

HDP Overview: Essentials

A Technical Understanding for Business Users and Decision Makers

Hortonworks. We do Hadoop.

A Hortonworks University Training Course

Course Objectives

- The case for Hadoop
- What is Hadoop and its ecosystem
- Moving data in and out of Hadoop
- Processing data in Hadoop
- The Hortonworks Data Platform (HDP)
- Preparing your Enterprise for Hadoop
- Key roles in a Hadoop environment



Course Outline

Unit 1 Understanding Big Data Unit 2 Understanding Hadoop Unit 3 Data Integration Unit 4 The Hadoop Ecosystem Unit 5 Adoption



Unit 1 – Understanding Big Data

Data, data, and more data





Who am I?

izhar@abyres.net





Who is Hortonworks?



The leaders of Hadoop's development

100% Open Source – Democratized Access to Data

Hortonworks do Hadoop.

Community driven, Enterprise Focused Drive Innovation in the platform – We lead the roadmap



Hortonworks do Hadoop successfully.

Contribute to Apache

Support Hadoop adopters Provide Consulting Services Deliver Training



Who is ABYRES





Support and service providers for community Open Source softwares

Local support and service providers for products from Open Source principals such as Hortonworks and Red Hat

We do Open Source Solutions ur customers move from proprietary lock-in into

Open Source

The only Hortonworks System Integrator and Training Partner in South East Asia



Our Clients





What is Hadoop? What is everyone talking about?





Big



"Big Data" is the marketing term of the decade in IT



What lurks behind the hype is the democratization of Data.



You need data.



But what do you do with your data now?



We are obsessive compulsive about collecting and structuring our data.



Put it away, delete it, tweet it, compress it, shred it, wikileak-it, put it in a database, put it in SAN/NAS, put it in the cloud, hide it in tape...



You need value from your data. You need to make decisions from your data.



So what are the problems with Big Data?



Volume

Volume

Volume

Volume





Volume Volume Volume Volume

Volume

Volume Volume Volume Volume Volume



Volume WAR HAMBE Volument /olume Volume Volume Volume Volume Volume Volume olume Volume Hortonworks © Hortonworks Inc. 2011 – 2014. All Rights Reser

Volume

Volume Volume Volume Volume Volume Volume Volume Volume Volu Volu**va**tum@olur Volume Volume Volume mevolume Volume Volume Volume ne Volume Ime Volume Volume √olume ™olume Volume Volume Volume 'me Volume Volume Volume Volume Volume Volume Volume Volume` Volume Volume Umae Volume Volume Volume Ime Volume Volume V**ର୍ବା**ଥିକିକ Volume Volume Volumene Volume Volume Volvoheme Volume ^{⊥m}¢olume Volume Volume Volume Volum Volume Volume Volumme VolumeVolume Volume Volume **Wheevolume** Volume Volume Volume Volume ie Volume **Voblima**e Volume Hortonw © Hortonworks Inc. 2011 – 2014. All Rights Re

Storage, Management, Processing all become challenges with Data at Volume



Traditional technologies adopt a divide, drop, and conquer approach





Another DB Data Data Data Data Data Data



The solution?









Another DB Data Data Data Data Jata Data Data





Analyzing the data usually raises more interesting questions...



...which leads to more data



Wait, you've seen this before.



Data begets Data.



What keeps us from our Data?


"Prices, Stupid passwords, and Boring Statistics." - Hans Rosling

http://www.youtube.com/watch?v=hVimVzgtD6w



Your data silos are lonely places.



... Data likes to be together.









New types of data don't quite fit into our pristine view of the world.



To resolve this, some of us have learned a lot from Lord Of The Rings...



...to create One-Schema-To-Rule-Them-All...







...but that has its problems too.







Fragile workflows make supporting the analytical models you want expensive and time-consuming.





What do you want to do with data?



Marketing Analytics needs data. Work with the population, not just a sample.





Middle Income Band

Your segmentation today.

Age: 25-30

Female

Male

Product Category Preferences



Spent 25 minutes looking at tea cozies

GPS coordinates

Walking into Starbucks right now...

Looking to start a business

Unhappy with his cell phone plan

\$65-68k per year

Product recommendations

Your segmentation with better data.

A depressed Chicago Blackhawks Fan

High tea party

Pregnant

Hippie

Age: 27 but feels old

Thinking about a new house

Gene Expression for Risk Taker

Products left in basket indicate drunk Amazon shopper

Female



Male

Pick up all of that data that was prohibitively expensive to store and use.



Why do viewer surveys...



...when raw data can tell you what button on the remote was pressed during what commercial for the entire viewer population?



To approach these use cases you need an affordable platform that stores, processes, and analyzes the data.



So what is the answer?



Enter the Hadoop.





Hadoop was created because traditional technologies never cut it for the Internet properties like Google, Yahoo, Facebook, Twitter, and LinkedIn



Traditional architecture didn't scale enough...

