

Hortonworks DataFlow

Planning Your Deployment

(February 28, 2018)

Hortonworks DataFlow: Planning Your Deployment

Copyright © 2012-2018 Hortonworks, Inc. Some rights reserved.



Except where otherwise noted, this document is licensed under
Creative Commons Attribution ShareAlike 4.0 License.
<http://creativecommons.org/licenses/by-sa/4.0/legalcode>

Table of Contents

1. Deployment Scenarios	1
2. HDF Cluster Types and Recommendations	3
3. Production Cluster Guidelines	4
4. Hardware Sizing Recommendations	5
5. Navigating the HDF Library	7

List of Tables

1.1. Installation Scenarios 1
1.2. Upgrade Scenarios 2

1. Deployment Scenarios

Your deployment scenario for installing, configuring, or upgrading your Hortonworks DataFlow (HDF) components depends on your particular use case.



- [Table 1.1, "Installation Scenarios" \[1\]](#)
- [Table 1.2, "Upgrade Scenarios" \[2\]](#)



Note

This document describes how to install Ambari and HDF on Intel x86 hardware. To install Ambari and HDF on IBM Power Systems, review your deployment options using [Planning Your Deployment for IBM Power Systems](#).

Table 1.1. Installation Scenarios

Scenario	Installation Scenario	Steps
Installing an HDF Cluster	<p>This scenario applies if you want to install the entire HDF platform, consisting of all flow management and stream processing components on a new cluster.</p> <p>The stream processing components include the new Streaming Analytics Manager (SAM) modules that are in GA (General Availability). This includes the SAM Stream Builder and Stream Operations modules but does not include installing the technical preview version of SAM Stream Insight, which is powered by Druid and Superset.</p> <p>This scenario requires that you install an HDF cluster.</p>	<ol style="list-style-type: none"> 1. Install Ambari. 2. Install databases. 3. Install the HDF management pack. 4. Install an HDF cluster using Ambari.
Installing HDF Services on a New HDP Cluster	<p>This scenario applies to you if you are both an Hortonworks Data Platform (HDP) and HDF customer and you want to install a fresh cluster of HDP and add HDF services.</p> <p>The stream processing components include the new (SAM) and all of its modules. This includes installing the technical preview version of the SAM Stream Insight module, which is powered by Druid and Apache Superset.</p> <p>This scenario requires that you install both an HDF cluster and an HDP cluster.</p> <p> Important</p> <p>You cannot install SAM and Schema Registry for HDF 3.1 on an HDP 2.6.4 cluster, and you cannot upgrade these services from a previous HDP cluster.</p> <p> Important</p> <p>You cannot upgrade your HDF Storm and Kafka versions if they exist on an HDP cluster.</p>	<ol style="list-style-type: none"> 1. Install Ambari. 2. Install databases. 3. Install an HDP cluster using Ambari. 4. Install the HDF management pack. 5. Update the HDF base URL. 6. Add HDF services to an HDP cluster
Installing HDF Services on an Existing HDP Cluster	<p>You have an existing HDP cluster with Apache Storm and or Apache Kafka services and want to install Apache NiFi or NiFi Registry modules on that cluster.</p> <p>This requires that you upgrade to the latest version of Apache Ambari and HDP, and then use Ambari to add HDF services to the upgraded HDP cluster.</p>	<ol style="list-style-type: none"> 1. Upgrade Ambari 2. Upgrade HDP 3. Install Databases




Scenario	Installation Scenario	Steps
	 <p>Important</p> <p>You cannot install SAM and Schema Registry for HDF 3.1 on an HDP 2.6.4 cluster, and you cannot upgrade these services from a previous HDP cluster.</p>  <p>Important</p> <p>You cannot upgrade your HDF Storm and Kafka versions if they exist on an HDP cluster.</p>	<ol style="list-style-type: none"> 4. Install HDF Management Pack 5. Update HDF Base URL 6. Add HDF Services to HDP cluster
<p>Performing any of the previous deployments by using a local repository</p> <p>See <i>Using Local Repositories</i> in the instructions appropriate for your scenario.</p>	<p>Local repositories are frequently used in enterprise clusters that have limited outbound internet access. In these scenarios, having packages available locally provides more governance and better installation performance.</p> <p>This requires that you perform several steps to create a local repository and update the Ambari repository configuration file.</p>	<ol style="list-style-type: none"> 1. Obtain the public repositories. 2. Configure the local repository. 3. Prepare the Ambari repository configuration file.

