

Chapter 1 : How-to: Deploy Apache Hadoop Clusters Like a Boss - Cloudera Engineering Blog

This course is intended for students with some experience with Hadoop and MapReduce, Python, and bash commands. You'll have to be able to work with HDFS and write MapReduce programs. You can learn about these in our Intro to Hadoop and MapReduce course.

Consequently, if the drives become unmounted, the processes writing to these directories will not fill up the OS mount. HDFS is an immutable filesystem that was designed for large file sizes with long sequential reads. This goal plays well with stand-alone SATA drives, as they get the best performance with sequential reads. These parity bits have to be written and read during standard operations and add significant overhead. Setting the drives up in RAID-5 or RAID-6 arrays will create a single array or a couple very large arrays of mount points depending on the drive configuration. RAID arrays will also affect other systems that expect numerous mount points. Impala, for example, spins up a thread per spindle in the system, which will perform favorably in a JBOD environment vs. For the same reasons, configuring your Hadoop drives under LVM is neither necessary nor recommended. Deploying Heterogeneously Many customers purchase new hardware in regular cycles; adding new generations of computing resources makes sense as data volumes and workloads increase. For such environments containing heterogeneous disk, memory, or CPU configurations, Cloudera Manager allows Role Groups, which allow the administrator to specify memory, YARN containers, and Cgroup settings per node or per groups of nodes. While Hadoop can certainly run with mixed hardware specs, we recommend keeping worker-node configurations homogenous, if possible. In distributed computing environments, workloads are distributed amongst nodes and optimizing for local data access is preferred. Nodes configured with fewer computing resources can become a bottleneck, and running with a mixed hardware configuration could lead to a wider variation in SLA windows. There are a few things to consider: Mixed spindle configuration – HDFS block placement by default works in a round-robin fashion across all the directories specified by dfs. If you have, for example, a node with six 1. Understand the implications of deploying drives in this fashion in advance. Furthermore, if you deploy nodes with more overall storage, remember that HDFS balances by percentage. Mixed memory configuration – Mixing available memory in worker nodes can be problematic as it does require additional configuration. It is important to be cognizant of the points above but remember that Cloudera Manager can help with allocating resources to different hosts; allowing you to easily manage and optimize your configuration. Cloudera Manager offers many valuable features to make life much easier. The Cloudera Manager documentation is pretty clear on this but in order to stamp out any ambiguity, below are the high-level steps to do a production-ready Hadoop deployment with Cloudera Manager. Set up an external database and pre-create the schemas needed for your deployment.

Chapter 2 : Enterprise Data Catalog: deployment and sizing

Deploy a Hadoop cluster in a cloud environment Approach This booklet is a step by step educational chocked with useful examples with a view to enable you to construct and deal with a Hadoop cluster in addition to its intricacies.

