

ENSF 612 Lecture

Big data platforms overview: Hadoop vs Spark Architecture

Dr. Gias Uddin, Assistant Professor
Electrical and Software Engineering
Schulich School of Engineering
University of Calgary
<https://giasuddin.ca/>

Topics

- Distributed vs Cluster Computing
- Hadoop vs Spark
- Hadoop Architecture
- Spark Architecture
- PySpark Programming

Distributed vs Cluster Computing

- Distributed Computing: Split a business into several sub-services and distribute those services on different machines
- Cluster Computing: Multiple servers are grouped together to achieve the same service.

Apache Hadoop: What is it?



- Wikipedia: “Apache Hadoop is Collection of open source software utilities that facilitate using a network of many computers to solve problems involving massive amounts of data and computation”
- It’s a software framework to store and process data in multiple computers across a cluster
- The storage and processing is handled by a novel algorithm “MapReduce Programming Model”. There will be separate lecture on this.
- Originally focused on clusters built on commodity computers (e.g., your or my laptops) but it also expanded into higher-end computers

Apache Hadoop: Framework



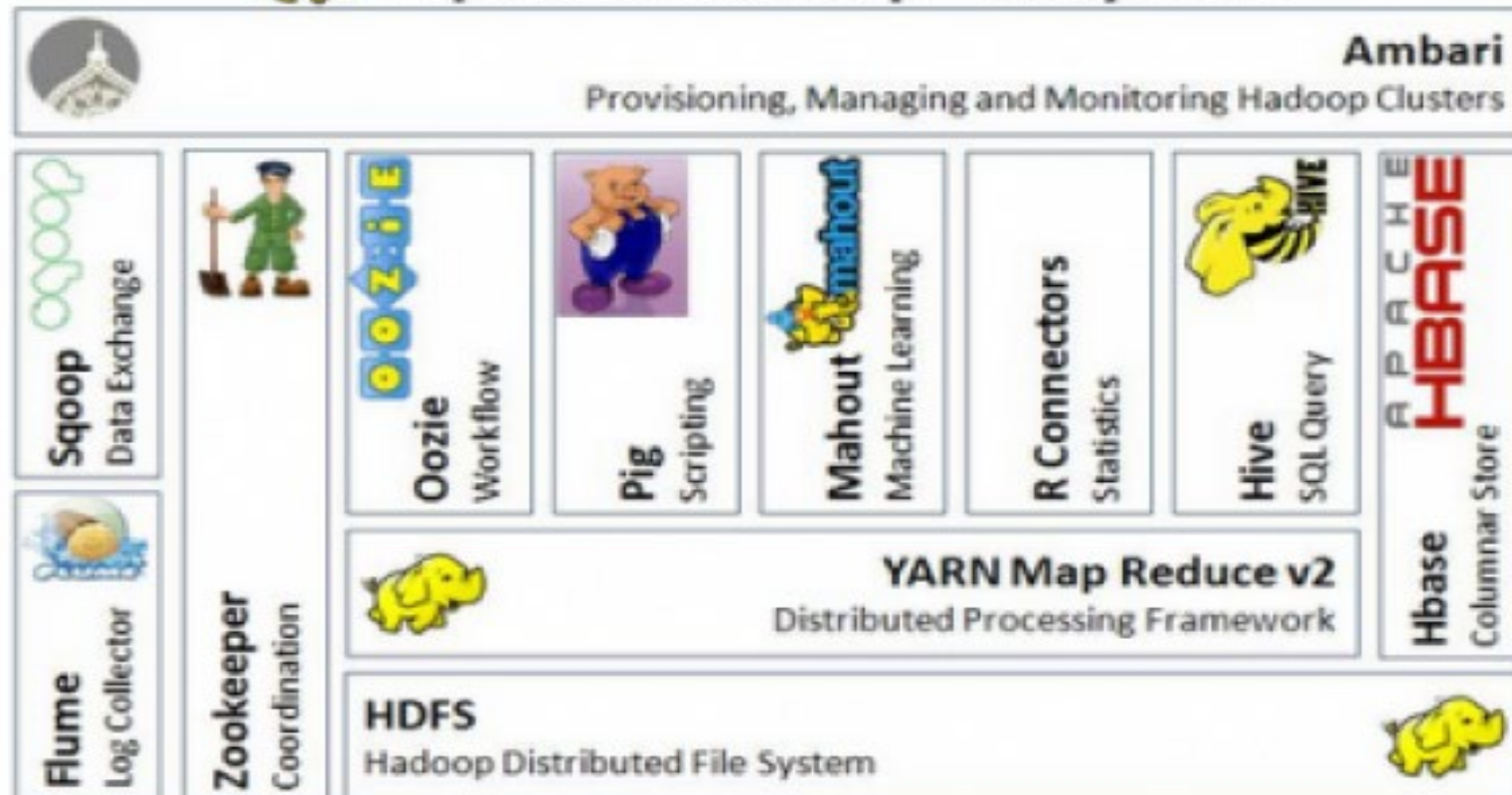
- Four base modules

Module	Description
Hadoop common	Contains libraries and utilities needed by other Hadoop modules
Hadoop distributed file system (HDFS)	A distributed file-system to store and process data on the commodity machines, by providing high bandwidth across the cluster. Inspired by Google File System (GFS). See research paper from Google on GFS: https://static.googleusercontent.com/media/research.google.com/en/archive/gfs-sosp2003.pdf
Hadoop Yarn	A resource management platform responsible for managing and scheduling the commodity machines across the cluster for computing jobs
Hadoop MapReduce	A programming model to support the storage and processing data across the cluster by utilizing HDFS. Inspired by Google research paper on MapReduce: https://research.google/pubs/pub62/

