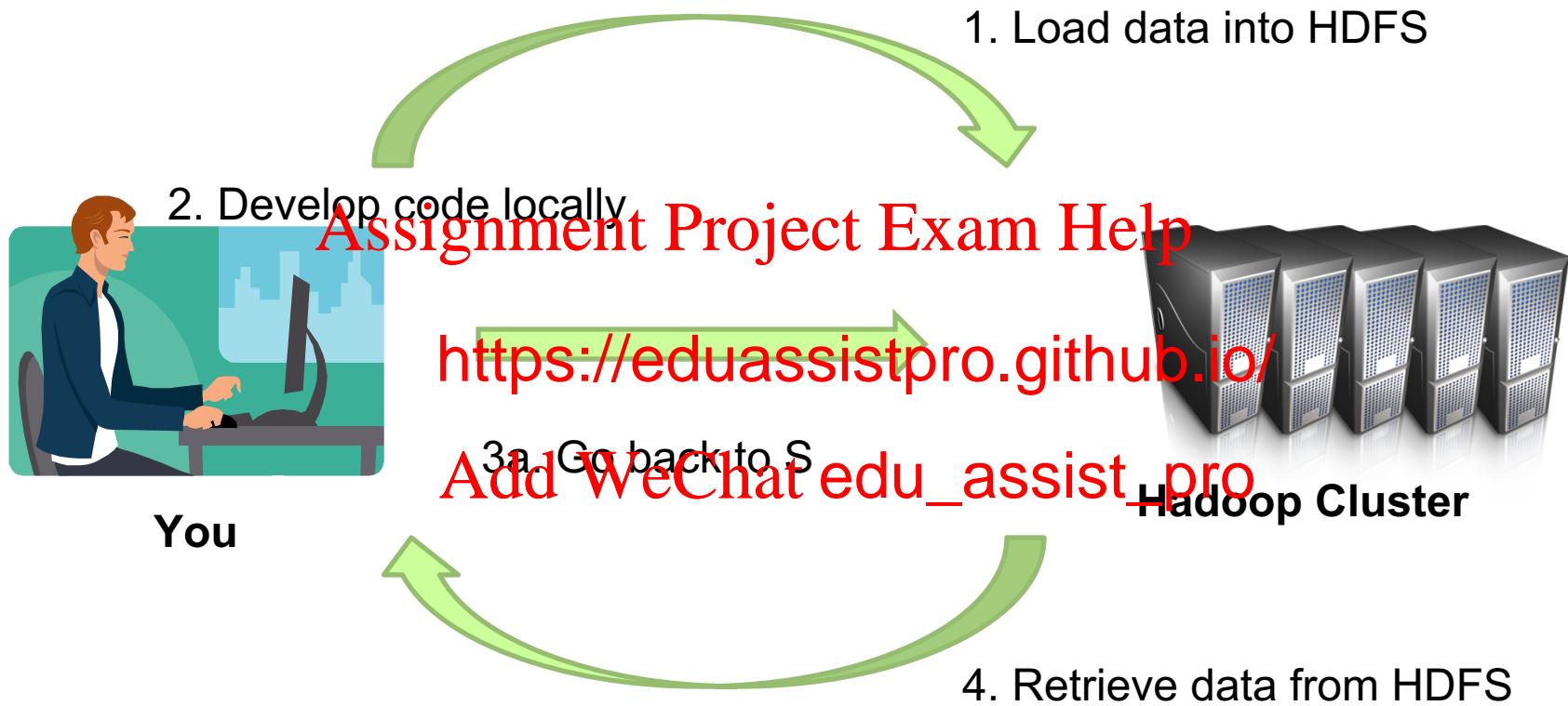
A simple illustration of MR process

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Hadoop Workflow



Larger View...

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

MR Patterns Examples

- Jimmy Lin's book
- Jeffrey Ulman's book
- An excellent blog: <https://highlyscalable.wordpress.com/2012/02/01/mapreduce-patterns/>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Uh Oh.. My Job Failed...Now what?

- First, take a deep breath
- Start small, start locally
- Strategies
 - Learn to use the webapp
 - Where does pr <https://eduassistpro.github.io/>
 - Don't use println
 - Throw RuntimeExceptions
- Logs are most easily accessible via the Job Tracker URL

Time for a Raise

- Finally you have mastered Hadoop Big Data
- Your applications are scaling.
 - You deserve a raise!!
- Boss **Assignment Project Exam Help**
 - Can we query t <https://eduassistpro.github.io/>
 - How long will t
- Problem **Add WeChat edu_assist_pro**
 - Remember this is still sequential access
 - To find a specific entity, you still need to read the entire data set.
- What now?
 - How is this solved in traditional systems?

