

COMP9313: Big Data Management

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Lecturer: Xin Cao

Course web site: <http://www.cse.unsw.edu.au/~cs9313/>

Assignment Project Exam Help
Chapter <https://eduassistpro.github.io/>
Introduction <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Course Info

- Lectures : 6 : 00 – 9:00 pm (Tuesday)
- Location:
 - Old Main Building 230 (K-K15-230)
 - Webstream
- Labs: Weeks 2-10
- Consultation (Wednesday) *lectures, course materials, assig* <https://eduassistpro.github.io/>
 - Time: 3:00 – 4:00 pm (Tuesday)
 - Place: 201D, K-17 **Add WeChat edu_assist_pro**
- TA:
 - Xuefeng Chen, *xuefeng.chen@student.unsw.edu.au*
- Tutors: Xuefeng Chen, Wei Li, You Peng, Yu Hao
- Discussion and QA: WebCMS3

Lecturer in Charge

- Lecturer: Xin Cao
 - Office: 201D K17 (outside the lift turn left)
 - Email: *xin.cao@unsw.edu.au*
 - Ext: 55932

Assignment Project Exam Help

- Research interests
 - Database <https://eduassistpro.github.io/>
 - Data Mining
 - Big Data Technologies
 - My homepage: <http://www.cse.unsw.edu.au/~z3515164/>
 - My publications list at google scholar:
<https://scholar.google.com.au/citations?user=kJlkUagAAAAJ&hl=en>

Course Aims

- This course aims to introduce you to the concepts behind Big Data, the core technologies used in managing large-scale data sets, and a range of technologies for developing solutions to large-scale data analytics problems.

Assignment Project Exam Help

- This course is in large-scale data and technologies, and will prepare systems as well as use them efficiently to address challenges in big data management.
 - understand modern range of topics
 - able to build such systems
 - rely to address challenges in big data management.
- *Not possible to cover every aspect of big data management.*

Lectures

- Lectures focusing on the frontier technologies on big data management and the typical applications
- Try to run in more interactive mode and provide more examples

Assignment Project Exam Help

- A few lectures
r (e.g., like a
lab/demo) to co <https://eduassistpro.github.io/>
- Lecture length varies slightly depending on progress (of that
lecture) □
- Note: attendance to every lecture is assumed

Resources

- Text Books
 - [Hadoop: The Definitive Guide](#). Tom White. 4th Edition - O'Reilly Media
 - [Mining of Massive Datasets](#). Jure Leskovec, Anand Rajaraman, Jeff Ullman. 2nd edition - Cambridge University Press
 - [Data-Intensive Text Processing with MapReduce](#). Jimmy Lin and Chris Dyer.
 - [Learning Sp](#) Park, au, Andy Konwinski, Patrick Wendell. O'Reilly Media
- Reference Books and other reading
 - [Apache MapReduce Tutorial](#)
 - [Apache Spark Quick Start](#)
 - Many other online tutorials
- Big Data is a relatively new topic (so no fixed syllabus)

Prerequisite

- Official prerequisite of this course is COMP9024 (Data Structures and Algorithms) and COMP9311 (Database Systems).
- Before commencing this course, you should:
 - have experiences and good knowledge of algorithm design (equivalent to COMP9024)
 - have a solid background in database systems (equivalent to COMP9311)
 - **have solid** <https://eduassistpro.github.io/>
 - **be familiar with** working on a **erating systems**
vector spaces, matrix multiplication), probability theory and statistics , and graph theory
- No previous experience necessary in
 - MapReduce/Spark
 - Parallel and distributed programming

Please do not enrol if you

- Don't have COMP9024/9311 knowledge
- Cannot produce correct Java program on your own
- Never worked on Unix-style operating systems
- Have poor time management
- Are too busy to attend lectures/labs

Assignment Project Exam Help

- *Otherwise, you* <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Learning outcomes

- After completing this course, you are expected to:
 - elaborate the important characteristics of Big Data
 - develop an appropriate storage structure for a Big Data repository
 - utilize the map/reduce paradigm and the to manipulate Big Data
 - utilize the Spark platform to manipulate Big Data
 - develop effi
Data blems involving Big

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Assessment

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Coding Projects

- Projects:
 - 1 warm-up programming project on Hadoop MapReduce
 - 1 harder project on Hadoop MapReduce
 - 1 project on Spark
 - 1 project on AWS (MapReduce/Spark)
- Both results and <https://eduassistpro.github.io/>
 - If not able to run your codes due to any reason, you will not lose all marks. Add WeChat **edu_assist_pro**

CSE Computing Environment

- Use Linux/command line (virtual machine image will be provided)
 - Projects marked on Linux servers
 - You need to be able to upload, run, and test your program under Linux

Assignment Project Exam Help

- Assignment sub
 - Use Give to <https://eduassistpro.github.io/> (Web page)
 - Classrun. Check your submission
https://wiki.cse.iitkgp.ac.in/Add_WeChat_edu_assist_pro Read

Final exam

- Final written exam (100 pts)
- If you are ill on the day of the exam, do not attend the exam – I will not accept any medical special consideration claims from people who already attempted the exam
- You need to add WeChat https://eduassistpro.github.io/the_final_exam
- No supplementary exam will be given

You May Fail Because ...