Apache Hadoop: Ecosystem



Apache Hadoop Ecosystem

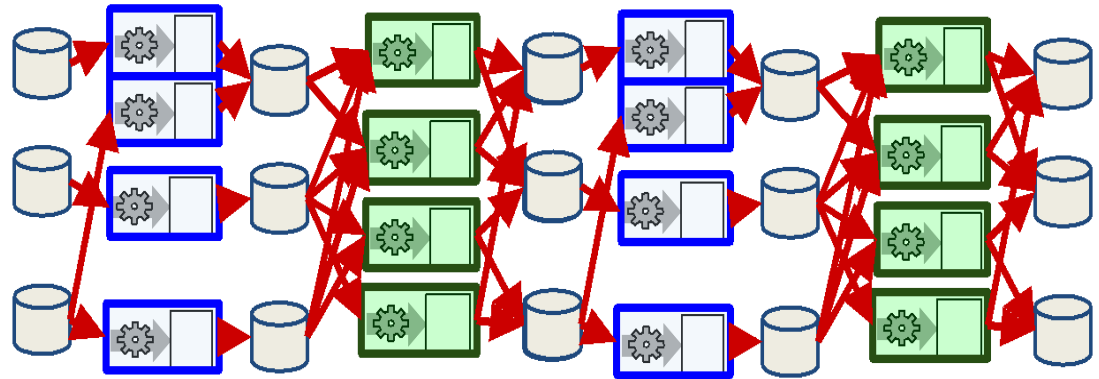


Disadvantages of Hadoop

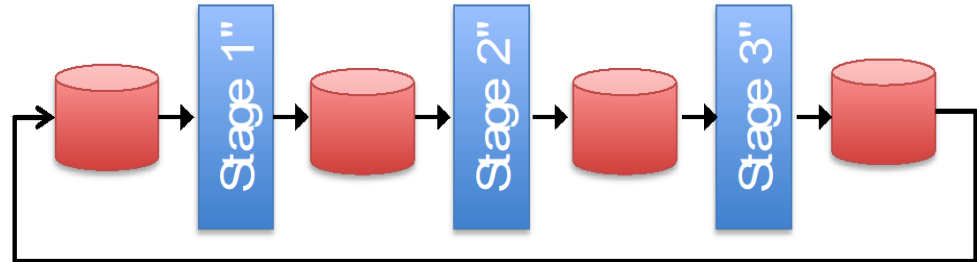
- ◆ Each stage involves disk

I/O

- ⑩ Map, group by key, and reduce phases read from disk



- ◆ Disk I/O is very slow!
- ◆ Problem gets magnified in iterative jobs

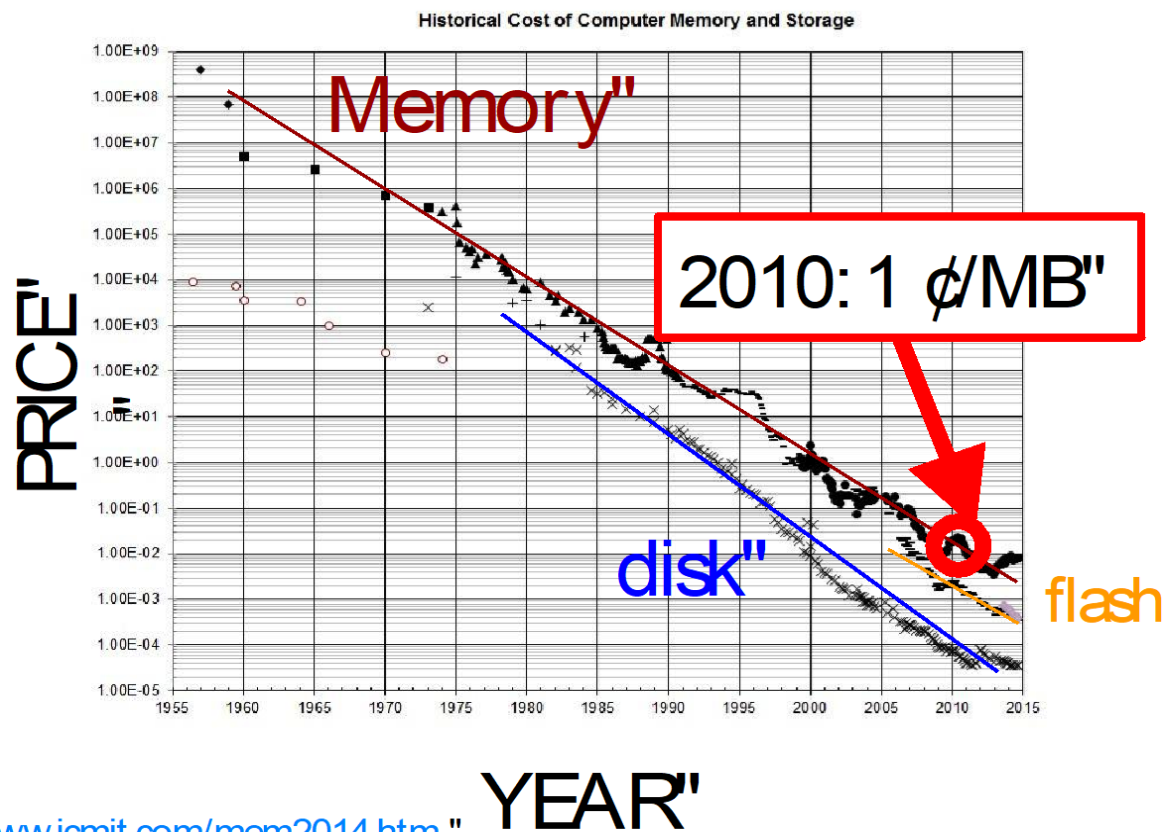


Disadvantages of Hadoop– cont'd

- ◆ MapReduce may not be well suited for complex jobs
- ◆ Jobs run at speed of disk – not speed of processors
- ◆ Need a different way of dealing with complex jobs
 - ⑩ Iterative jobs – e.g., BFS, matrix multiplication, clustering
 - ⑩ Interactive jobs – querying a dataset in multiple ways
 - ⑩ Streaming jobs – dealing with fast arriving data
- ◆ **Avoid going to disk as much as possible!**

