# Micro Data Center
# &
# Hadoop  Big Data WareHouse

V0.01, r 2018

# Open Source Platform

- Micro Data Center 25 TB, Small Business Solution (Plug & Play)

- Hadoop Open Source Technology

- Hive Data Warehouse

- Hadoop Testing Data model

- Software & Tools Library

- Business Intelligence  report

# Infrastructure V0.07



PLUG&PLAY

Micro Data Center

ISP

Cisco Small Business Router

Wifi Cisco Small Business Router

Cisco Small Business Smart Switch

WD Sentinel Server - Xeon E3

Seagate Personal Cloud NAS Server
(Backup - Set 1)

Seagate Backup Plus External Hard Drive
(Backup - Set 2)
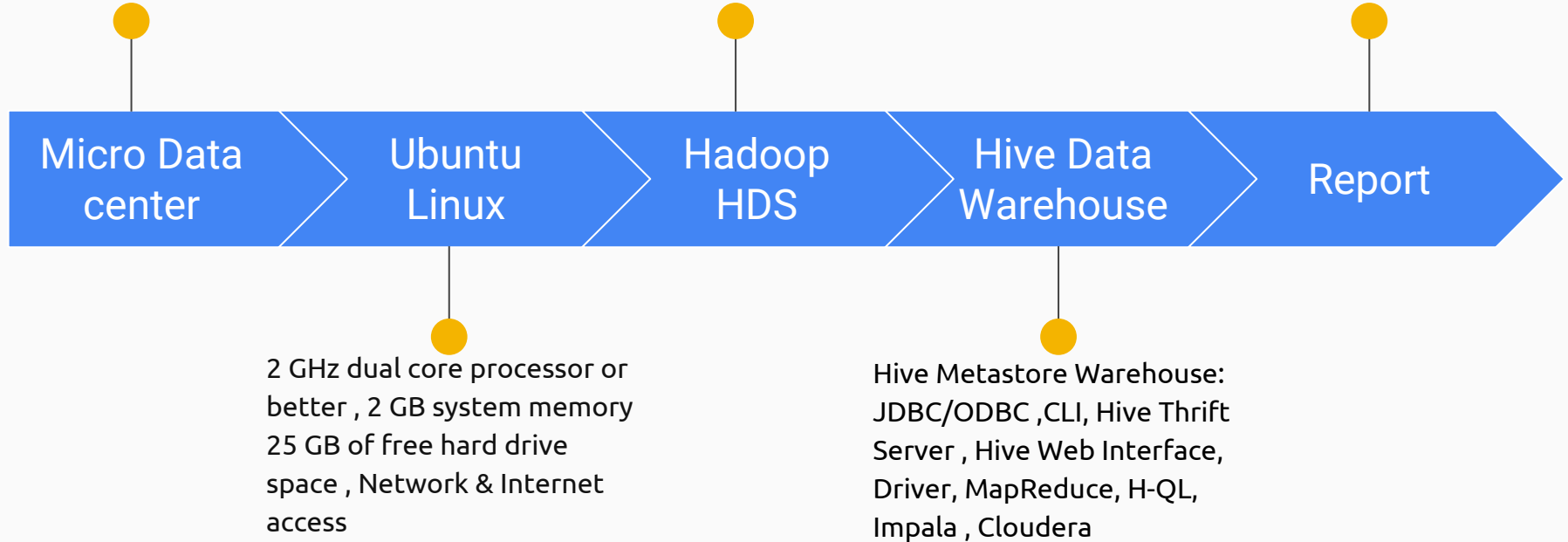
cloudera

amazon
web services™

ubuntu    Linux

HIVE

Impala

hadoop

The Data Warehouse

SQL

Excel

Full Redundant Feature : ISP
RV325Router ,RV340W
Wireless Router,SF200-24
Switch , NAS Server  Xeon
Storage, NAS Server, UPS

Java Package,Hadoop 2.7.3
Package,bash file (.bashrc),
NameNode, DataNode,
ResourceManager,NodeManager.