Table 1.2. Upgrade Scenarios

Scenario	Upgrade Scenario	Steps
<p>Upgrading an HDF cluster</p>	<p>You have an existing Ambari-managed HDF cluster and want to upgrade it using either Rolling or Express upgrade.</p>  <p>Note</p> <p>Rolling Upgrade is not supported for NiFi. During the Rolling Upgrade, each NiFi instance is stopped, upgraded, and restarted.</p>	<ol style="list-style-type: none"> 1. Upgrade Ambari and your HDF management pack. 2. Upgrade HDF.
<p>Upgrading an HDP cluster with only NiFi installed</p>	<p>You have an existing Ambari-managed HDP cluster with NiFi installed. You want to upgrade the NiFi service installed on your HDP cluster.</p>	<ol style="list-style-type: none"> 1. Upgrade Ambari and your HDF management pack. 2. Upgrade NiFi.
<p>Upgrading an HDP cluster with all HDF services installed</p>	<p>You have an existing Ambari-managed HDF cluster with all your HDF services installed and you want to upgrade the NiFi service, and migrate SAM and Schema Registry to an HDF cluster.</p>	<ol style="list-style-type: none"> 1. Upgrade Ambari and your HDF management pack. 2. Upgrade NiFi. 3. Migrate SAM and Schema Registry to an HDF cluster.

2. HDF Cluster Types and Recommendations

Cluster Type	Description	Number of VMs or Nodes	Node Specification	Network
Single VM HDF Sandbox	Evaluate HDF on local machine. Not recommended to deploy anything but simple applications.	1 VM	At least 4 GB RAM	
Evaluation Cluster	Evaluate HDF in a clustered environment. Used to evaluate HDF for simple data flows and streaming applications.	3 VMs or no des	<ul style="list-style-type: none"> • 16 GB of RAM • 8 cores/vCores 	
Small Development Cluster	Use this cluster in development environments.	6 VM s/No des	<ul style="list-style-type: none"> • 16 GB of RAM • 8 cores or vCores 	
Medium QE Cluster	Use this cluster in QE environments.	8 VMs/Nodes	<ul style="list-style-type: none"> • 32 GB of RAM • 8 to16 cores or vCores 	
Small Production Cluster	Use this cluster in small production environments.	15 VMs/Nodes	<ul style="list-style-type: none"> • 64 - 128 GB of RAM • 8 - 16 cores of RAM 	1 GB Bonded Nic
Medium Production Cluster	Use this cluster in a medium production environment.	24 VMs/Nodes	<ul style="list-style-type: none"> • 64 - 128 GB of RAM • 8 - 16 cores of RAM 	10 GB bonded network interface card (NIC)
Large Production Cluster	Use this cluster in a large production environment.	32 VMs/Nodes	<ul style="list-style-type: none"> • 64 - 128 GB of RAM • 16 cores of RAM 	10 GB Bonded Nic

More Information

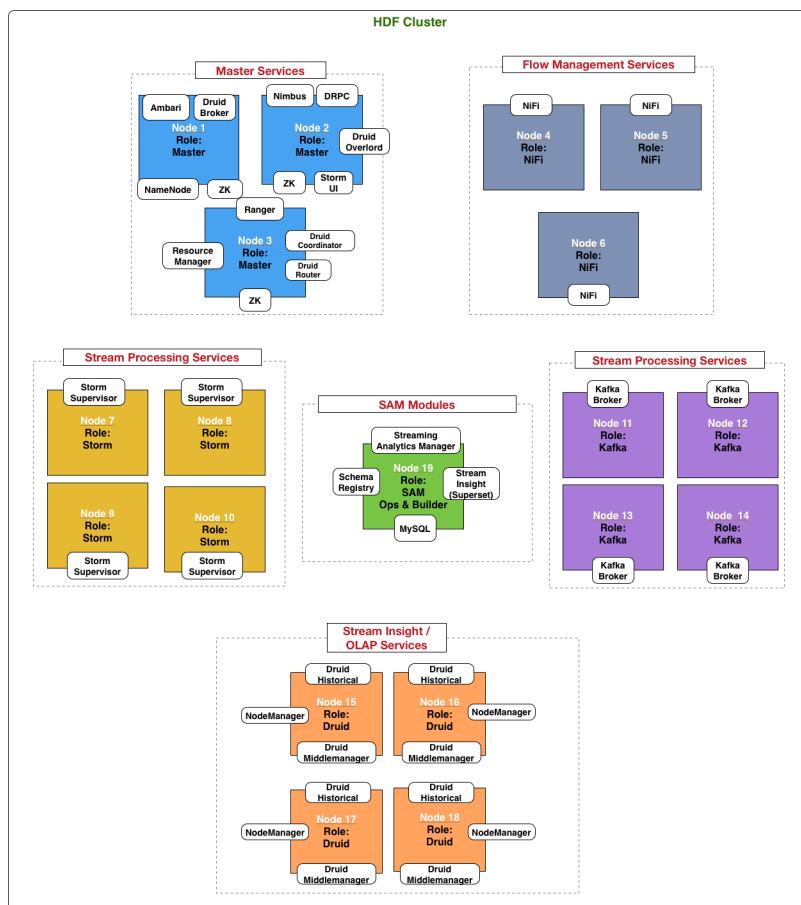
[Download the Sandbox](#)

3. Production Cluster Guidelines

General guidelines for production guidelines for service distribution :

- NiFi, Storm, and Kafka should not be located on the same node or virtual machine.
- NiFi, Storm, and Kafka must have a dedicated ZK cluster with at least three nodes.
- If the HDF SAM is being used in an HDP cluster, the SAM should not be installed on the same node as the Storm worker node.