- *Plagiarism*
- Code failed to compile due to some mistakes
- Late submission
 - 1 sec late = 1 day late
 - submit wrong files
- Program did not
 - <https://eduassistpro.github.io/>
- I am unlikely to accept the following
 - “Too busy”
 - “It took longer than I thought it would take”
 - “It was harder than I initially thought”
 -

Tentative Course Schedule

Week	Topic	Assignment
1	Course info and introduction to big data	
2	Hadoop MapReduce 1	
3	Hadoop MapReduce 2	Proj1
4	Hadoop MapReduce 3 Assignment Project Exam Help	
5	Graph	Proj2
6	Spark https://eduassistpro.github.io/	
7	Spark 2 Add WeChat edu_assist_pro	Proj3
8	Data stream mining	
9	Finding Similar Items	Proj4
10	Recommender Systems	
11	NoSQL and High Level MapReduce Tools	
12	Revision and exam preparation	

Labs

- 5 labs on MapReduce
- 3 labs on Spark
- 1 lab on higher level MapReduce tools
- 1 lab on AWS <https://eduassistpro.github.io/>
- 1 lab on big data machine learning
[Add WeChat edu_assist_pro
ive]

Virtual Machine

- Software: Virtualbox

- Images:

- Pure Xubuntu 14.04:

- http://www.cse.unsw.edu.au/~z3515164/Raw_Xubuntu.zip

- Xubuntu 14.04 with pre-installed Hadoop and Eclipse plugin:
<http://mirror.cse.unsw.edu.au/pub/cs9313/Xubuntu.zip>

- Download "xubunusit" and rename the file "xubunusit-disk2.vmdk"
 - Open VirtualBox, File->Add **WeChat** **edu_assist_pro** **vf**" file
 - Browse the image folder,
 - The image will be imported to your computer, which may take 10 minutes
 - comp9313 is used as both username and password. The hadoop installation path is the same as in the virtual machine on lab computers.

Your Feedbacks Are Important

- Big data is a new topic, and thus the course is tentative
- The technologies keep evolving, and the course materials need to be updated correspondingly

Assignment Project Exam Help

- Please advise w
discussio
cturer, at the
<https://eduassistpro.github.io/>
- myExperience system
Add WeChat edu_assist_pro

Why Attend the Lectures?

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Assignment Project Exam Help

Part 2: <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

What is Big Data?

- No standard definition! here is from Wikipedia:
 - Big data is a term for data sets that are so voluminous or complex that traditional data processing application software are inadequate to deal with them
 - Challenges include capture, storage, analysis, data curation, search, sharing, transfer, visualization, querying, updating and information
 - The term "big data" refers to the use of predictive analytics, user behaviour analysis, and other advanced data analytics methods that extract value from a particular size of data set
 - Analysis of data sets can find new correlations to "spot business trends, prevent diseases, combat crime and so on."

Dan Ariely 
January 7, 2013 · 

Follow 

Big data is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it...

Instead of Talking about “Big Data”...

- Let's talk about a crowded application ecosystem:
 - Hadoop MapReduce
 - Spark
 - NoSQL (e.g., HBase, MongoDB, Neo4j)
 - Pregel
 - Assignment Project Exam Help
 -
- Let's talk about data science and d
nt:
Add WeChat edu_assist_pro
 - Finding similar items
 - Graph data processing
 - Streaming data processing
 - Machine learning technologies
 -

Who is generating Big Data?

Social

User Tracking & Engagement

Homeland Security



Assignment Project Exam Help

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

Add WeChat edu_assist_pro

eCommerce

Financial Serv

Real Time Search



Google



Big Data Characteristics: 3V

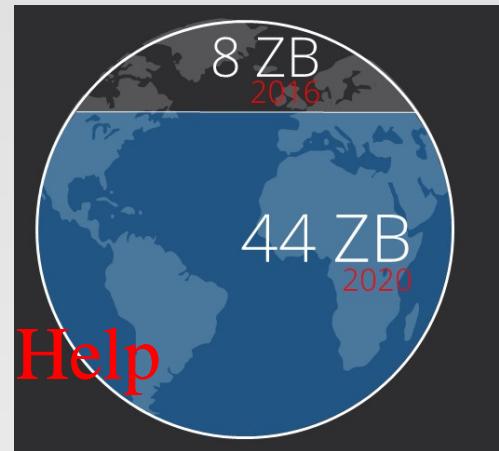
Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Volume (Scale)

