

Getting Started with HDP Sandbox

Loading Sensor Data into HDFS

Ready to Get Started? [Download Sandbox](#)

Getting Started with HDP Sandbox

OVERVIEW

1. Concepts

2. Loading Sensor Data into HDFS

3. Hive - Data ETL

4. Spark - Risk Factor

5. Data Reporting With Zeppelin

NOTICE

As of January 31, 2021, this tutorial references legacy products that no longer represent Cludera's current product offerings.

Please visit recommended tutorials:

- [How to Create a CDP Private Cloud Base Development Cluster](#)
- All [Cludera Data Platform \(CDP\)](#) related tutorials

Introduction

In this section, you will download the sensor data and load that into HDFS using Ambari User Views. You will get introduced to the Ambari Files User View to manage files. You can perform tasks like create directories, navigate file systems and upload files to HDFS. In addition, you'll perform a few other file-related tasks as well. Once you get the basics, you will create two directories and then load two files into HDFS using the Ambari Files User View.

Prerequisites

The tutorial is a part of series of hands on tutorial to get you started on HDP using Hortonworks sandbox. Please ensure you complete the prerequisites before proceeding with this tutorial.

- Downloaded and deployed the [Hortonworks Data Platform \(HDP\)](#) Sandbox
- [Learning the Ropes of the HDP Sandbox](#)

Outline

- [HDFS backdrop](#)
- [Download and Extract Sensor Data Files](#)
- [Load the Sensor Data into HDFS](#)
- [Summary](#)
- [Further Reading](#)

HDFS backdrop

A single physical machine gets saturated with its storage capacity as the data grows. This growth drives the need to partition your data across separate machines. This type of File system that manages storage of data across a network of machines is called Distributed File Systems. **HDFS** is a core component of Apache Hadoop and is designed to store large files with streaming data access patterns, running on clusters of commodity hardware. With Hortonworks Data Platform (HDP), HDFS is now expanded to support [heterogeneous storage](#) media within the HDFS cluster.

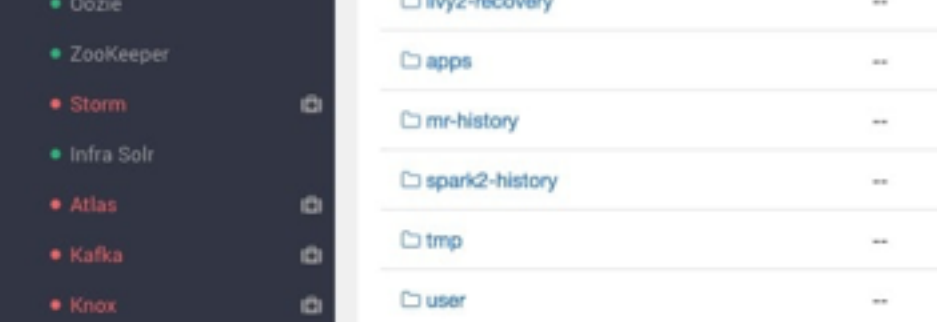
Download and Extract Sensor Data Files

- Download the sample sensor data contained in a compressed (.zip) folder here: [Geolocation.zip](#)
- Save the **Geolocation.zip** file to your computer, then extract the files. You should see a Geolocation folder that contains the following files:

- geolocation.csv** – This is the collected geolocation data from the trucks. It contains records showing *truck location, date, time, type of event, speed, etc.*
- trucks.csv** – This is data was exported from a relational database and it shows information on *truck models, driverid, truckid, and aggregated mileage info.*

Load the Sensor Data into HDFS

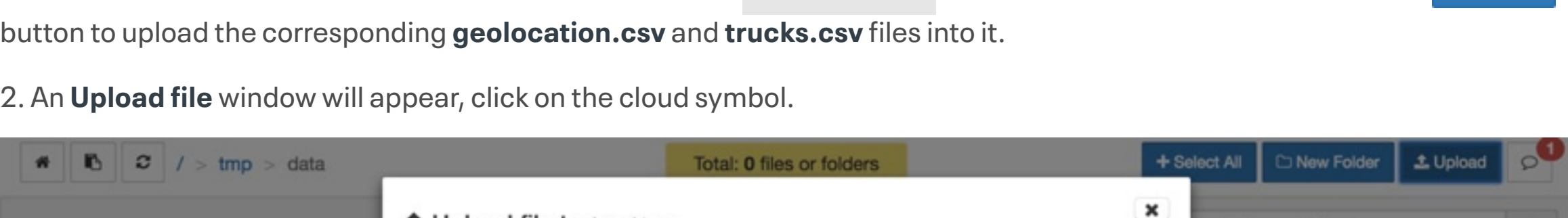
- Logon to Ambari using: **maria_dev/maria_dev**
- Go to Ambari Dashboard and open **Files View**.