Арр Арр Арр Арр



SAN

Арр Арр Арр Арр



SAN



SAN

App App

App



App

Databases can become bloated and create problems rather than help solve them.





Traditional architectures cost too much at that volume...







So what is the answer?



If you could design a system that would handle this, what would it look like?



It would probably need a highly resilient, self-healing, cost-efficient, distributed file system.

Storage Storage Storage

Storage Storage Storage

Storage Storage Storage



It would probably need a completely parallel processing framework that took *tasks* to the *data*...

Processing	Processing	Processing
Storage	Storage	Storage
Processing	Processing	Processing
Storage	Storage	Storage
Processing	Processing	Processing
Storage	Storage	Storage



It would probably run on commodity hardware, virtualized machines, and common OS platforms

Processing	Processing	Processing
Storage	Storage	Storage
Processing	Processing	Processing
Storage	Storage	Storage
Processing	Processing	Processing
Storage	Storage	Storage



It would probably be **open source** so innovation could happen as quickly as possible



It would need a critical mass of users



Hadoop 2 hit the ground recently: Introducing **YARN**



{Processing + Storage}

{MapReduce/YARN + HDFS}

{Hadoop}



YARN lets you run more data apps than ever before



YARN turns Hadoop into a smart phone: An App Ecosystem

hortonworks.com/yarn/



YARN: Yeah, we did that too.

hortonworks.com/yarn/






Unit 2 – Understanding Hadoop

Storage and Processing







What is HDFS?



The NameNode

1.When the NameNode starts, it reads 2. The transactions in **edits** are merged with fsimage and edits files from disk. fsimage, and edits is emptied. 3. A client application creates a new file in HDFS. hadoop fs -put foo.log bar/foo.log fsimage edits hdfs snapshots hdfs journals 4. The NameNode logs that transaction in the edits file. NameNode



The DataNodes







1. The first DataNode pipelines the replication to the next DataNode in the list.





Demonstration Understanding HDFS & Block Storage



Hortonworks

What is MapReduce? Break a large problem into sub-solutions





WordCount in MapReduce

constitution.txt

HDFS



The mappers read the file's blocks from HDFS line-by-line

We the people, in order to form a...



The reducers add up the "1's" and output the word and its count

<We,4>

<the.265>

<people,5>

<form,1>

<We, (1,1,1,1)> <the, (1,1,1,1,1,1,1,...)> <people,(1,1,1,1,1,1)>

<form, (1)>







Shuffle/Sort

Reduce Phase

SELECT word, COUNT(*) FROM constitution WHERE....

GROUP BY

Hive/Pig compile as Reduce side function

ORDER BY JOIN DISTINCT

Examples of more Reduce side functions





Spill files are merged into a single file

Records are sorted and spilled to disk when the buffer reaches a threshold



© Hortonworks Inc. 2013



Lifecycle of a YARN Application



A Cluster View Example





Demonstration WordCount MapReduce



Hortonworks







Welcome to the Hortonworks Data Platform (HDP) 2.1 Sandbox

HDP 2.1 is a major release that delivers required enterprise functionality for data management, data access, data governance, integration, security and operations developed and delivered completely in the open. Incorporating the very latest community innovations across all Apache Hadoop projects, HDP 2.1 provides the foundational platform for organizations looking to incorporate Hadoop in a modern data architecture.



Get started with Hadoop

NEW

Try New Features



Dive right in





Demonstration The Hortonworks Sandbox



Hortonworks



Unit 3 – Data Integration

Loading data into Hadoop





Options for Data Input



The Hadoop Client

- The **put** command to uploading data to HDFS
- Perfect for inputting local files into HDFS
 - -Useful in batch scripts
- Usage:

hadoop fs -put mylocalfile /some/hdfs/path

- POSIX utility commands such as 1s, mv, cp, touch, cat, mkdir are also supported
- Full list of commands hadoop fs



WebHDFS

- REST API for accessing all of the HDFS file system interfaces:
 - -http://host:port/webhdfs/v1/test/mydata.txt?op=OPEN
 - -http://host:port/webhdfs/v1/user/train/data?op=MKDIRS
 - -http://host:port/webhdfs/v1/test/mydata.txt?op=APPEND



Overview of Flume: Data Streaming





Overview of Sqoop: Database Import/Export







Demonstration Using Sqoop



Hortonworks



Unit 4 – The Hadoop Ecosystem

Tour of the Hadoop Ecosystem





Hadoop Ecosystem: Pig

- An engine for executing programs on top of Hadoop
- It provides a language, Pig Latin, to specify these programs





Why use Pig?