- Data Volume
 - Growth 40% per year
 - From 8 zettabytes (2016) to 44zb (2020)
- Data volume is increasing exponentially



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

Add WeChat edu_assist_pro

Number of Tweets

Recent Twitter Statistics

Assignment Project Exam Help

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

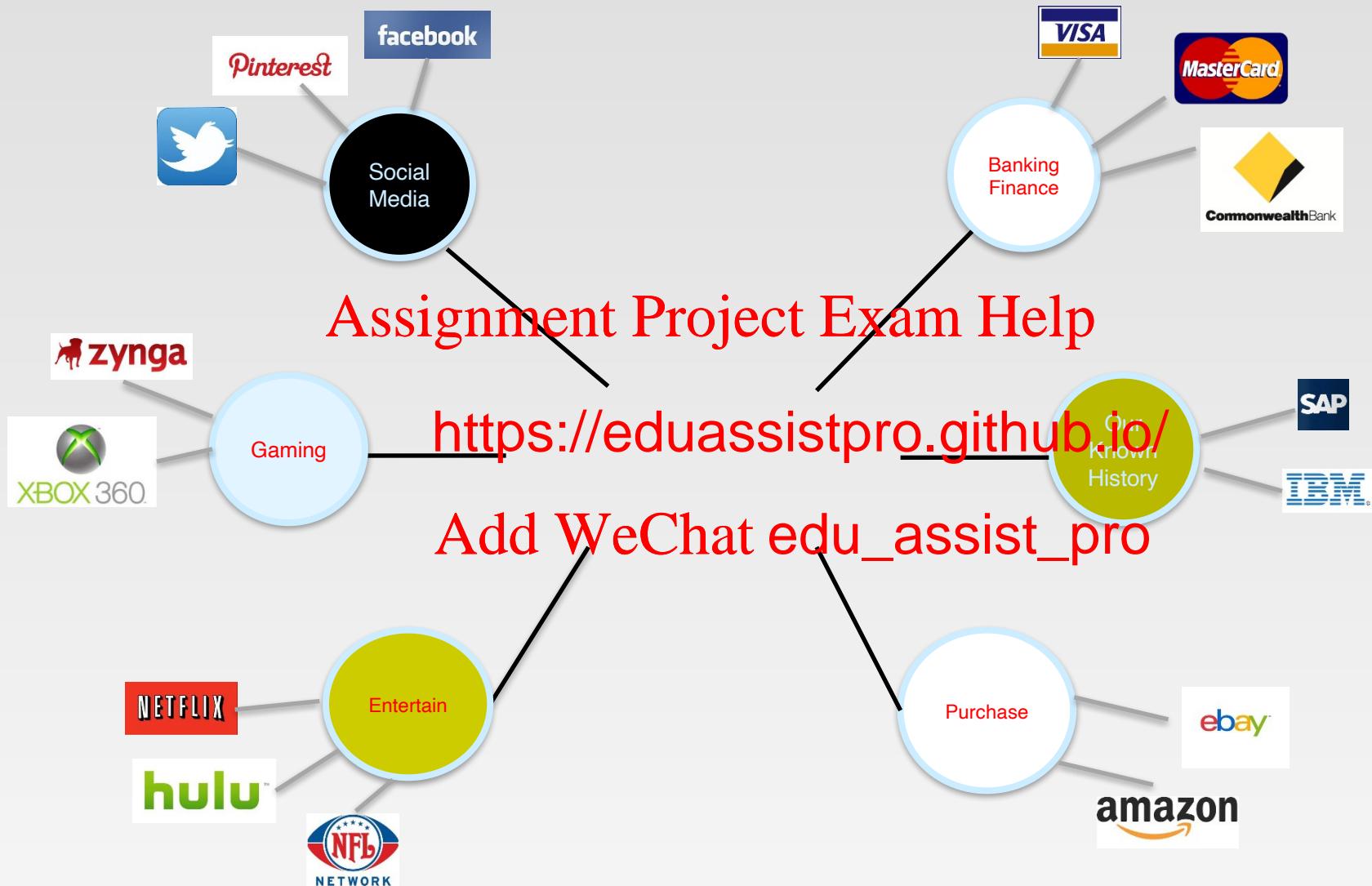
Add WeChat edu_assist_pro

Variety (Complexity)

- Different Types:
 - Relational Data (Tables/Transaction/Legacy Data)
 - Text Data (Web)
 - Semi-structured Data (XML)
 - Graph Data
 - ▶ Social Networks
 - Streaming Data
 - ▶ You can only scan the data
 - A single application can be generating many types of data
- Different Sources :
 - Movie reviews from IMDB and Rotten Tomatoes
 - Product reviews from different provider websites

To extract knowledge → all these types of data need to be linked together

A Single View to the Customer



A Global View of Linked Big Data

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



Diversified social network

Velocity (Speed)