- Start from the top root of the HDFS file system, you will see all the files the logged in user (**maria_dev** in this case) has access to see:

Name	Size	Last Modified	Owner	Group	Permission	Erasure Coding	Encrypted
mapred	--	2018-09-20 08:09	mapred	hdfs	drwxr-xr-x	No	
app-logs	--	2018-09-20 08:43	yam	hadoop	drwxrwxr-x	No	
ats	--	2018-09-20 08:09	yam	hadoop	drwxr-xr-x	No	
ats2	--	2018-09-20 08:10	hdfs	hdfs	drwxr-xr-x	No	
hdp	--	2018-09-20 08:09	hdfs	hdfs	drwxr-xr-x	No	
hdp2-recovery	--	2018-09-20 08:41	ivy	hdfs	drwx-----	No	
apps	--	2018-09-20 09:43	hdfs	hdfs	drwxr-xr-x	No	
re-history	--	2018-09-20 08:10	mapred	hadoop	drwxrwxr-x	No	
spark2-history	--	2018-09-20 16:09	spark	hadoop	drwxrwxr-x	No	
tmp	--	2018-09-20 09:43	hdfs	hdfs	drwxrwxr-x	No	
user	--	2018-09-20 09:46	hdfs	hdfs	drwxr-xr-x	No	
warehouse	--	2018-09-20 08:35	hdfs	hdfs	drwxr-xr-x	No	

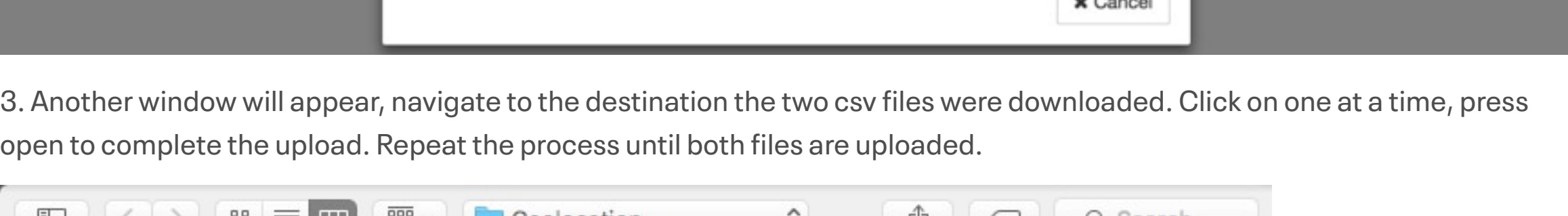
- Navigate to `/tmp/` directory by clicking on the directory links.
- Create directory `data`. Click the **New Folder** button to create that directory. Then navigate to it. The directory path you should see: `/tmp/data`