Business Intelligence, Excel,
Report Archiving, Backup &
Recovery, Cloud Storage

Micro Data
center

Ubuntu
Linux

Hadoop
HDS

Hive Data
Warehouse

Report

2 GHz dual core processor or
better , 2 GB system memory
25 GB of free hard drive
space , Network & Internet
access

Hive Metastore Warehouse:
JDBC/ODBC ,CLI, Hive Thrift
Server , Hive Web Interface,
Driver, MapReduce, H-QL,
Impala , Cloudera

# Micro Data Center (Hardware Specification)

**Networking :**
1. Cisco Small Business RV325 Router - 14-port - Gigabit Ethernet
2. Cisco Small Business RV340W Wireless Router - 2.4 GHz / 5 GHz
3. Cisco Small Business Smart SF200-24 Switch - 24 Ethernet Ports

**Central Storage :**
1. WD Sentinel DS5100 WDBYVE0080KBK Server  Xeon - 15 TB
2. Seagate Personal Cloud STCR3000101 NAS Server - 
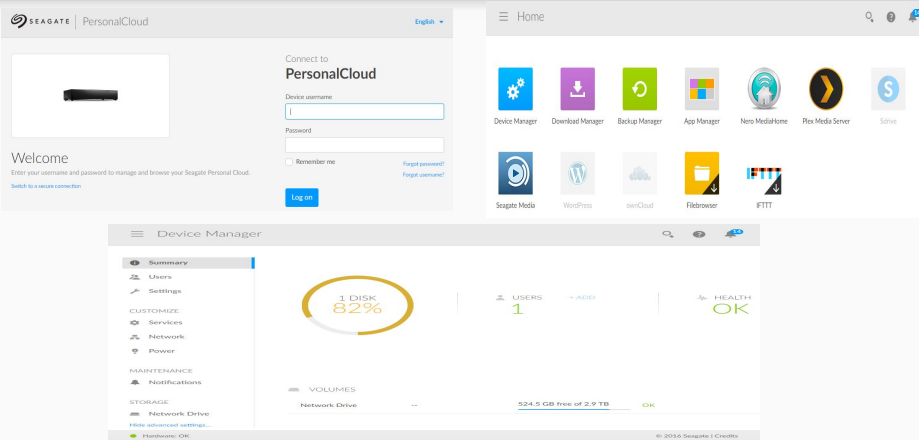3. Seagate 5TB Backup Plus External Hard Drive - 5TB

**Redundant Power UPS :**
1.  OL1000RTXL2U, Runtime @ 450 W: 20 min

# Network Connectivity

# Micro Data Center storage v0.07
# & Linux Ubuntu Workstation v18.04 LTS



## Connect to Micro Data Center  Storage:
1. Connect Network/Wifi router
2. PersonalCloude : http://192.168.1.82/
   Device user's name & PW
3. Network Configuration for Micro Datacenter Storage
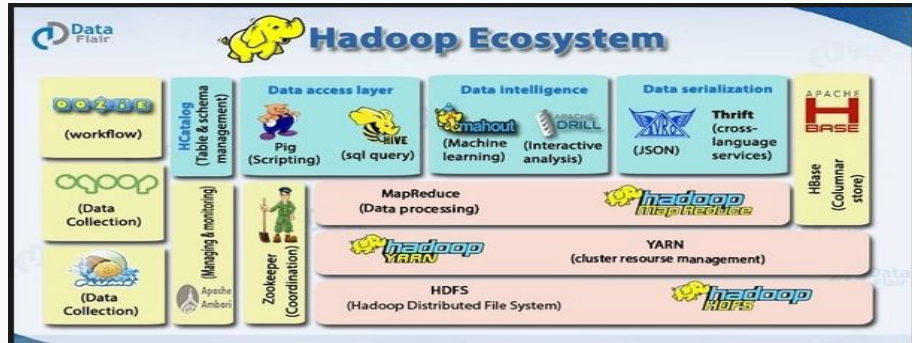
## Linux Ubuntu (Workstation) to Micro Data Center:
1. Boot from USB/DVD
2. Prepare to install Ubuntu
3. Allocate drive space
4. Begin installation
5. Login as User Admin for Storage