- Data is being generated fast and need to be processed fast
- Online Data Analytics
- Late decisions → missing opportunities
- Examples
 - **E-Promotions:** Based on your current location, your purchase history, where you are right now for store next to you <https://eduassistpro.github.io/>
 - **Healthcare monitoring:** sensors in your body → any abnormality Add WeChat **edu_assist_pro** require immediate reaction
 - **Disaster management and response**

Velocity in Real-world

- Every second, on average, around **6,000** tweets are tweeted on Twitter (visualize them here), which corresponds to over **350,000** tweets sent per minute, **500 million** tweets per day and around **200 billion** tweets per year.
- **Assignment Project Exam Help**

<http://www.inter>

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

Add WeChat **edu_assist_pro**

Extended Big Data Characteristics: 6V

- Volume: In a big data environment, the amounts of data collected and processed are much larger than those stored in typical relational databases.
- Variety: Big data consists of a rich variety of data types.
- Velocity: Big data arrives to the organization at high speeds and from multiple sources simultaneously.

- Veracity: Data quality in a big data context.
<https://eduassistpro.github.io/> challenging in a big data context.
- Visibility/Visualization: After big data is collected, we need a way of presenting the data in a manner that's readable and accessible.
- Value: Ultimately, big data is meaningless if it does not provide value toward some meaningful goal.

Veracity (Quality & Trust)

- *Data = quantity + quality*
- When we talk about big data, we typically mean its quantity:
 - What capacity of a system provides to cope with the sheer size of the data?
 - Is a query feasible on big data within our available resources?
 - How can we handle big data?
 - ... <https://eduassistpro.github.io/>
- Can we trust the answers to our queries?
 - Dirty data routinely lead to misinformed reports, strategic business planning decision ⇒ loss of revenue, credibility and customers, disastrous consequences
- *The study of data quality is as important as data quantity*

Data in real-life is often dirty

81 million National Insurance numbers but only 60 million eligible citizens

Assignment Project Exam Help

98000 deaths each year caused by errors in medical data

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

Add WeChat edu_assist_pro

500,000 dead people retain active Medicare cards

Visibility/Visualization

- Visibility: the state of being able to see or be seen is implied.
 - Big Data – visibility = Black Hole?
- Visualization: Making all that vast amount of data comprehensible in a manner that is easy to understand and read

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

A visualization of Divvy bike rides across Chicago

- Big data visualization tools:



Value

- Big data is meaningless if it does not provide value toward some meaningful goal

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Big Data: 6V in Summary

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Transforming Energy and Utilities through Big Data & Analytics. By Anders Quitzau@IBM

Other V's

- Variability
 - Variability refers to data whose meaning is constantly changing. This is particularly the case when gathering data relies on language processing.
- Viscosity
 - This term is sometimes used to describe the latency or lag time in the data relative to the event being described. We found that this is just as easily understood a
- Volatility
 - Big data volatility refers to how lid and how long should it be stored. You need to Add WeChat **edu_assist_pro** that point is data no longer relevant to the current analysis.
- More V's in the future ...
 - How many v's are there in big data?
<http://www.clc-ent.com/TBDE/Docs/vs.pdf>

Tag Clouds of Big Data

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



Why Study Big Data Technologies?

- The hottest topic in both research and industry
- Highly demanded in real world
- A promising future career
 - Research and development of big data systems:
Assignment Project Exam Help
distributed systems (eg, Hadoop), visualization tools, data wareh
ata quality control, ...
 - Big data ap
social marketing, healthc
https://eduassistpro.github.io/
Add WeChat edu_assist_pro
 - Data analysis: to get values out of big data
discovering and applying patterns, predicative analysis, business intelligence, privacy and security, ...
- Get enough credits

Big Data Open Source Tools

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

What will the course cover

- Topic 1. Big data management tools
 - Apache Hadoop
 - MapReduce
 - YARN/HDFS/HBase/Hive/Pig (briefly introduced)
 - Spark
 - AWS
 - Mahout [<https://eduassistpro.github.io/>]

Add WeChat edu_assist_pro

- Topic 2. Big data typical application
 - Finding similar items
 - Graph data processing
 - Data stream mining
 - Recommender Systems

Distributed processing is non-trivial

- How to assign tasks to different workers in an efficient way?
- What happens if tasks fail?
- How do workers exchange results?
- How to synchronize distributed tasks allocated to different workers?

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Big data storage is challenging

- Data Volumes are massive
- Reliability of Storing PBs of data is challenging
- All kinds of failures: Disk/Hardware/Network Failures
- Probability of failures simply increase with the number of machines ...

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

What is Hadoop

- Open-source data storage and processing platform
- Before the advent of Hadoop, storage and processing of big data was a big challenge
- Massively scalable, automatically parallelizable
 - Based on work from Google
- Assignment Project Exam Help
 - Google: (Not open)
 - Hadoop: [https://eduassistpro.github.io/
HBase opensource](https://eduassistpro.github.io/HBase)
- Named by Doug Cutting in 2006 (*at that time*), after his son's toy elephant.