Microsoft will provide an optimized experience for Hadoop running on Azure and on Windows Servers. Read the official announcement for more information. This blog post has more information on the Hadoop announcement. In this post I will demonstrate how to create a typical cluster with a Name Node, a Job Tracker and a customizable number of Slaves. You will also be able to dynamically change the number of Slaves using the Azure Management Portal. I will save the explanation of the mechanics for another post. Follow these steps to create an Azure package for your Hadoop cluster: Download all dependencies This Visual Studio project is pre-configured with Roles for each Hadoop component. The cluster configuration templates. Install the latest Azure SDK. As of this writing the latest version was 1. I used version 0. Hadoop is distributed in a tar. You can use 7-zip for the task. Now install Cygwin and package it in a single ZIP file. There is an on-going effort to remove this dependency for Hadoop 0. Just run the Cygwin install and accept all defaults. You should end up with Cygwin installed in c: Create a compressed folder of c: If you have your JVM installed under C: Configure your cluster The cluster-config. You will find the familiar [core hdfs mapped]-site. Make sure to only add and not change any of the properties. Create a new cluster-config. Upload all dependencies to your Azure Storage account Create a container called bin and upload all zip files to it. You should end up with these files in the bin container Configure your Azure Deployment Unzip the Visual Studio project. You can either use Visual Studio from here or update the required files using any text editor. I included a batch file to package the deployment if you are going down the command line route. The projects that require this file have links to it. Get an access key from your storage account and construct a connection string then paste it in the first line replacing [your connection string]. An Azure connection string has this format: If you are ok with that configuration skip to the next step. Deploy your cluster Create a new service to host your Hadoop cluster. The project is pre-configured for remote access to the machines in the cluster. The certificate password is hadoop. If you are using Visual Studio you can deploy by right-clicking on the Cloud project and selecting deploy. If you are not just run the buildPackage. You will get the Azure package Hadoop. Deploy your service you can ignore the warning message. Wait for it to complete, you should see something like this: Using your Hadoop cluster Now that everything is up and running you can navigate to the Name Node Summary page. The URL is http: As it is configured right now you must log in to the Job Tracker to start a new job. I will present alternatives in a future post hint Azure Connect. The username is hadoop and the password is H1Dooop. After you log in open a command prompt window and execute the following commands: The syntax is the same as the regular hadoop scripts. Congratulations, you just ran your first Hadoop job in Azure! What can I do with my Hadoop cluster? The cluster is fully operational. You can run any job you would like. You can also use the Azure Management Portal to dynamically change the number of Slaves. Hadoop will discover the new node s or find out nodes were removed and reconfigure the cluster accordingly. I added an extra Slave node And my cluster changed to If you used Hadoop in production you know you must take extra steps to prepare the Name Node, mostly around high availability. Azure Drive is probably another piece. Let me know your experiences using Hadoop in Azure.

Big Data is the hottest trend in the IT industry at the moment. Companies are realizing the value of collecting, retaining, and analyzing as much data as possible. They are therefore rushing to implement the next generation of data platform, and Hadoop is the centerpiece of these platforms. This.

I am very passionate about teaching, primarily helping developers understand, search, and bake data. And at the moment, less than. Imagine the possibilities of what you can discover with the help of baked data. Some of the major topics that we will cover include preparing the prerequisites in AWS to deploy Hadoop. The cloud has many features, but there is only a small subset that we need to know. Planning required before deploying, this includes security, capacity planning, and understanding best practices for the different workload types. Then, we will deploy CDH manually, a similar process to deploying on-prem, but I will highlight the different steps. By the end of this course, you will be prepared to take your baked data to the cloud, taking advantage of the flexibility and power that AWS has to offer. Of all the thousands of things available in AWS, there are about 10 things that you really need to know, and we will cover them in this module. It is time to continue our journey with understanding the cloud, an AWS mini crash course, part two of two. Here are the topics that I would like to cover in this module: Always remember to take security very seriously unless you want to be in the news in a very bad way, or even in a court of law. Capacity Planning, or how to size your cluster properly. Then, we will talk about Architectural Best Practices. Basically, how to get it done right, and finally, Preparing Cluster Deployment, namely, the prerequisites to start setting up your cluster. First, we decided how to install the cluster, more on this soon. Then, we deploy Cloudera Manager, and deploy the agents. In previous modules and trainings, we first learned how to set up the prerequisites to deploy a Hadoop cluster with CDH and AWS, and then we took the necessary steps to get our cluster up and running. And with so many steps, it is easy to make a mistake. And we want to avoid making mistakes. We want our clusters to be deployed with the precision of a hand-made expensive watch. Automation is the way to go, and we do it using Cloudera Director. Cloudera Director makes it simple, like an easy button for provisioning, managing, and de-provisioning one or many clusters in a predictable and efficient way. Director starts by deploying Cloudera Manager. It takes care of the prerequisites too, and then deploys one or many clusters, taking advantage of parallelism, making the whole process fast and efficient. With the installation paths, we have a set of instructions that you can follow. You complete them and manually deploy a cluster, and manually perform any steps necessary to modify the cluster. Since you are doing a lot of steps, it is easier to make a mistake. On the other hand with Cloudera Director, you get a tool to manage your cloud infrastructure. Best of all, it works with multiple cloud providers.

