



THE  
POWER  
TO KNOW.

# **SAS<sup>®</sup> High-Performance Analytics Infrastructure 2.94**

Installation and Configuration Guide

The correct bibliographic citation for this manual is as follows: SAS Institute Inc. 2015. *SAS® High-Performance Analytics Infrastructure 2.94: Installation and Configuration Guide*. Cary, NC: SAS Institute Inc.

**SAS® High-Performance Analytics Infrastructure 2.94: Installation and Configuration Guide**

Copyright © 2015, SAS Institute Inc., Cary, NC, USA

All rights reserved. Produced in the United States of America.

**For a hard-copy book:** No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

**For a web download or e-book:** Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a) and DFAR 227.7202-4 and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR 52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513-2414.

May 2015

SAS provides a complete selection of books and electronic products to help customers use SAS® software to its fullest potential. For more information about our offerings, visit [support.sas.com/bookstore](http://support.sas.com/bookstore) or call 1-800-727-3228.

SAS® and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.

---

# Contents

<i>What's New in Installation and Configuration for SAS High-Performance Analytics Infrastructure 2.94</i> . . . . .	<i>v</i>
<i>Accessibility</i> . . . . .	<i>vii</i>
<b>Chapter 1 • Introduction to Deploying the SAS High-Performance Analytics Infrastructure</b> . .	<b>1</b>
What Is Covered in This Document? . . . . .	1
Which Version Do I Use? . . . . .	2
Experimental Software . . . . .	2
What is the Infrastructure? . . . . .	2
Where Do I Locate My Analytics Cluster? . . . . .	4
Deploying the Infrastructure . . . . .	9
<b>Chapter 2 • Preparing Your System to Deploy the SAS High-Performance Analytics Infrastructure</b> . . . . .	<b>13</b>
Infrastructure Deployment Process Overview . . . . .	14
System Settings for the Infrastructure . . . . .	14
List the Machines in the Cluster or Appliance . . . . .	15
Review Passwordless Secure Shell Requirements . . . . .	16
Preparing for Kerberos . . . . .	16
Preparing to Install SAS High-Performance Computing Management Console . . . . .	24
Preparing to Deploy Hadoop . . . . .	25
Preparing to Deploy the SAS High-Performance Analytics Environment . . . . .	28
Recommended Database Names . . . . .	30
Pre-installation Ports Checklist for SAS . . . . .	30
<b>Chapter 3 • Deploying SAS High-Performance Computing Management Console</b> . . . . .	<b>33</b>
Infrastructure Deployment Process Overview . . . . .	33
Benefits of the Management Console . . . . .	34
Overview of Deploying the Management Console . . . . .	34
Installing the Management Console . . . . .	35
Configure the Management Console . . . . .	36
Create the Installer Account and Propagate the SSH Key . . . . .	38
Create the First User Account and Propagate the SSH Key . . . . .	41
<b>Chapter 4 • Deploying Co-Located Hadoop</b> . . . . .	<b>45</b>
Infrastructure Deployment Process Overview . . . . .	45
Overview of Deploying Hadoop . . . . .	46
Deploying SAS High-Performance Deployment of Hadoop . . . . .	47
Configuring Existing Hadoop Clusters . . . . .	60
<b>Chapter 5 • Deploying the SAS High-Performance Analytics Environment</b> . . . . .	<b>71</b>
Infrastructure Deployment Process Overview . . . . .	71
Overview of Deploying the Analytics Environment . . . . .	72
Encrypting SASHDAT Files . . . . .	75
Install the Analytics Environment . . . . .	76
Configuring the Analytics Environment for SASHDAT Encryption . . . . .	79
Validating the Analytics Environment Deployment . . . . .	80
Resource Management for the Analytics Environment . . . . .	81