Hadoop offers

- Redundant, Fault-tolerant data storage
- Parallel computation framework
- Job coordination



Programmers

Assignment Project Exam Help

N

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

Add WeChat edu_assist?pro



Q: Where file is located?

: How to handle failures & data

How to divide computation?

Q: How to program for scaling?

Why Use Hadoop?

- Cheaper
 - Scales to Petabytes or more easily
- Faster
 - Parallel data processing
- Better
 - Assignment Project Exam Help
 - Suited for p ms

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

Add WeChat edu_assist_pro

Companies Using Hadoop



Assignment Project Exam Help

The New York Times

JPMorganChase

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



Add WeChat edu_assist_pro



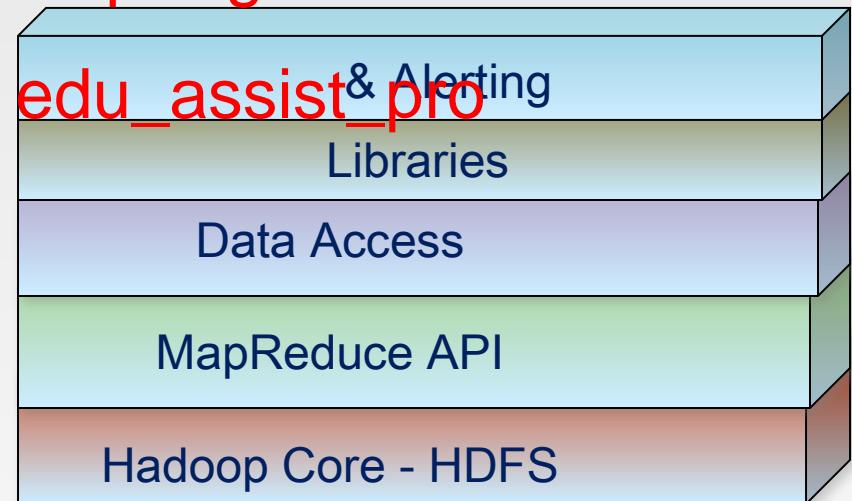
NING



YAHOO!

Hadoop is a set of Apache Frameworks and more...

- Data storage (**HDFS**)
 - Runs on commodity hardware (usually Linux)
 - Horizontally scalable
- Processing (**MapReduce**)
 - Parallelized (scalable) processing
 - Fault Tolerant
- Other Tools / Frameworks
 - https://eduassistpro.github.io/
 - Data Access
 - ▶ HBase, Hive, Pig, Mahout
 - Tools
 - ▶ Hue, Sqoop
 - Monitoring
 - ▶ Greenplum, Cloudera



What are the core parts of a Hadoop distribution?

HDFS Storage

Redundant (3 copies)

For large files – large blocks

64 or 128 MB / block

Can scale to 1000s of nodes

MapReduce API

Assignment Project Exam Help

Other Libraries

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

Add WeChat edu_assist_pro

Fault-tolerant (auto retries)

Adds high availability and more

Others

Hadoop 2.0

- Single Use System
 - Batch apps
- Multi-Purpose Platform
 - Batch, Interactive, Online, Streaming

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Hadoop YARN (Yet Another Resource Negotiator): a resource-management platform responsible for managing computing resources in clusters and using them for scheduling of users' applications

Hadoop Ecosystem

A combination of technologies which have proficient advantage in solving business problems.