Databases

Other projects based on Hadoop

- HBase
- Hive
- PIG
- Spark
- Mahout

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Hive – a SQL-like data warehouse on Hadoop

<https://cwiki.apache.org/confluence/display/Hive/Tutorial>

- Supports a SQL-like data warehouse on top of Hadoop
 - began at Facebook
- Provides SQL users the capability of big data without requiring lower level programming for a wide range of tasks
- Fewer lines of code!
- /bin/hive –help

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

```
hive> LOAD DATA LOCAL INPATH 'cite75_99.txt'
      > OVERWRITE INTO TABLE cite;
Copying data from file:/root/cite75_99.txt
Loading data to table cite
OK
Time taken: 9.51 seconds
```

Wordcount Example

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Hive – SQL-like data warehouse on top of Hadoop

- A data warehouse with SQL like (HiveQL) interface on Hadoop for processing and managing structured data
- Hive hides the complexity of MR programming details from user with interest for processing structured data

Assignment Project Exam Help

<https://eduassistpro.github.io/>

INSERT OVERWRITE TABLE user_active

SELECT user.*

Add WeChat edu_assist_pro

FROM user

WHERE user.active =1;

- Key components

- Metastore
- Parser/planner/optimizer
- Interface

Hive Installation and config

- Just download recent Hive tarball and extract it.
- Need to set up few directories
 - bin/hadoop fs -mkdir /tmp
 - bin/hadoop fs -mkdir /user/hive/warehouse
- Hive manages all the data under this directory
- Hive stores metadata at <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Hive examples

- 11.1.2

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Data model – column partition

- Tables as fundamental data model – stored under /usr/hive/warehouse
- Hive uses concept of partition columns – partitions data by column index, i.e., state column, date column partition..

/user/hive/warehouse/users/date=20090901/state=CA

/user/hive/warehouse/users/date=20090901/state=NY

/user/hive/warehouse/users/date=20090901/state=TX ...

Assignment Project Exam Help

/user/hive/warehouse/users/d

/user/hive/warehouse/users/d

<https://eduassistpro.github.io/>

/user/hive/warehouse/users/d

/user/hive/warehouse/users/date=20090903/state=CA

/user/hive/warehouse/users/date=20090903/state=NY

/user/hive/warehouse/users/date=20090903/state=TX ...

Add WeChat edu_assist_pro

Hive Data Model – partition and cluster

- Tables stored under user/hive/warehouse in HDFS
- Partition columns

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Buckets – allows to create smaller range partitions

Add WeChat edu_assist_pro

```
CREATE TABLE page_view(viewTime INT, userid BIGINT,
                      page_url STRING, referrer_url STRING,
                      ip STRING COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY (dt STRING, country STRING)
CLUSTERED BY (userid) INTO 32 BUCKETS
ROW FORMAT DELIMITED
  FIELDS TERMINATED BY '\t'
  LINES TERMINATED BY '\n'
STORED AS SEQUENCEFILE;
```

Data model - buckets

- concept of *buckets*, which provide efficiency to queries that can work well on a random sample of data
- Bucketing divides data into a specified number of files based on the hash of the bucket column

Assignment Project Exam Help

/user/hive/warehouse/users/date=20090901	state=CA/part-00000
/user/hive/warehouse	00031
/user/hive/warehouse	https://eduassistpro.github.io/00000...
/user/hive/warehouse/users/date=20090901	31
/user/hive/warehouse/users/date=20090901	30

Add WeChat edu_assist_pro

- Example - pageview

```
CREATE TABLE page_view(viewTime INT, userid BIGINT,
                      page_url STRING, referrer_url STRING,
                      ip STRING COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY (dt STRING, country STRING)
CLUSTERED BY (userid) INTO 32 BUCKETS
ROW FORMAT DELIMITED
  FIELDS TERMINATED BY '\t'
  LINES TERMINATED BY '\n'
STORED AS SEQUENCEFILE;
```

Resources

- Papers
 - Google File System, 2003
 - Google MapReduce, 2004
 - Google Bigtable, 2006
- URLs
 - Apache Hadoop <https://eduassistpro.github.io/>
- Available Hadoop Distribution
 - Apache, IBM, Cloudera, Hortonworks

Hive Architecture

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Main components
 - SQL interface
 - Parser/Planner
 - Metastore
 - Driver