<b>Chapter 6 • Deploying SAS Embedded Process for Hadoop</b>	<b>85</b>
Infrastructure Deployment Process Overview	85
Important Note	86
In-Database Deployment Package for Hadoop	86
Hadoop Installation and Configuration	89
SASEP-SERVERS.SH Script	98
Hadoop Permissions	104
Documentation for Using In-Database Processing in Hadoop	105
<b>Chapter 7 • Configuring the Analytics Environment for a Remote Parallel Connection</b>	<b>107</b>
Infrastructure Deployment Process Overview	107
Overview of Configuring the Analytics Environment for a Remote Parallel Connection	108
Preparing for a Remote Parallel Connection	109
Configuring for Access to a Data Store with a SAS Embedded Process	113
<b>Appendix 1 • Updating the SAS High-Performance Analytics Infrastructure</b>	<b>117</b>
Overview of Updating the Analytics Infrastructure	117
Updating the SAS High-Performance Computing Management Console	117
Updating SAS High-Performance Deployment of Hadoop	118
Update the Analytics Environment	127
<b>Appendix 2 • SAS High-Performance Analytics Infrastructure Command Reference</b>	<b>129</b>
<b>Appendix 3 • SAS High-Performance Analytics Environment Client-Side Environment Variables</b>	<b>131</b>
<b>Appendix 4 • Deploying on SELinux and IPTables</b>	<b>133</b>
Overview of Deploying on SELinux and IPTables	133
Prepare the Management Console	134
Prepare Hadoop	134
Prepare the Analytics Environment	135
Analytics Environment Post-Installation Modifications	135
iptables File	136
<b>Recommended Reading</b>	<b>137</b>
<b>Glossary</b>	<b>139</b>
<b>Index</b>	<b>145</b>

## What's New

# What's New in Installation and Configuration for SAS High-Performance Analytics Infrastructure 2.94

---

## Overview

The *SAS High-Performance Analytics Infrastructure: Installation and Configuration Guide* explains how to install and initially configure the SAS High-Performance Analytics infrastructure. This infrastructure consists of the following products:

- SAS High-Performance Computing Management Console 2.6
- SAS High-Performance Deployment of Hadoop 2.8
- SAS High-Performance Analytics environment 2.94

(also referred to as the SAS High-Performance Node Installation)

SAS High-Performance Analytics Infrastructure 2.94 includes the following changes and enhancements:

- Encryption for SASHDAT

---

## Encryption for SASHDAT

In release 2.94, the SAS High-Performance Analytics environment supports reading and writing files using AES encryption with 256-bit keys. For more information, see [“Encrypting SASHDAT Files” on page 75](#).



# Accessibility

For information about the accessibility of any of the products mentioned in this document, see the usage documentation for that product.





## 1

# Introduction to Deploying the SAS High-Performance Analytics Infrastructure

<i>What Is Covered in This Document?</i> .....	<b>1</b>
<i>Which Version Do I Use?</i> .....	<b>2</b>
<i>Experimental Software</i> .....	<b>2</b>
<i>What is the Infrastructure?</i> .....	<b>2</b>
<i>Where Do I Locate My Analytics Cluster?</i> .....	<b>4</b>
Overview of Locating Your Analytics Cluster .....	4
Analytics Cluster Co-Located with Your Data Store .....	6
Analytics Cluster Remote from Your Data Store (Serial Connection) .....	7
Analytics Cluster Remote from Your Data Store (Parallel Connection) .....	8
<i>Deploying the Infrastructure</i> .....	<b>9</b>
Overview of Deploying the Infrastructure .....	9
Step 1: Create a SAS Software Depot .....	9
Step 2: Check for Documentation Updates .....	10
Step 3: Prepare Your Analytics Cluster .....	10
Step 4: (Optional) Deploy SAS High-Performance Computing Management Console .....	10
Step 5: (Optional) Deploy Hadoop .....	10
Step 6: Deploy the SAS High-Performance Analytics Environment .....	11
Step 7: (Optional) Deploy the SAS Embedded Process for Hadoop .....	11
Step 8: (Optional) Configure the Analytics Environment for a Remote Parallel Connection .....	11

## What Is Covered in This Document?

This document covers tasks that are required after you and your SAS representative have decided what software you need and on what machines you will install the software. At this point, you can begin performing some pre-installation tasks, such as creating a SAS Software Depot if your site already does not have one and setting up the operating system user accounts that you will need.