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

<http://www.edupristine.com/blog/hadoop-ecosystem-and-components>

Common Hadoop Distributions

- Open Source

 - Apache

- Commercial

 - Cloudera

Assignment Project Exam Help

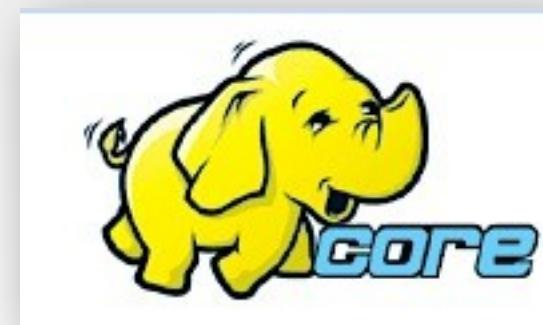
 - Hortonwork

 - MapR

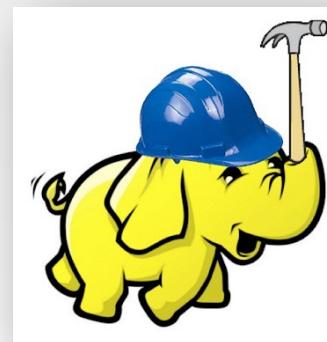
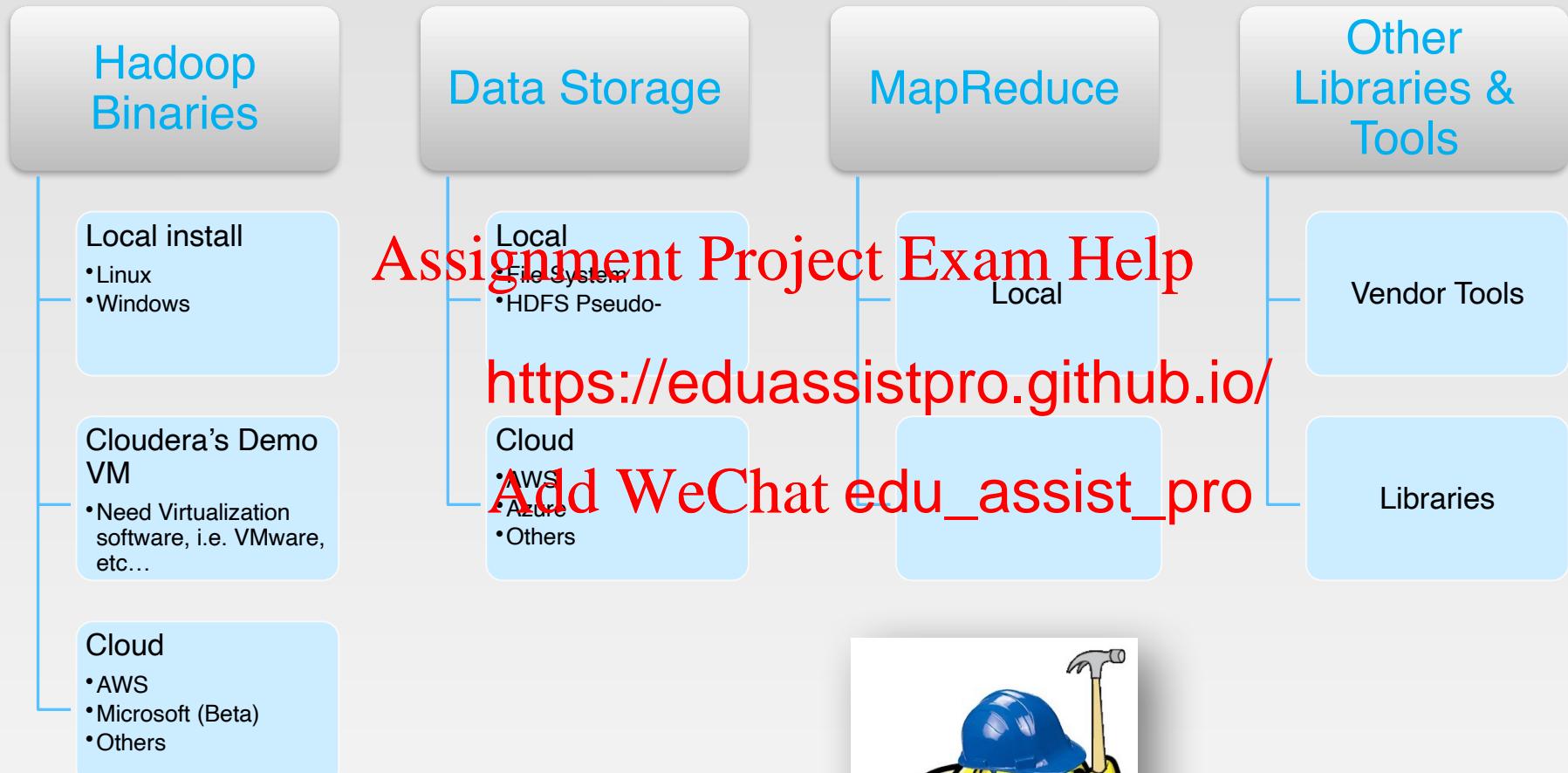
<https://eduassistpro.github.io/>

 - AWS MapReduce

 - Microsoft Azure HDInsight (Bet



Setting up Hadoop Development

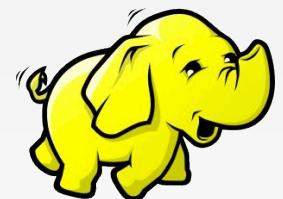


Comparing: RDBMS vs. Hadoop

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



The Changing Data Management Landscape

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Philosophy to Scale for Big Data Processing

Divide Work

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



**Combine
Results**

MapReduce

- Typical big data problem
 - Iterate over a large number of records
 - Extract something of interest from each **Map**
 - Shuffle and sort intermediate results
 - Aggregate intermediate results
 - Generate final output
- Programmers specify two functions:
 - map** (k_1, v_1) \rightarrow [k_2, v_2]
 - reduce** ($k_2, [v_2]$) \rightarrow [k_3, v_3]
 - All values with the same key are sent to the same reducer
- The execution framework handles everything else...