Disadvantages of Hadoop– cont'd

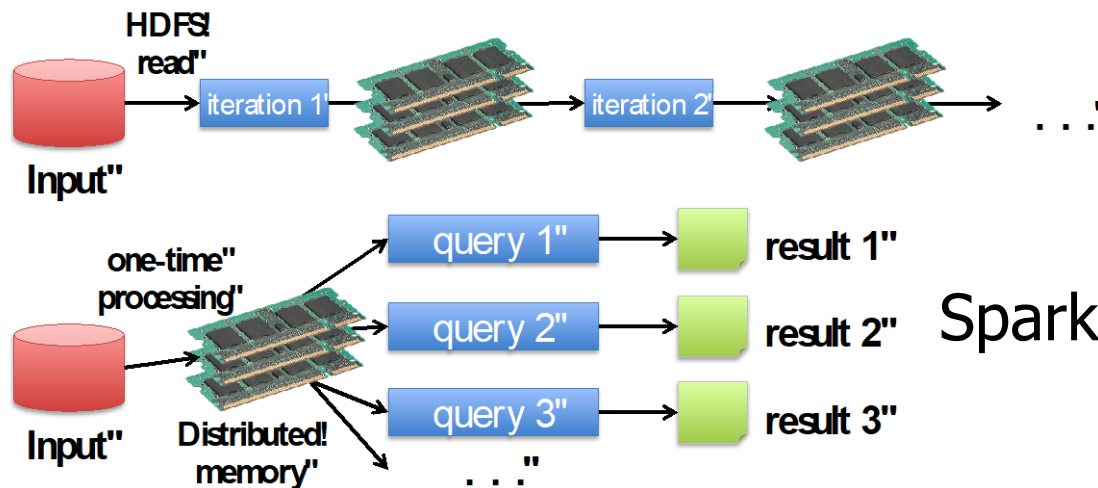
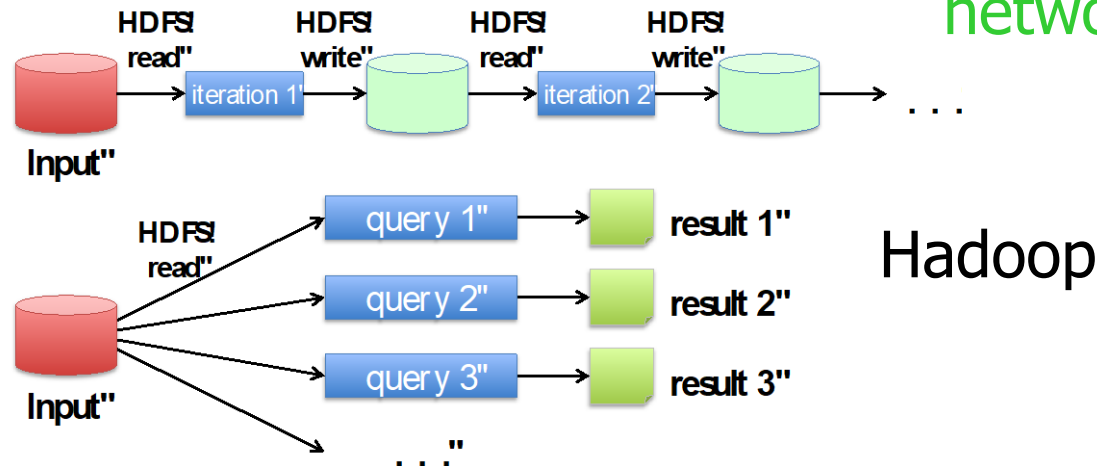
- ◆ **Solution: keep data in memory**
- ◆ Memory costs are falling
- ◆ Equip servers with more memory