By the end of this document, you will have deployed the SAS High-Performance Analytics environment, and optionally, SAS High-Performance Computing Management Console, and SAS High-Performance Deployment of Hadoop.

You will then be ready to deploy your SAS solution (such as SAS Visual Analytics, SAS High-Performance Risk, and SAS High-Performance Analytics Server) on top of the SAS High-Performance Analytics infrastructure. For more information, see the documentation for your respective SAS solution.

---

## Which Version Do I Use?

This document is published for each major release of the SAS High-Performance Analytics infrastructure, which consists of the following products:

- SAS High-Performance Computing Management Console, version 2.6
- SAS High-Performance Deployment for Hadoop, version 2.8
- SAS High-Performance Analytics environment, version 2.94

(also referred to as the SAS High-Performance Node Installation)

Refer to your order summary to determine the specific version of the infrastructure that is included in your SAS order. Your order summary resides in your SAS Software Depot for your respective order under the `install_doc` directory (for example, `C:\SAS Software Depot\install_doc\my-order\ordersummary.html`).

---

## Experimental Software

Experimental software is sometimes included as part of a production-release product. It is provided to (sometimes targeted) customers in order to obtain feedback. All experimental uses are marked Experimental in this document.

The design and implementation of experimental software might change before any production release. Experimental software has been tested prior to release, but it has not necessarily been tested to production-quality standards, and so should be used with care.

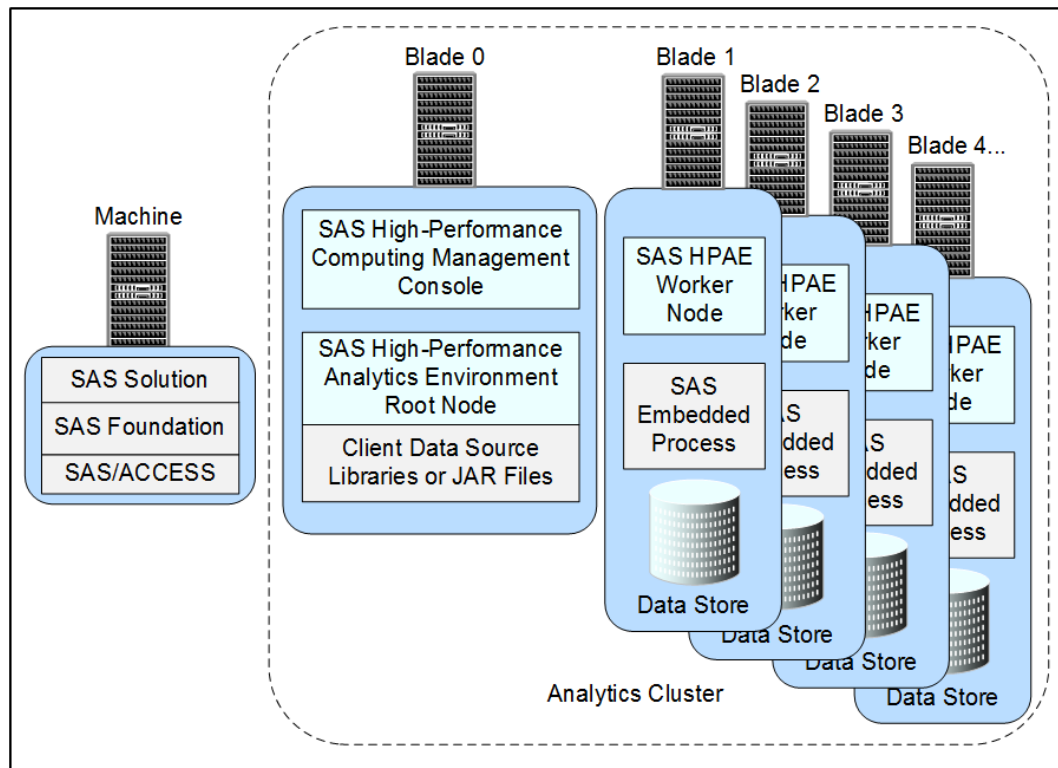
---

## What is the Infrastructure?

The SAS High-Performance Analytics infrastructure consists of software that performs analytic tasks in a high-performance environment, which is characterized by massively parallel processing (MPP). The infrastructure is used by SAS products and solutions that typically analyze big data that resides in a distributed data storage appliance or Hadoop cluster.