# Hadoop Installation v2.7.3



```
hduser@ubuntu:~/hadoop-2.8.2/sbin$ jps
3218 ResourceManager
3062 SecondaryNameNode
2648 NameNode
3357 NodeManager
2762 DataNode
3676 Jps
hduser@ubuntu:~/hadoop-2.8.2/sbin$ █
```



## Install Hadoop

**1. Download & Installation** :

a. Linux Ubuntu

b. Java JDK c. Vim CLI (command line interface) d.hadoop-2.6.5.tar.gz

**2. Group & Admin for Hadoop User**

**3. Configuration :** a. sysctl.conf (ipv6)  b. Generating public/private rsa key pair c. ssh localhost d..bashrc (Hadoop Variables), e. Hadoop Core conf files (hadoop-env.sh, core-site.xml,mapred-site.xml & hdfs-site.xml) f. Namenode, Datanode & hadoop_store g. Namenode format h. Start-all.sh jps

**4. Hadoop Daemons JPS** (ResorceManager, SecondaryNamenode, NodeManager, & Datanode)

# Hadoop Web Interface



- **http://localhost:50070/ of the NameNode daemon :**
  Namenode Summary report,Security , Safemode status, DFS Used%, DFS Remaining%, Block Pool Used, DataNodes usages%, Live Nodes, Dead Nodes, Decommissioning Nodes, Number of Under-Replicated Blocks, NameNode Journal Status, Journal Manager, NameNode Storage

- **Datanode Information :**
  Node, Admin State, Capacity, Used, Non DFS Used, Remaining, Block pool used, Failed Volumes.

- **Browsing HDFS :**
  Browse Directory, Permission, Owner, Group, Size, Block Size, Folder Name.

# Hive
# Data Warehouse Implementation v2.0



- Data warehouse built on top of Hadoop
- Provides an SQL like interface to analyze data
- An open source project under apache
- Works on high throughput and high latency principle (same as Hadoop)
- Ability to plug-in custom Map Reduce programs
- Mainly targeted for structured data
- Hides Map Reduce program complexities to end user

Step1. **Hive Installation**
• Download the Hive
• Configure  ~/.bashrc and set the environment variables

Step2. **Hive Warehouse Directory Creation**
• Hive is based on Hadoop platform in Hadoop in PATH
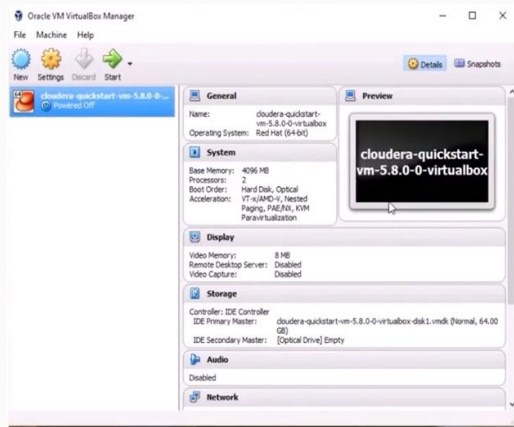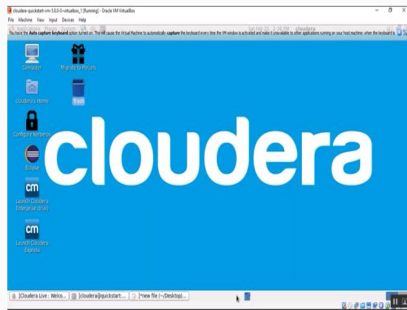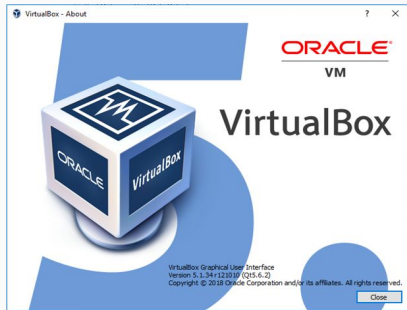• HDFS create the Hive Warehouse Directory