Apache Spark architecture

- ◆ New cluster computing platform
- ◆ Allows in-memory computations

10-100x faster than
network and disk



Apache Spark architecture – cont'd

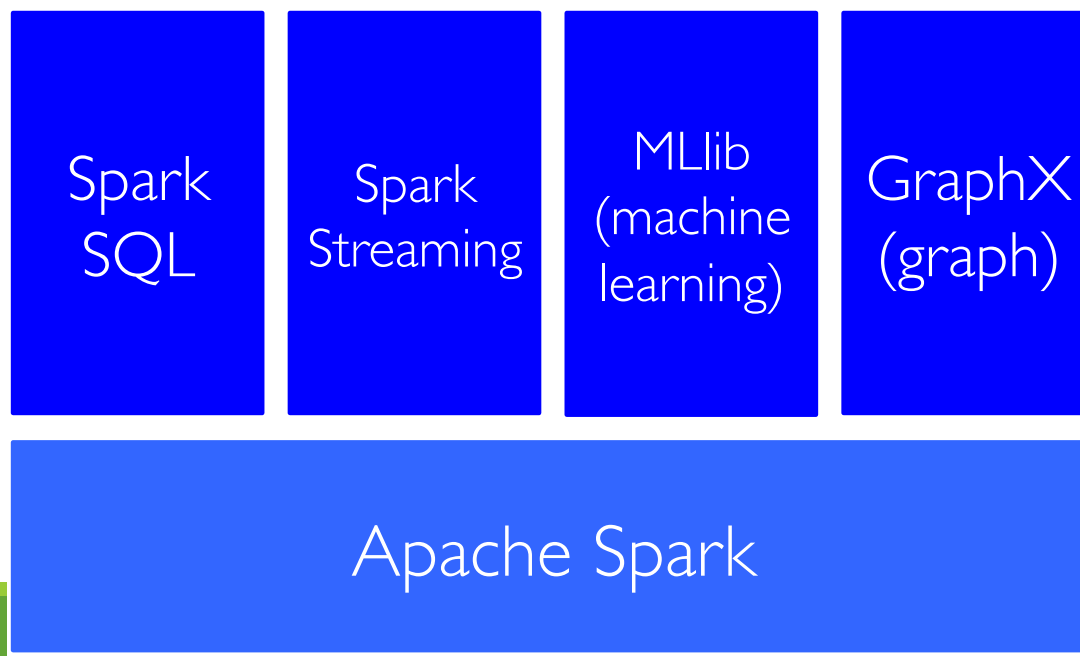
- ◆ Data may not fit memory of one node
- ◆ Need to distribute across memories of many nodes
- ◆ Spark uses Resilient Distributed Datasets (**RDDs**) for this
- ◆ RDD – partitioned collection of objects
 - ⑩ Stored across disks or memories of many nodes
- ◆ Can build and manipulate RDDs
 - ⑩ Use parallel **transformations** and **actions**
- ◆ RDDs automatically rebuilt on machine failures

Apache Spark architecture– cont'd

- ◆ Spark hides complexities from programmer
 - ⑩ Parallel execution, fault tolerance
- ◆ Programmer - “Here is an operation. Run it on my data”
- ◆ Spark parallelizes the operation and runs it on many nodes
- ◆ Spark re-runs jobs to handle failures

Apache Spark architecture– cont'd

- ◆ Spark provides core functionality
 - ⑩ RDDs, transformations, actions
- ◆ Spark SQL – provides SQL-like interface
- ◆ Spark Streaming – deals with rapid data streams
- ◆ MLlib – helps build machine learning models from data
- ◆ GraphX – simplifies graph algorithm development