Understanding MapReduce

- Map>>
 - $(K1, V1) \rightarrow$
 - Info in
 - Input Split
 - $\text{list}(K2, V2)$
 - Key / Value out (intermediate values)
 - One list per local node
 - Can implement local Reducer (or Combiner)
- Shuffle/Sort>>
 - $(K2, \text{list}(V2)) \rightarrow$
 - Shuffle / Sort phase precedes Reduce phase
 - Combines Map output into a list
- Reduce
 - $\text{list}(K3, V3)$
 - Usually aggregates intermediate values

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

(input) $\langle k1, v1 \rangle \rightarrow \text{map} \rightarrow \langle k2, v2 \rangle \rightarrow \text{combine} \rightarrow \langle k2, \text{list}(V2) \rangle \rightarrow \text{reduce} \rightarrow \langle k3, v3 \rangle$ (output)

WordCount - Mapper

- Reads in input pair $\langle k1, v1 \rangle$
- Outputs a pair $\langle k2, v2 \rangle$
 - Let's count number of each word in user queries (or Tweets/Blogs)
 - The input to the mapper will be $\langle \text{queryID}, \text{QueryText} \rangle$:

~~<Q1, Assignment Project Exam Help
store opens in the store was closed; the
store opens in the .>~~

- The output <https://eduassistpro.github.io/>

~~<The, 1> <teacher, 1> <the, 1> <store, 1>
<the, 1> <store, 1> <was, 1> <the, 1> <store, 1>
<opens, 1> <in, 1> <the, 1> <store, 1>
<opens, 1> <at, 1> <9am, 1> <the, 1> <store, 1>~~

WordCount - Reducer

- Accepts the Mapper output (k_2, v_2), and aggregates values on the key to generate (k_3, v_3)

- For our example, the reducer input would be:

```
<The, 1> <teacher, 1> <went, 1> <to, 1> <the, 1> <store, 1>  
<the, 1> <store, 1> <was, 1> <closed, 1> <the, 1> <store, 1>  
<opens, 1> <in, 1> <the, 1> <morning, 1> <the, 1> <store, 1>  
<opens, 1>
```

- The output <https://eduassistpro.github.io/>

```
<The, 6> <teacher, 1> <went, 1> <store, 4> <was, 1>  
<closed, 1> <opens, 1> <in, 1> <at, 1> <9am, 1>
```

Assignment Project Exam Help

Add WeChat edu_assist_pro

MapReduce Example - WordCount

Assignment Project Exam Help

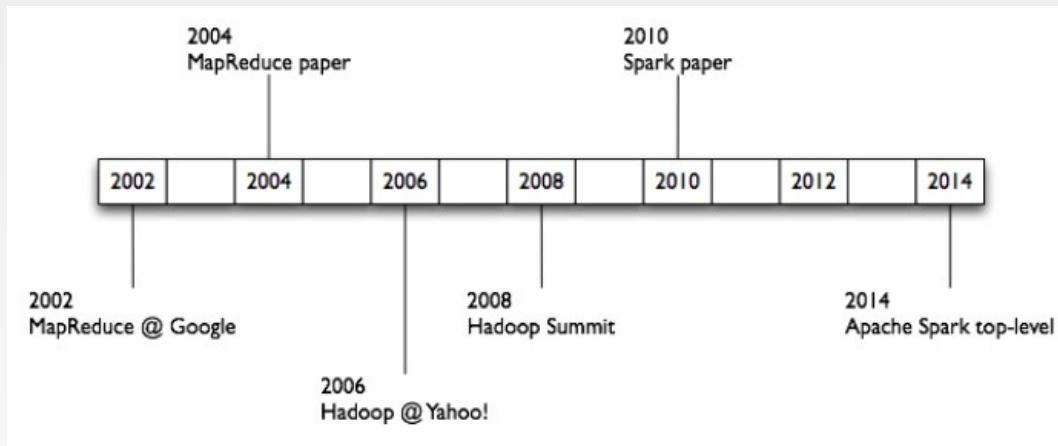
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Hadoop MapReduce is an implementation of MapReduce
 - MapReduce is a computing paradigm (Google)
 - Hadoop MapReduce is an open-source software

Spark

- One popular answer to “What’s beyond MapReduce?”
- Open-source engine for large-scale data processing
 - Supports generalized dataflows
 - Written in Scala, with bindings in Java and Python
- Brief history:
Assignment Project Exam Help
 - Developed
 - Open-sourc <https://eduassistpro.github.io/>
 - Became top-level Apache proje 2014
 - Commercial support provided b Add WeChat **edu_assist_pro**



Spark

- Fast and expressive cluster computing system interoperable with Apache Hadoop
 - Improves efficiency through:
 - In-memory computing primitives
 - General computation graphs
 - Improves usability
 - Rich APIs in <https://eduassistpro.github.io/>
 - Interactive shell
 - **Spark is not**
 - a modified version of Hadoop
 - dependent on Hadoop because it has its own cluster management
 - Spark uses Hadoop for storage purpose only
- Up to 100× faster
(2-10× on disk)
- Assignment Project Exam Help
- Add WeChat  edu_assist_pro 5x less code

Spark Platform

- Spark is the basis of a wide set of projects in the Berkeley Data Analytics Stack (BDAS)

Shark
(SQL)

Assignment Project Exam Help
S
(<https://eduassistpro.github.io/>)

MLlib
(machine learning)

...