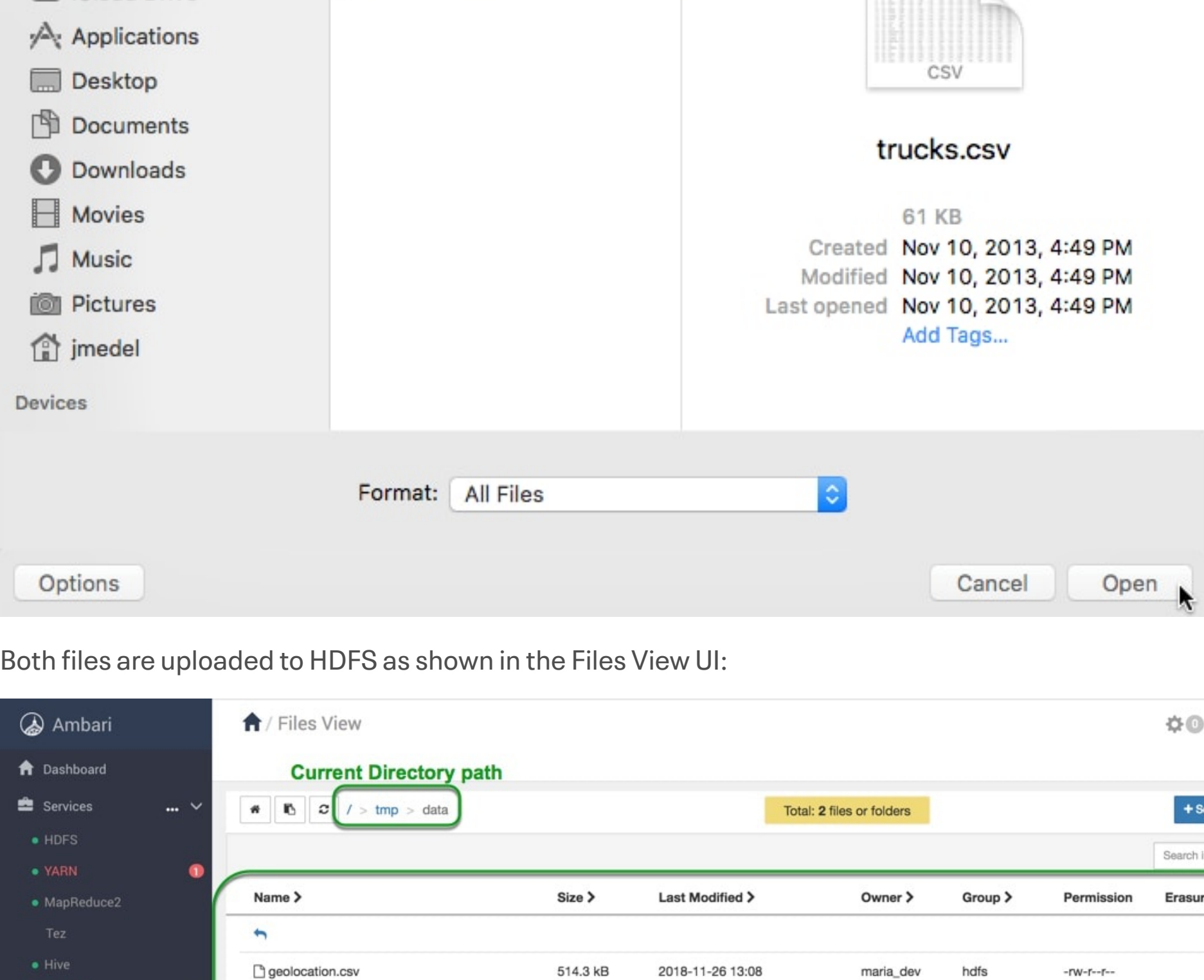
Upload Geolocation and Trucks CSV Files to data Folder

- If you're not already in your newly created directory path `/tmp/data`, go to the **data** folder. Then click on the **Upload** button to upload the corresponding **geolocation.csv** and **trucks.csv** files into it.

- An **Upload file** window will appear, click on the cloud symbol.



- Another window will appear, navigate to the destination the two csv files were downloaded. Click on one at a time, press open to complete the upload. Repeat the process until both files are uploaded.



Both files are uploaded to HDFS as shown in the Files View UI:

Current Directory path							
/tmp/data							
Name	Size	Last Modified	Owner	Group	Permission	Erasure Coding	Encrypted
geolocation.csv	514.3 kB	2018-11-26 13:08	maria_dev	hdfs	-rw-r--r--	No	
trucks.csv	59.9 kB	2018-11-26 13:08	maria_dev	hdfs	-rw-r--r--	No	