Apache Spark architecture– cont'd

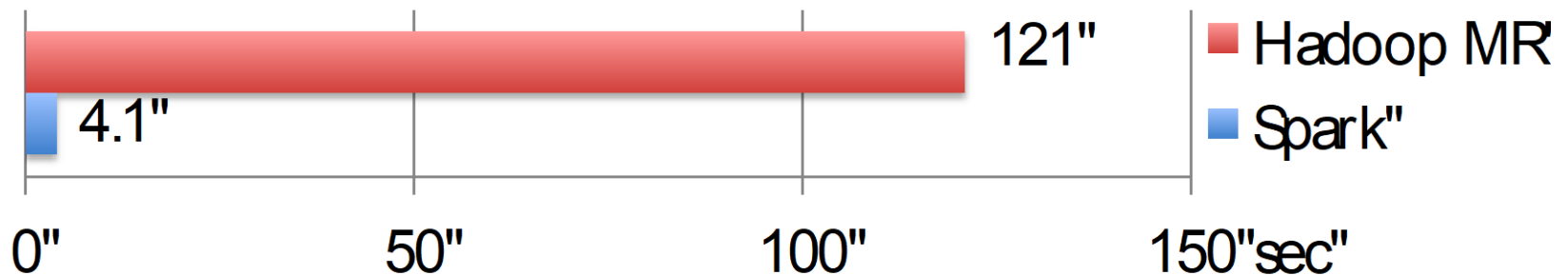
◆ Summary: Hadoop vs. Spark

	Hadoop Map Reduce	Spark
Storage	Disk only	In-memory or on disk
Operations	Map and Reduce	Map, Reduce, Join, Sample, etc...
Execution model	Batch	Batch, interactive, streaming
Programming environments	Java	Scala, Java, R, and Python

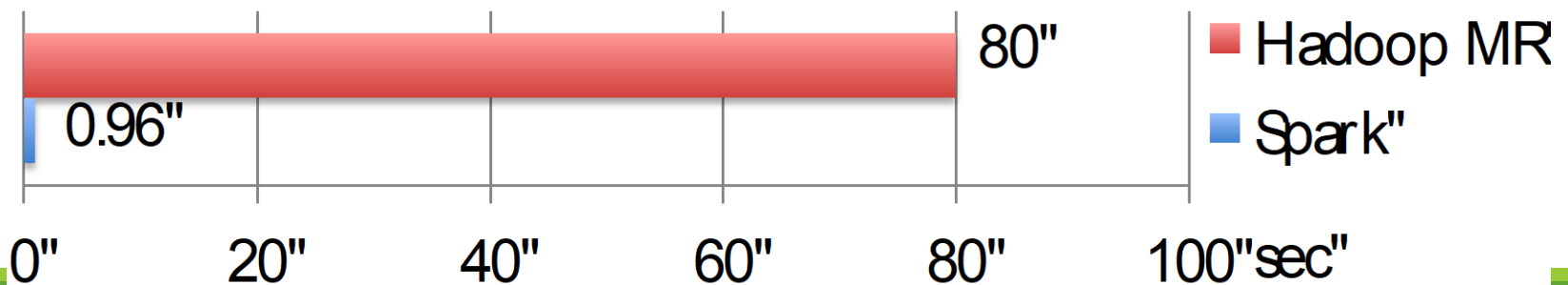
Apache Spark architecture– cont'd

- ◆ In-memory computation can improve performance
- ◆ Spark outperforms Hadoop MapReduce for iterative jobs