The following figure depicts the SAS High-Performance Analytics infrastructure in its most basic topology:

**Figure 1.1** SAS High-Performance Analytics Infrastructure Topology (Simplified)



The SAS High-Performance Analytics infrastructure consists of the following components:

- SAS High-Performance Analytics environment

The SAS High-Performance Analytics environment is the core of the infrastructure. The environment performs analytic computations on an analytics cluster. The analytics cluster is a Hadoop cluster or a data appliance.

- (Optional) SAS High-Performance Deployment of Hadoop

Some solutions, such as SAS Visual Analytics, rely on a SAS data store that is co-located with the SAS High-Performance Analytics environment on the analytics cluster. One option for this co-located data store is the SAS High-Performance Deployment for Hadoop. This is an Apache Hadoop distribution that is easily configured for use with the SAS High-Performance Analytics environment. It adds services to Apache Hadoop to write SASHDAT file blocks evenly across the HDFS filesystem. This even distribution provides a balanced workload across the machines in the cluster and enables SAS analytic processes to read SASHDAT tables at very impressive rates.

Alternatively, these SAS high-performance analytic solutions can use a pre-existing, supported Hadoop deployment.

- (Optional) SAS High-Performance Computing Management Console

The SAS High-Performance Computing Management Console is used to ease the administration of distributed, high-performance computing (HPC)

environments. Tasks such as configuring passwordless SSH, propagating user accounts and public keys, and managing CPU and memory resources on the analytics cluster are all made easier by the management console.

Other software on the analytics cluster include the following:

- SAS/ACCESS Interface and SAS Embedded Process

Together the SAS/ACCESS Interface and SAS Embedded Process provide a high-speed parallel connection that delivers data from the co-located SAS data source to the SAS High-Performance Analytics environment on the analytics cluster. These components are contained in a deployment package that is specific for your data source.

For more information, refer to the *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf> and the *SAS/ACCESS for Relational Databases: Reference*, available at <http://support.sas.com/documentation/cdl/en/acreldb/67473/PDF/default/acreldb.pdf>.

**Note:** For deployments that use Hadoop for the co-located data provider and access SASHDAT tables exclusively, SAS/ACCESS and SAS Embedded Process is *not* needed.

- Database client libraries or JAR files

Data vendor-supplied client libraries—or in the case of Hadoop, JAR files—are required for the SAS Embedded Process to transfer data to and from the data store and the SAS High-Performance Analytics environment.

- SAS solutions

The SAS High-Performance Analytics infrastructure is used by various SAS High-Performance solutions such as the following:

- SAS High-Performance Analytics Server

For more information, refer to <http://support.sas.com/documentation/onlinedoc/hpa>.

- SAS High-Performance Marketing Optimization

For more information, refer to <http://support.sas.com/documentation/onlinedoc/mktopt/index.html>.

- SAS High-Performance Risk

For more information, refer to <http://support.sas.com/documentation/onlinedoc/hprisk/index.html>.

- SAS Visual Analytics

For more information, refer to <http://support.sas.com/documentation/onlinedoc/va/index.html>.

---

## Where Do I Locate My Analytics Cluster?

### Overview of Locating Your Analytics Cluster

You have two options for where to locate your SAS analytics cluster:

- Co-locate SAS with your data store.
- Separate SAS from your data store.

When your SAS analytics cluster is separated (remote) from your data store, you have two basic options for transferring data:

- ☐ Serial data transfer using SAS/ACCESS.
- ☐ Parallel data transfer using SAS/ACCESS in conjunction with the SAS Embedded Process.