The following diagram illustrates how services could be distributed for a small production cluster across 19 nodes:



4. Hardware Sizing Recommendations

Recommendations for Kafka

- Kafka Broker Node: eight cores, 64 GB to 128 GB of RAM, two or more 8-TB SAS/SSD disks, and a 10-GigE Nic.
- Minimum of three Kafka broker nodes
- Hardware Profile: More RAM and faster speed disks are better; 10-GigE Nic is ideal.
- 75 MB/sec per node is a conservative estimate (can go much higher if more RAM and reduced lag between writing/reading and therefore 10-GigE Nic is required).

With a minimum of three nodes in your cluster, you can expect 225 MB/sec data transfer.

You can perform additional further sizing by using the following formula:

$$\text{num_brokers} = \text{desired_throughput (MB/sec)} / 75$$

Recommendations for Storm

- Storm Worker Node: 8 core, 64 GB RAM, 1 GigE Nic
- Minimum of 3 Storm worker nodes
- Nimbus Node: Minimum 2 nimbus nodes, 4 core, 8 GB RAM
- Hardware profile: disk I/O is not that important; more cores are better.
- 50 MB/sec per node with low to moderate complexity topology reading from Kafka and no external lookups. Medium-complexity and high-complexity topologies might have reduced throughput.

With a minimum 2 nimbus, 2 worker cluster, you can expect to run 100 MB/sec of low to medium complexity topology.

Further sizing can be done as follows. Formula: $\text{num_worker_nodes} = \text{desired_throughput(MB/sec)} / 50$

Recommendations for NiFi

NiFi is designed to take advantage of:

- all the cores on a machine
- all the network capacity
- all the disk speed
- many gigabytes of RAM (although usually not all) on a system

Hence, it is important that NiFi be running on dedicated nodes. Following are the recommended server and sizing specifications for NiFi:

- Minimum of 3 nodes
- 8+ cores per node (more is better)
- 6+ disks per node (SSD or Spinning)
- At least 8 GB

If you want this sustained throughput...	Then provide this minimum hardware ...
50 MB and thousands of events per second	<ul style="list-style-type: none"> • 1 or 2 nodes • 8 or more cores per node, although more is better • 6 or more disks per node (SSD or spinning) • 2 GB memory per node • 1 GB bonded NICs
100 MB and tens of thousands of events per second	<ul style="list-style-type: none"> • 3 or 4 nodes • 16 or more cores per node, although more is better • 6 or more disks per node (SSD or spinning) • 2 GB of memory per node • 1 GB bonded NICs
200 MB and hundreds of thousands of events per second	<ul style="list-style-type: none"> • 5 to 7 nodes • 24 or more cores per node (effective CPUs) • 12 or more disks per node (SSD or spinning) • 4 GB of memory per node • 10 GB bonded NICs
400 to 500 MB/sec and hundreds of thousands of events per second	<ul style="list-style-type: none"> • 7 - 10 nodes • 24 or more cores per node (effective CPUs) • 12 or more disks per node (SSD or spinning) • 6 GB of memory per node • 10 GB bonded NICs

5. Navigating the HDF Library

To navigate the Hortonworks DataFlow (HDF) documentation library, begin by deciding your current goal.

If you want to...	See this document...
Install or upgrade an HDF cluster using Apache Ambari	<ul style="list-style-type: none">• Release Notes• Support Matrix• Planning Your Deployment• Ambari Upgrade• MiNiFi Java Agent Quick Start
Get started with HDF	<ul style="list-style-type: none">• Getting Started with Apache NiFi• Getting Started with Stream Analytics
Use and administer HDF Flow Management capabilities	<ul style="list-style-type: none">• Apache NiFi User Guide• Apache NiFi Administration Guide• Apache NiFi Developer Guide• Apache NiFi Expression Language Guide• MiNiFi Java Agent Administration Guide
Use and administer HDF Stream Analytics capabilities	<ul style="list-style-type: none">• Streaming Analytics Manager User Guide• Schema Registry User Guide• Apache Storm Component Guide• Apache Kafka Component Guide