Add WeChat edu_assist_pro
Spar

- Spark SQL (SQL on Spark)
- Spark Streaming (stream processing)
- GraphX (graph processing)
- MLlib (machine learning library)

Spark

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

WordCount in Spark (Scala)

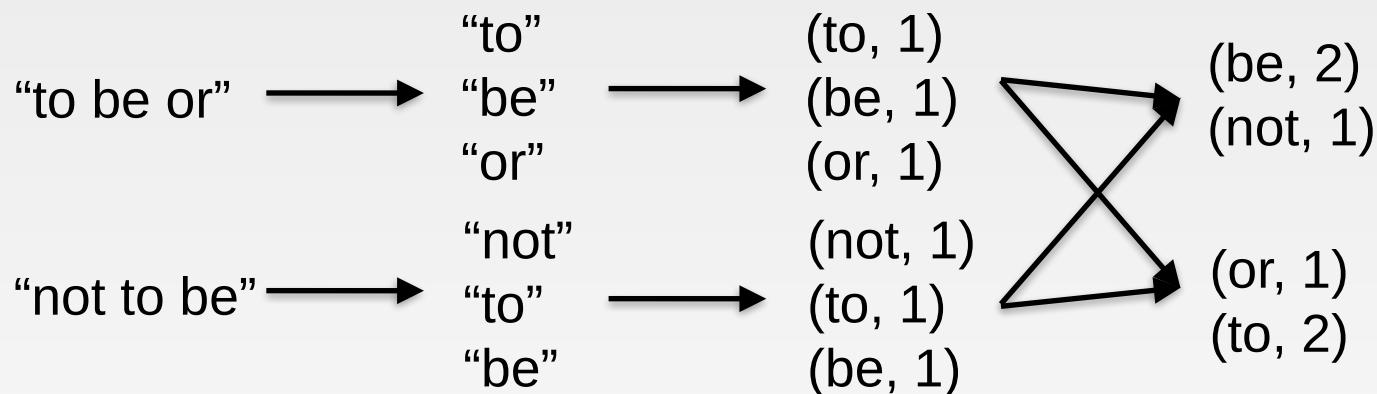
Transformation

Assignment Project Exam Help

Action

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

Add WeChat edu_assist_pro



AWS (Amazon Web Services)

□ Amazon

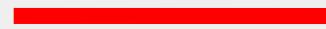
From Wikipedia 2006

From Wikipedia 2017

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



AWS (Amazon Web Services)

- AWS is a subsidiary of Amazon.com, which offers a suite of cloud computing services that make up an on-demand computing platform.
- Amazon Web Services (AWS) provides a number of different services, including:
 - Amazon Elastic Compute Cloud (EC2)
Virtual machines for running custom software
 - Amazon Simple Key-Value Store (SimpleDB)
NoSQL database service
 - Amazon Elastic MapReduce (EMR)
Scalable MapReduce computation
 - Amazon DynamoDB
Distributed NoSQL database, one of several in AWS
 - Amazon SimpleDB
Simple NoSQL database
 - ...

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Cloud Computing Services in AWS

- IaaS
 - EC2, S3, ...
 - Highlight: EC2 and S3 are two of the **earliest** products in AWS
- PaaS
 - Aurora, Redshift, ...
 - Highlight: A **e fastest growing** products in <https://eduassistpro.github.io/>
- SaaS
 - WorkDocs, WorkMail
 - Highlight: May not be the main focus of AWS

Setting up an AWS account

aws.amazon.com

Assignment Project Exam Help

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

Add WeChat edu_assist_pro



- Sign up for an account on aws.amazon.com
 - You need to choose an username and a password
 - These are for the management interface only
 - Your programs will use other credentials (RSA keypairs, access keys, ...) to interact with AWS

Signing up for AWS Educate

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

- Complete the web form on
<https://aws.amazon.com/education/awseducate/>
 - Assumes you already have an AWS account
 - Use your UNSW email address!
 - Amazon says it should only take 2-5 minutes (but don't rely on this!!)
- This should give you \$100/year in AWS credits. **Be careful!!!**

Big Data Applications

- Finding similar items
- Graph data processing
- Data stream mining
- Recommender <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

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

Add WeChat edu_assist_pro