

# Hortonworks DataFlow

## Planning Your Deployment

(August 9, 2017)

## Hortonworks DataFlow: Planning Your Deployment

Copyright © 2012-2017 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
  - 1.1. Identify your Deployment Scenarios ..... 1
  - 1.2. HDF Cluster Types and Recommendations ..... 1
  - 1.3. Production Cluster Guidelines ..... 2
  - 1.4. Hardware Sizing Recommendations ..... 3
- 2. Where to Go Next? ..... 6

# 1. Deployment Scenarios

## 1.1. Identify your Deployment Scenarios

Depending on your use case, your deployment scenario for installing and configuring HDF components is different. These scenarios are covered in the following table.

Scenario	Deployment Scenario	Scenario Steps
<a href="#">Installing HDF Services on a New HDP Cluster</a>	<p>This scenario applies to you if you are both an 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 Streaming Analytics Manager (SAM) and <b>all</b> of its modules. This includes installing the technical preview version of SAM's Stream Insight module which is powered by Druid and SuperSet.</p> <p>This requires that you install both an HDF and an HDP cluster.</p>	<ol style="list-style-type: none"> <li>1. Install Ambari</li> <li>2. Install Databases</li> <li>3. Install HDP Cluster using Ambari</li> <li>4. Install HDF Management Pack</li> <li>5. Update HDF Base URL</li> <li>6. Add HDF Services to HDP cluster</li> </ol>
<a href="#">Installing an HDF Cluster</a>	<p>You want to install the entire HDF platform consisting of all flow management and stream processing components on a <b>new</b> cluster.</p> <p>The stream processing components include the new Streaming Analytics Manager (SAM) modules that are <b>GA</b>. This includes the SAM's Stream Builder and Stream Operations modules but <b>does not</b> include installing the technical preview version of SAM's Stream Insight module which is powered by Druid and SuperSet.</p> <p>This requires that you install an HDF cluster.</p>	<ol style="list-style-type: none"> <li>1. Install Ambari</li> <li>2. Install Databases</li> <li>3. Install HDF Management Pack</li> <li>4. Install HDF cluster using Ambari</li> </ol>
<a href="#">Installing HDF Services on an Existing HDP Cluster</a>	<p>You have an <b>existing</b> HDP cluster with Storm and or Kafka services and want to install NiFi or SAM's modules on that cluster.</p> <p>This requires that you upgrade to the latest version of Ambari and HDP, and use Ambari to add HDF services to the upgraded HDP cluster.</p>	<ol style="list-style-type: none"> <li>1. Upgrade Ambari</li> <li>2. Upgrade HDP</li> <li>3. Install Databases</li> <li>4. Install HDF Management Pack</li> <li>5. Update HDF Base URL</li> <li>6. Add HDF Services to HDP cluster</li> </ol>
<p>Performing any of the above deployment scenarios using a local repository.</p> <p><i>See <a href="#">Using Local Repositories</a> in the instructions appropriate for your scenario.</i></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 prepare the Ambari repository configuration file.</p>	<ol style="list-style-type: none"> <li>1. Obtain the Public Repositories</li> <li>2. Set Up the Local Repository</li> <li>3. Prepare the Ambari Repository Configuration File</li> </ol>