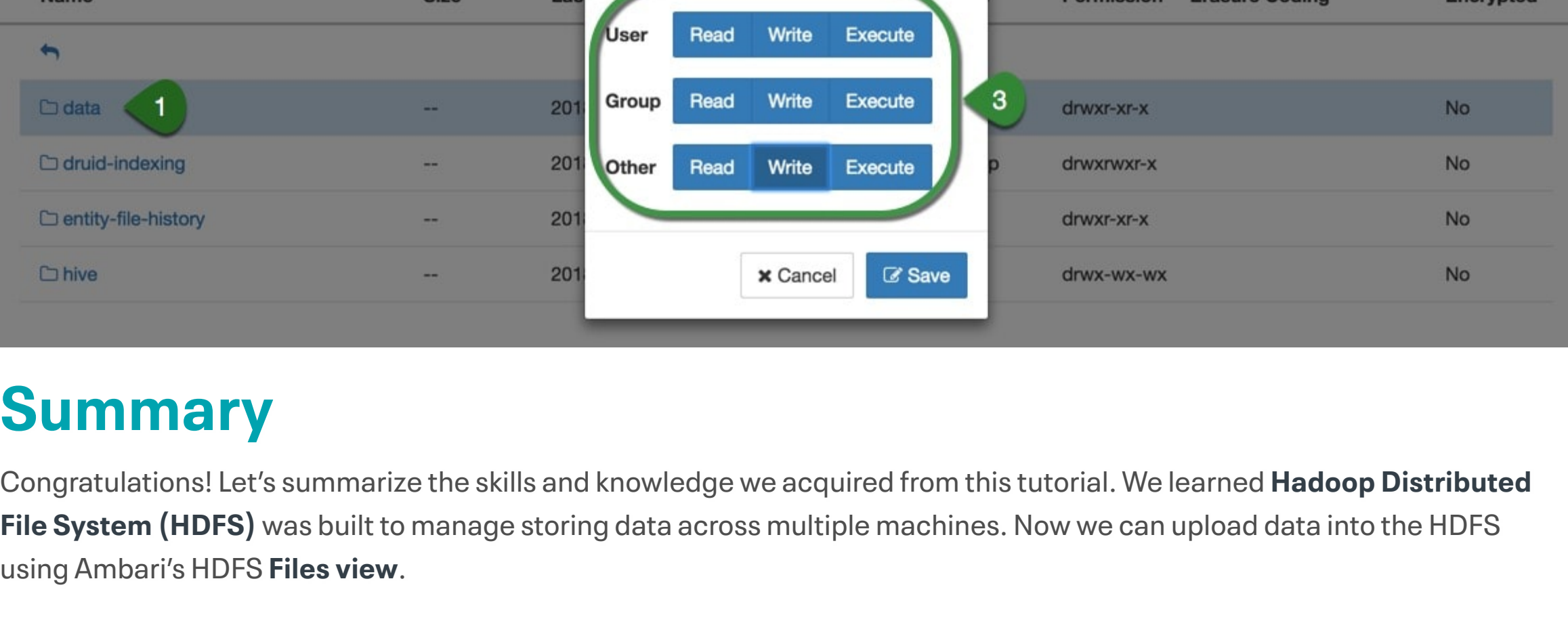
Newly uploaded files

You can also perform the following operations on a file or folder by clicking on the entity's row: **Open**, **Rename**, **Permissions**, **Delete**, **Copy**, **Move**, **Download** and **Concatenate**.

Set Write Permissions to Write to data Folder

- click on the `data` folder's row, which is contained within the directory path `/tmp/`.
- Click **Permissions**.
- Make sure that the background of all the **write** boxes are checked (**blue**).

Refer to image for a visual explanation.



Summary

Congratulations! Let's summarize the skills and knowledge we acquired from this tutorial. We learned **Hadoop Distributed File System (HDFS)** was built to manage storing data across multiple machines. Now we can upload data into the HDFS using Ambari's HDFS **Files view**.

Further Reading

- [HDFS](#)
- [HDFS User Guide](#)
- [HDFS Architecture Guide](#)
- [HDP OPERATIONS: HADOOP ADMINISTRATION](#)

[Previous](#)

[Next](#)

Contact Us

US: +1 888 789 1488
Outside the US: +1 650 362 0488

English

Company

About us
Careers
Diversity, Equality & Inclusion
Events
Leadership
Locations
Newsroom

Get started

Certification
Contact sales
Downloads
Find a partner
Find a solution
Training
Tutorials

Resources

Blog
CDP resources
CDP Trust Center
Community
Documentation
Resources library
Support