Step3. **Hive Configuration**
• Configure Hive with Hadoop
• Congigure "hive-env.sh" file
• Configure to external database to configure Metastore

Step4. **Hive Data Warehouse Files Location**
• $hadoop fs-ls /user/hive/warehouse

# VirtualBox Installation v5.1.34r





**R&D platform (Cloudera QuickStart)**

- Download from Oracle Virtual Box.org
- Configure Network Interfaces
- Open VM virtualBox Manager
- Appliance to import
- Appliance Settings : Name, Guest OS Type, CPU 2, RAM 10GB, DVD, Network Adapter
- Enable Network :
- Adapter 1: Inter PRO/1000 MT Desktop (NAT)
- Adapter 2:  Configure Host only Adapter ( VirtualBox Host-only Ethernet Adapter 2)
- System : Motherboard, Base Memory :10GB ,Processor 2 CPU
- VirtualBox running : Booting CentOS 6 (2.6.32-573.e16.x86_64)

# Tools & Software Library



**WINSCP Ftp infterace between Window & Linux**
- SSH and SPC code based on Putty
- Login : New Site > File protocol :FSTP, Host name IP : 192.168.56.101, Port number:22 and Username/PW.
- Open Two different OS windows with Window OS:c:\Users\Document and Linux OS: /home/cloudera/
- Upload & Download file : File Upload to Linux/Window OS and File Download to Windows/Linux OS

**PuTTY Key Generator :**
- Private key file for authentication
- Public key for pasting into OpenSSH authorized File,
- Type of Key Parameters RSA, Save public key

**PuTTY release 0.70** :
- Host Name IP : 3.17.0.143 & Port 22
- SSH authentication : Private key file for authentication
- Controlling session logging : Open remote terminal 3.17.0.143-Putty

**ETL (extract, transform, load) ELT (extract, load, transform) :** SQOOP

# Cloudera CDH5.3.

## Business Data Testing & Analysis 12k+ Customers

Hadoop, Hive & Impala (SQL) , Source Cloudera@quickstart

### MySQL (retail_db)

mysql> show tables;

| Table | Records |
|---|---|
| categories | 58 |
| customers | 12,435 |
| Departments | 6 |
| order_items | 68,883 |
| orders | 1,72,198 |
| products | 1,345 |

### BigData/Hive/Impala

hive> show tables;

| Table | Records |
|---|---|
| categories | 58 |
| customers | 12,435 |
| Departments | 6 |
| order_items | 68,883 |
| orders | 1,72,198 |
| products | 1,345 |

MySQL (retail_db) :
mysql> show databases;
mysql> use retail_db;
mysql> select count(*) from customers;

Sqoop :
[cloudera@quickstart ~]$ sqoop import-all-tables \

Hive (retail_db) :
hive> show databases;
hive> use default;
hive> show tables;
hive> select count(*) from customers;

Hive Data Warehouse :
[cloudera@quickstart ~]$  hadoop fs -ls /user/hive/warehouse/

Business Intelligence Report :
Most popular product categories
Top 10 revenue generating products

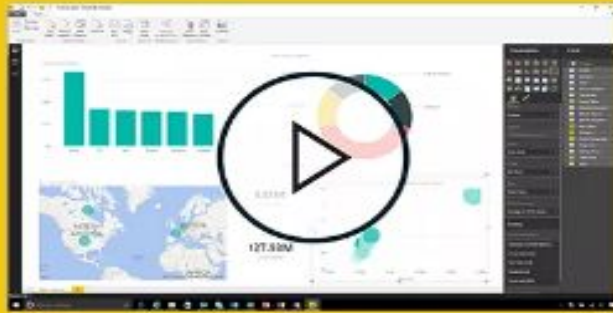# AWS Cloud Services Ubuntu Server 18.04 LTS