K-means Clustering

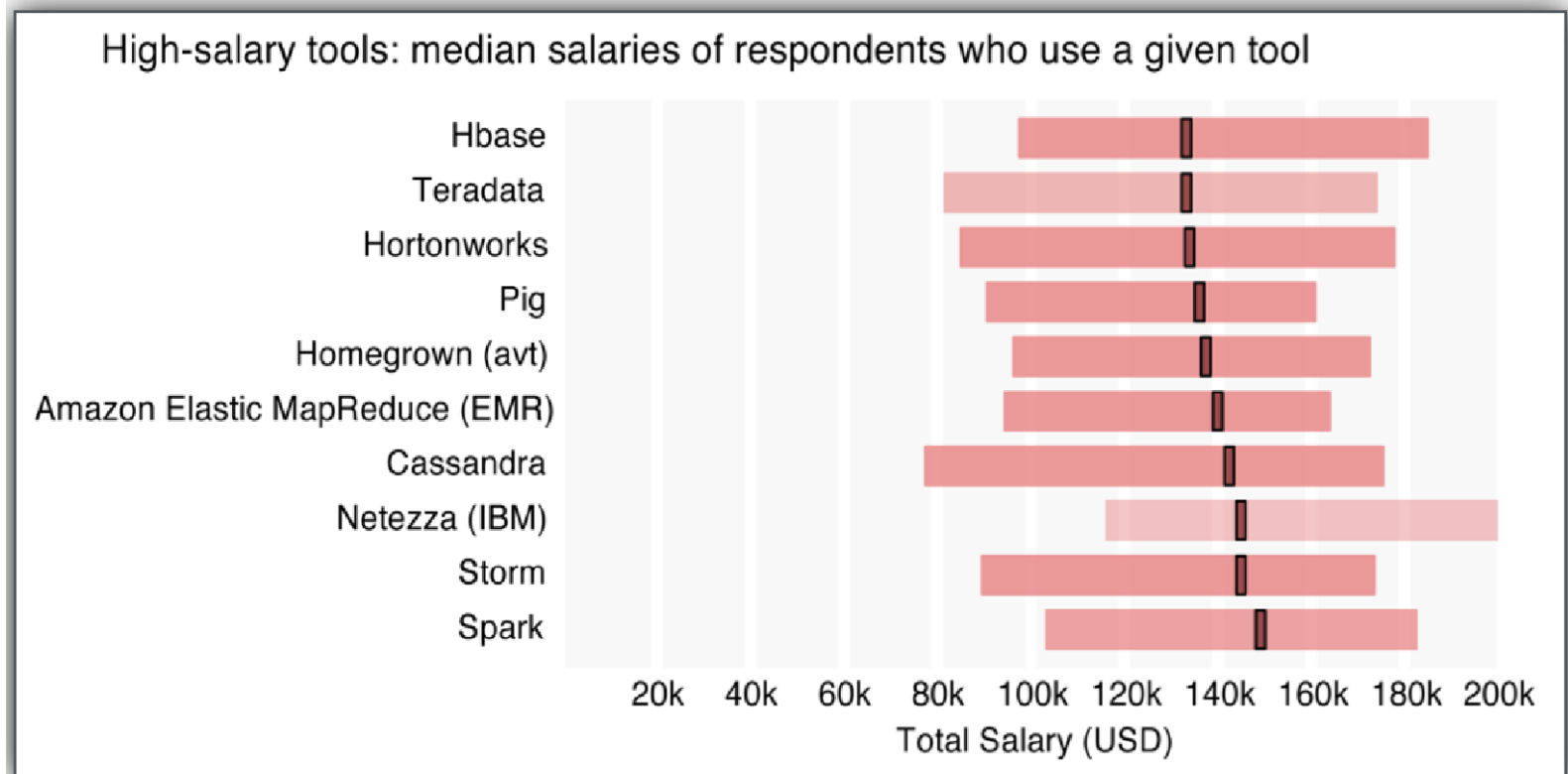


Logistic Regression



Apache Spark architecture– cont'd

Still not convinced you need to learn Spark? ☺



Over 800 respondents across 53 countries and 41 U.S. states"

<http://www.oreilly.com/data/free/2014-data-science-salary-survey.csp>

Topics

- Distributed vs Cluster Computing
- Hadoop vs Spark
- Hadoop Architecture
- Spark Architecture

Acknowledgements

Portions of these slides were adapted from external material available under creative commons license CC-BY-NC-SA 4.0. This license grants the ability to share and adapt the material for non-commercial purposes.

Name of the creator: Dr. Anthony Joseph, University of Berkeley and team

License notice: <https://creativecommons.org/licenses/by-nc-sa/4.0/legalcode>

Copyright notice: CC-BY-NC-SA 4.0

Link to material: <https://courses.edx.org/courses/BerkeleyX/CS100.1x/1T2015/info>