## 1.2. HDF Cluster Types and Recommendations

Cluster Type	Description	Number of 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/Nodes	<ul style="list-style-type: none"> <li>• 16 GB of RAM</li> <li>• 8 cores/vCores</li> </ul>	
Small Development Cluster	Use this cluster in development environments.	6 VMs/Nodes	<ul style="list-style-type: none"> <li>• 16 GB of RAM</li> <li>• 8 cores/vCores</li> </ul>	
Medium QE Cluster	Use this cluster in QE environments.	8 VMs/Nodes	<ul style="list-style-type: none"> <li>• 32 GB of RAM</li> <li>• 8 - 16 cores/vCores</li> </ul>	
Small Production Cluster	Use this cluster in small production environments.	15 VMs/Nodes	<ul style="list-style-type: none"> <li>• 64 - 128 GB of RAM</li> <li>• 8 - 16 cores of RAM</li> </ul>	1 GB Bonded Nic
Medium Production Cluster	Use this cluster in a medium production environment.	24 VMs/Nodes	<ul style="list-style-type: none"> <li>• 64 - 128 GB of RAM</li> <li>• 8 - 16 cores of RAM</li> </ul>	10 GB Bonded Nic
Large Production Cluster	Use this cluster in a large production environment.	32 VMs/Nodes	<ul style="list-style-type: none"> <li>• 64 - 128 GB of RAM</li> <li>• 16 cores of RAM</li> </ul>	10 GB Bonded Nic

### More Information

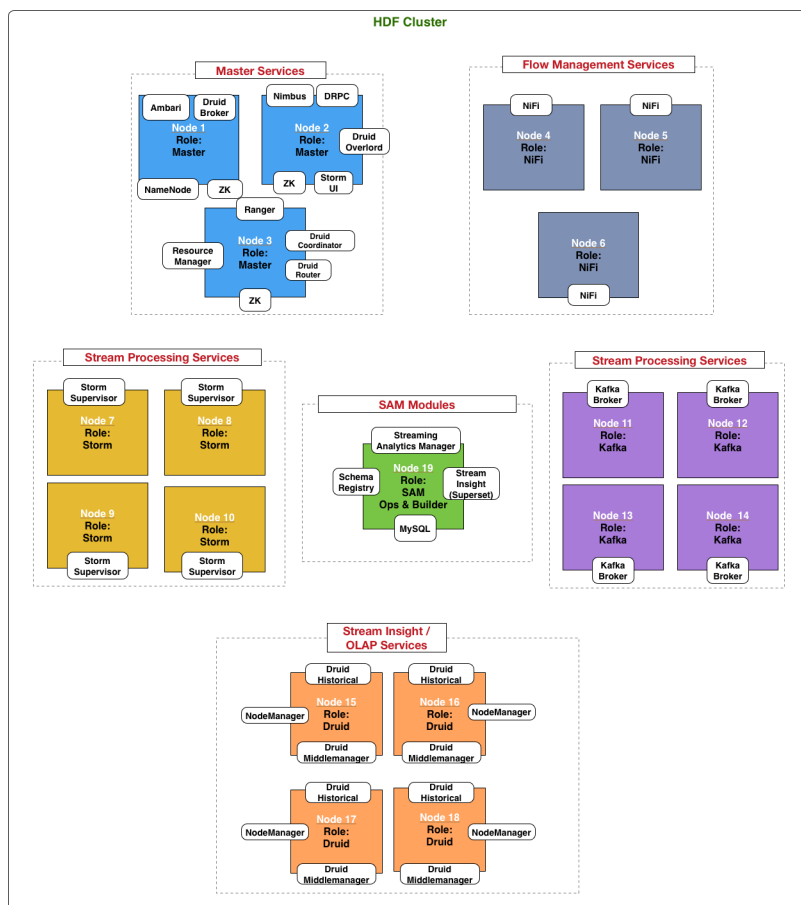
[Download the Sandbox](#)

## 1.3. Production Cluster Guidelines

General guidelines for production guidelines for service distribution:

- NiFi, Storm and Kafka should not be collocated on the same Node/VM.
- NiFi, Storm and Kafka have at least a dedicated 3 Node ZK cluster.
- If HDF's SAM is being used in an HDP cluster, SAM should not installed on the same node as Storm worker node.

The below diagram illustrates how services could be spread out for small production cluster across 15 nodes.



## 1.4. Hardware Sizing Recommendations

### Recommendations for Kafka

- Kafka Broker Node: 8 core, 64-128GB RAM, 2+ 8TB SAS/SSD disk, 10Gige Nic.
- Minimum of 3 Kafka Broker Nodes
- Hardware Profile: More RAM and faster speed disks are better, 10Gige 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 10GB Nic is required).