Chapter 4 : Deploying HDFS on a Cluster | x | Cloudera Documentation

Deploy a Hadoop cluster in a cloud environment Approach This book is a step-by-step tutorial filled with practical examples which will show you how to build and manage a Hadoop cluster along with its intricacies.

Choose the version of HDInsight for this cluster. For more information, see Supported HDInsight versions. This package provides option to have a more secure cluster setup by using Apache Ranger and integrating with Azure Active Directory. The default user name is admin. It uses the basic configuration on the Azure portal. Sometimes it is called "Cluster user. Used to connect to the cluster through SSH. Multiple users can be created using the Enterprise security package. The cluster is in the same location as the default storage. For a list of supported regions, click the Region drop-down list on HDInsight pricing. Storage endpoints for clusters Although an on-premises installation of Hadoop uses the Hadoop Distributed File System HDFS for storage on the cluster, in the cloud you use storage endpoints connected to cluster. Warning Using an additional storage account in a different location from the HDInsight cluster is not supported. During configuration, for the default storage endpoint you specify a blob container of an Azure Storage account or a Data Lake Store. The default storage contains application and system logs. Optionally, you can specify additional linked Azure Storage accounts and Data Lake Store accounts that the cluster can access. The HDInsight cluster and the dependent storage accounts must be in the same Azure location. Note The Secure transfer required feature enforces all requests to your account through a secure connection. This feature is only supported by HDInsight cluster version 3. Optional metastores You can create optional Hive or Oozie metastores. This can cause the cluster creation process to fail. Hive metastore If you want to retain your Hive tables after you delete an HDInsight cluster, use a custom metastore. You can then attach the metastore to another HDInsight cluster. Oozie metastore To increase performance when using Oozie, use a custom metastore. A metastore can also provide access to Oozie job data after you delete your cluster. Important You cannot reuse a custom Oozie metastore.

Chapter 5 : Hadoop using YARN & Dremio

Apache Ambari started as a sub-project of Hadoop but currently it enjoys the distinction of being a top-level Apache project. Due to the increasing size and complexities of Hadoop clusters with each passing day, the management of these Hadoop frameworks becomes a highly challenging task. This is.

Chapter 6 : Hadoop Cluster Deployment | PACKT Books

information Article Hadoop Cluster Deployment: A Methodological Approach Ronaldo Celso Messias Correia 1, Gabriel Spadon 2 ID, Pedro Henrique De Andrade Gomes 1, Danilo Medeiros Eler 1, ID, Rog rio Eduardo Garcia 1 ID and Celso Olivete Junior 1.*

Chapter 7 : Information | Free Full-Text | Hadoop Cluster Deployment: A Methodological Approach

Read "Hadoop Cluster Deployment" by Danil Zburivsky with Rakuten Kobo. This book is a step-by-step tutorial filled with practical examples which will show you how to build and manage a Hadoop.

Chapter 8 : Hadoop in Azure – Distributed Development

Welcome to the Hadoop Deployment Manual for Bright Cluster Manager About This Manual This manual is aimed at helping cluster administrators install, understand, configure, and manage the.

Chapter 9 : Cluster setup for Hadoop, Spark, Kafka, HBase, or R Server - Azure HDInsight | Microsoft Doc

DOWNLOAD PDF HADOOP CLUSTER DEPLOYMENT

HDInsight version. Choose the version of HDInsight for this cluster. For more information, see Supported HDInsight versions.. Enterprise security package. For Hadoop, Spark, and Interactive Query cluster types, you can choose to enable the Enterprise Security Package.