 Maybe we want to join two datasets, from different sources, on a common value, and want to filter, and sort, and get top 5

```
users = LOAD 'input/users' USING PigStorage(',')
           AS (name:chararray, age:int);
 2
 3
  filtrd = FILTER users BY age >= 18 and age <= 25;
 4
 5
   pages = LOAD 'input/pages' USING PigStorage(',')
           AS (user:chararray, url:chararray);
 8
 9
   jnd = JOIN filtrd BY name, pages BY user;
10
   grpd = GROUP jnd BY url;
11
12
   smmd = FOREACH grpd GENERATE group, COUNT(jnd) AS clicks;
13
14
  srtd = ORDER smmd BY clicks DESC;
15
16
17
  top5 = LIMIT srtd 5;
18
19 STORE Top5 INTO 'output/top5sites' USING PigStorage(',');
```



Hadoop Ecosystem: Hive

Use existing SQL tools and existing SQL processes





Mobile



Hortonworks

What is Hive?

- Data warehouse system for Hadoop
- Create schemas/table definitions that point to data in Hadoop
- Treat your data in Hadoop as tables
- SQL 92
- Interactive queries at scale



HiveQL

SQL Semantics
SELECT, LOAD, INSERT from query
Expressions in WHERE and HAVING
GROUP BY, ORDER BY, SORT BY
CLUSTER BY, DISTRIBUTE BY
Sub-queries in FROM clause
GROUP BY, ORDER BY
ROLLUP and CUBE
UNION
LEFT, RIGHT and FULL INNER/OUTER JOIN
CROSS JOIN, LEFT SEMI JOIN
Windowing functions (OVER, RANK, etc.)
Sub-queries for IN/NOT IN, HAVING
EXISTS / NOT EXISTS
INTERSECT, EXCEPT



Hive Architecture





Hadoop Ecosystem: HCatalog

- Hive component
- Glue between Pig & Hive
- Schema visibility to

 Pig Scripts & MapReduce
- REST API to
 - -Access Hive schemas
 - -Submit DDL
 - -Launch Hive queries
 - -Launch Pig jobs
 - -Launch MR
 - -Notifications to message broker




Overview of the Hive ODBC Driver





http://hortonworks.com/hdp/addons



ROI from Your Cluster



MapReduce applies to *a lot* of data processing problems





MapReduce goes a long way, but not all data processing and analytics are solved the same way



Sometimes your data application needs parallel processing *and* inter-process communication





...like complex event processing in Apache Storm



Sometimes your machine learning data application needs to process in memory and iterate



...like in Machine Learning in Spark



Overview of Stinger



Performance Optimizations

100X+ Faster Time to Insight

Deeper Analytical Capabilities

Hortonworks







Back to YARN: Taking Hadoop Beyond Batch

• With YARN, applications run natively *in* Hadoop





Hadoop Security

Kerberos Authentication All services can be principals HBase, Hive & HDFS Authorization ACLs for each Wire encryption Basic Auditing HDFS, Shuffle, JDBC Knox: Hadoop REST Gateway Logging of events **Centralized Security Administration** Centralized Audit Reporting Perimeter security for Hadoop cluster Delegated Policy Administration

XA Secure



Defining an Oozie Workflow



Falcon: Data Lifecycle Management

Falcon At A Glance



- > Falcon provides the key services data processing applications need.
- > Complex data processing logic handled by Falcon instead of hard-coded in apps.
- > Faster development and higher quality for ETL, reporting and other data processing apps on Hadoop.







Unit 5 – Hadoop Adoption

Charting a path to adopt Hadoop





Hadoop Adoption Jumpstart

- Download & run the Sandbox
- http://www.youtube.com/hortonworks
 Sandbox tutorials have related videos on YouTube

POC on Sandbox

Easily transfer sample data onto Sandbox



Hadoop Proof of Concept

- Start small
 - For example, a 4-node cluster is a good start
- A POC cluster can be scaled and productionized
- Empower yourself & team with knowledge
- What more do I need to learn:
 - To use the Hadoop Ecosystem?
 - To *manage* the Hadoop Ecosystem?



Use Cases Run tutorials on Sandbox

Videos of tutorials youtube.com/hortonworks



WHAT'S YOUR PATH?

I'm a **Developer**

Interested in: Architecture & Fundementals MapReduce Programming Real-time Analytics

Developer Classes:

Developing Apps with Java > 4 Days

Data Analysis with Hive & Pig > 4 Days

Developing Solutions on Windows > 4 Days

Developing Custom YARN Applications > 2 Days

I'm a System Admin

Interested in: Cluster Monitoring Security & Governance Certification

Admin Classes:

Operations Management > 4 Days

÷

I'm a Data Analyst

Interested in: SQL & Scripting Languages Large Scale Data Sets Creating Value & Opportunity

Data Classes:

Applying Data Science with Hadoop > 2 Days

Data Analysis with Hive & Pig > 4 Days







hwuniversity@hortonworks.com

izhar@abyres.net