AWS Management Console:

- Step 1: Amazon Machine Image Ubuntu Server 18.04 LTS
- Step 2: Build an Instance
- Step 3: Configure Instance Details
- Step 4: Add Storage
- Step 5: Add Tags
- Step 6: Configure Security Group
- Step 7: Review Instance Launch

Connect AWS  Management Console:

- Connect ubuntu@ip-172.47.106:~$
- Generate Private Key by PPuttygen
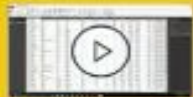- Connect AWS from Putty FTP
- File Transfer by WinSCP SFTP

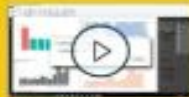# Business Intelligence Reporting Tools v2.65



Getting started with Power BI Desktop

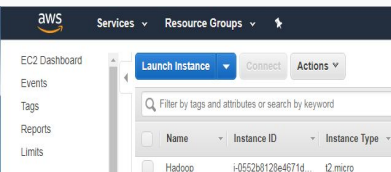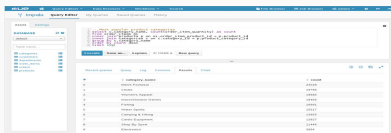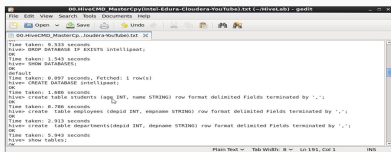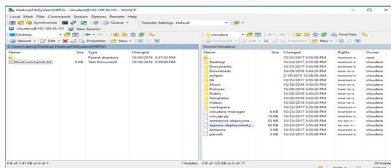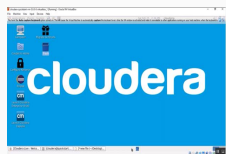Building reports  Query view concepts  Uploading your reports

## Report Generation

- Opening BI window
- Run Power BI Desktop
- Import Data from different source
- Connecting Dataset
- Load Data into BI
- Management Data as per Query
- Export as BI / Export to PDF

# Prototype Demo Cloudera (Remote Login)



Connected over Cloud Through :
- Team Viewer
- Windows10
- VirualBox
- Start Cloudera Desktop
- Cloudera CLI Terminal
- Run Mysql database
- Run HIVE open source database
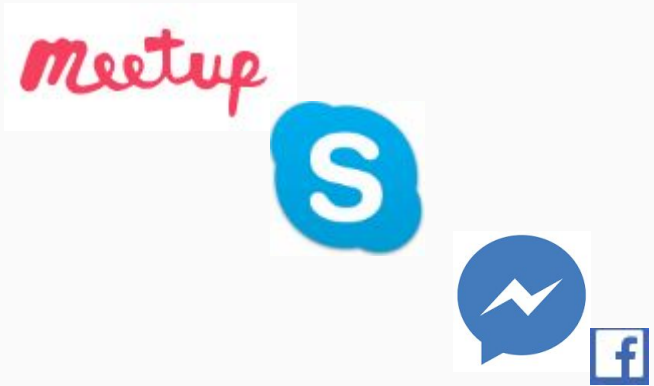- Cloudera QuickStart Hive/Impala SQL terminal
- SQL Data Analysis

FTP (File Transfer Protocol)
- Run WinSCP
- Connect Window Desktop to Linux Desktop

Business Intelligence Report
- Visual Analytics at your fingertips and creating interactive data visualizations and reports.

# Meetup

**Registration :**

- Free Orientation

- Prototype demo

- Consultancy

**Info & Registration :**
**Micro DataCenter & Data Warehouse**
**MDCDWH@gmail.com**
**https://goo.gl/forms/SuCTolEeZNNIL35V2**