With a minimum 3 node cluster, you can expect 225 MB/second data transfer.

Further sizing can be done as follows. Formula:  $\text{num\_brokers} = \text{desired\_throughput}(\text{MB/sec}) / 75$

### 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 io 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 to high complexity topologies may see 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 GB of RAM (though usually not all) on a system

Hence is important that NiFi be running on dedicated nodes. The below are the recommended server and sizing specs 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 ...	Recommended hardware sizing ...
50 MB/second sustained throughput and thousands of events per second	<ul style="list-style-type: none"> <li>• 1 - 2 nodes</li> <li>• 8 or more cores per node, although more is better</li> <li>• 6 or more disks per node (solid state or spinning)</li> <li>• 2 GB memory per node</li> <li>• 1 GB bonded NICs</li> </ul>
100 MB/second sustained throughput and tens of thousands of events per second	<ul style="list-style-type: none"> <li>• 3 - 4 nodes</li> <li>• 8 or more cores per node, although more is better</li> <li>• 6 or more disks per node (solid state or spinning)</li> <li>• 2 GB of memory per node</li> <li>• 1GB bonded NICs</li> </ul>
200 MB/second sustained throughput and hundreds of thousands of events per second	<ul style="list-style-type: none"> <li>• 5 - 7 nodes</li> <li>• 24 or more cores per node (effective CPUs)</li> <li>• 12 or more disks per node (solid state or spinning)</li> </ul>

If you want ...	Recommended hardware sizing ...
	<ul style="list-style-type: none"><li>• 4 GB of memory per node</li><li>• 10 GB bonded NICs</li></ul>
400 - 500 MB/second sustained throughput and hundreds of thousands of events per second	<ul style="list-style-type: none"><li>• 7 - 10 nodes</li><li>• 24 or more cores per node (effective CPUs)</li><li>• 12 or more disks per node (solid state or spinning)</li><li>• 6 GB of memory per node</li><li>• 10 GB bonded NICs</li></ul>



## 2. Where to Go Next?

Use this table to help you navigate the HDF documentation library.

If you want to ...	See this document ...
Install or Upgrade an HDF cluster using Ambari	<ul style="list-style-type: none"><li>• <a href="#">Release Notes</a></li><li>• <a href="#">Support Matrices</a></li><li>• <a href="#">Planning Your Deployment</a></li><li>• <a href="#">Ambari Upgrade</a></li></ul>
Manually install or upgrade HDF components. This option is not available for Streaming Analytics Manager and Schema Registry.	<ul style="list-style-type: none"><li>• <a href="#">Command Line Installation</a></li><li>• <a href="#">MiNiFi Java Agent Quick Start</a></li><li>• <a href="#">Manual Upgrade</a></li></ul>
Get Started with HDF	<ul style="list-style-type: none"><li>• <a href="#">Getting Started with Apache NFi</a></li><li>• <a href="#">Getting Started with Stream Analytics</a></li></ul>
Use and administer HDF Flow Management capabilities	<ul style="list-style-type: none"><li>• <a href="#">Apache NiFi User Guide</a></li><li>• <a href="#">Apache NiFi Administration Guide</a></li><li>• <a href="#">Apache NiFi Developer Guide</a></li><li>• <a href="#">Apache NiFi Expression Language Guide</a></li><li>• <a href="#">MiNiFi Java Agent Administration Guide</a></li></ul>
Use and administer HDF Stream Analytics capabilities	<ul style="list-style-type: none"><li>• <a href="#">Streaming Analytics Manager User Guide</a></li><li>• <a href="#">Schema Registry User Guide</a></li><li>• <a href="#">Apache Storm Component Guide</a></li><li>• <a href="#">Apache Kafka Component Guide</a></li></ul>