The topics in this section contain simple diagrams that describe each option for analytics cluster placement:

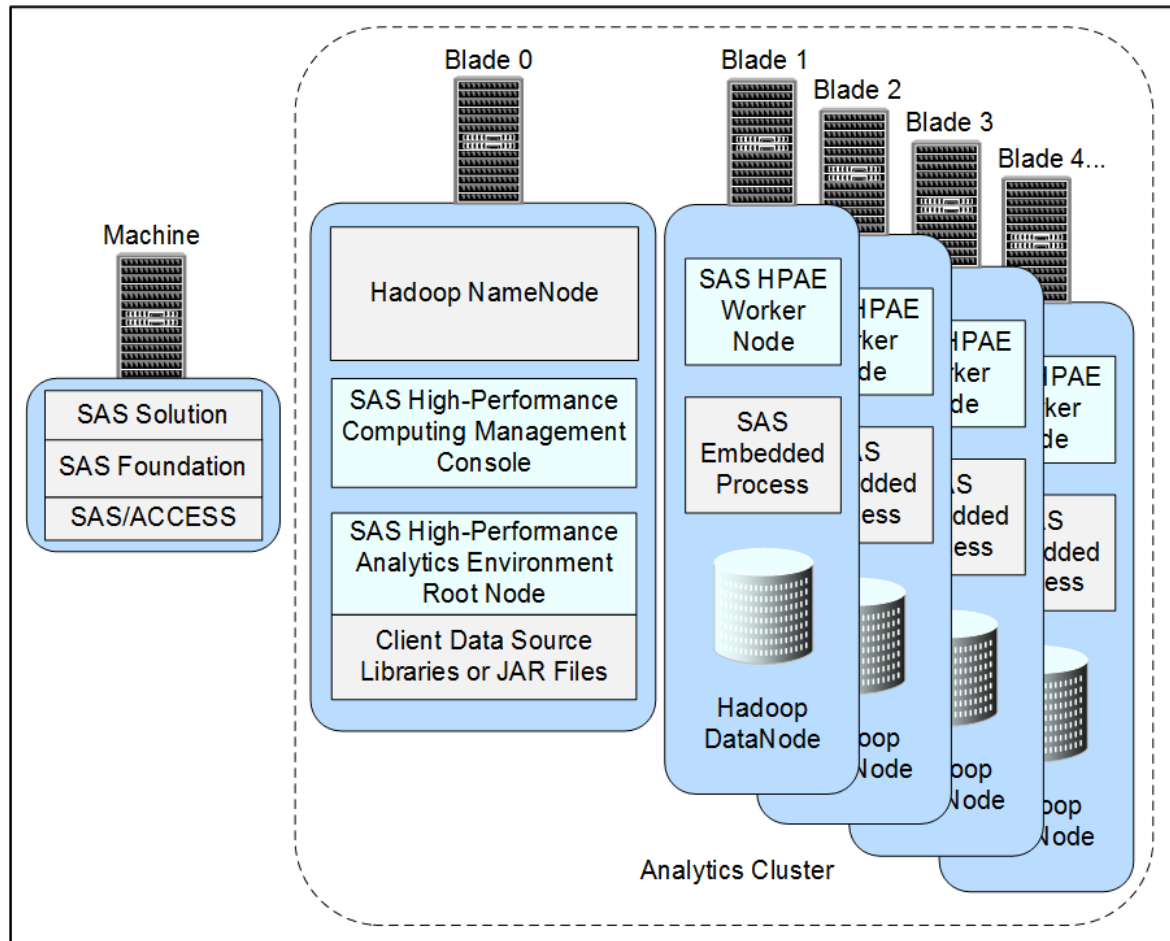
- [Co-Located with the data store](#)
- [Remote from the data store \(serial connection\)](#)
- [Remote from the data store \(parallel connection\)](#)

**TIP** Where you locate your cluster depends on a number of criteria. Your SAS representative will know the latest supported configurations, and can work with you to help you determine which cluster placement option works best for your site. Also, there might be solution-specific criteria that you should consider when determining your analytics cluster location. For more information, see the installation or administration guide for your specific SAS solution.

## Analytics Cluster Co-Located with Your Data Store

The following figure shows the analytics cluster co-located on your Hadoop cluster:

**Figure 1.2** Analytics Cluster Co-Located on the Hadoop Cluster

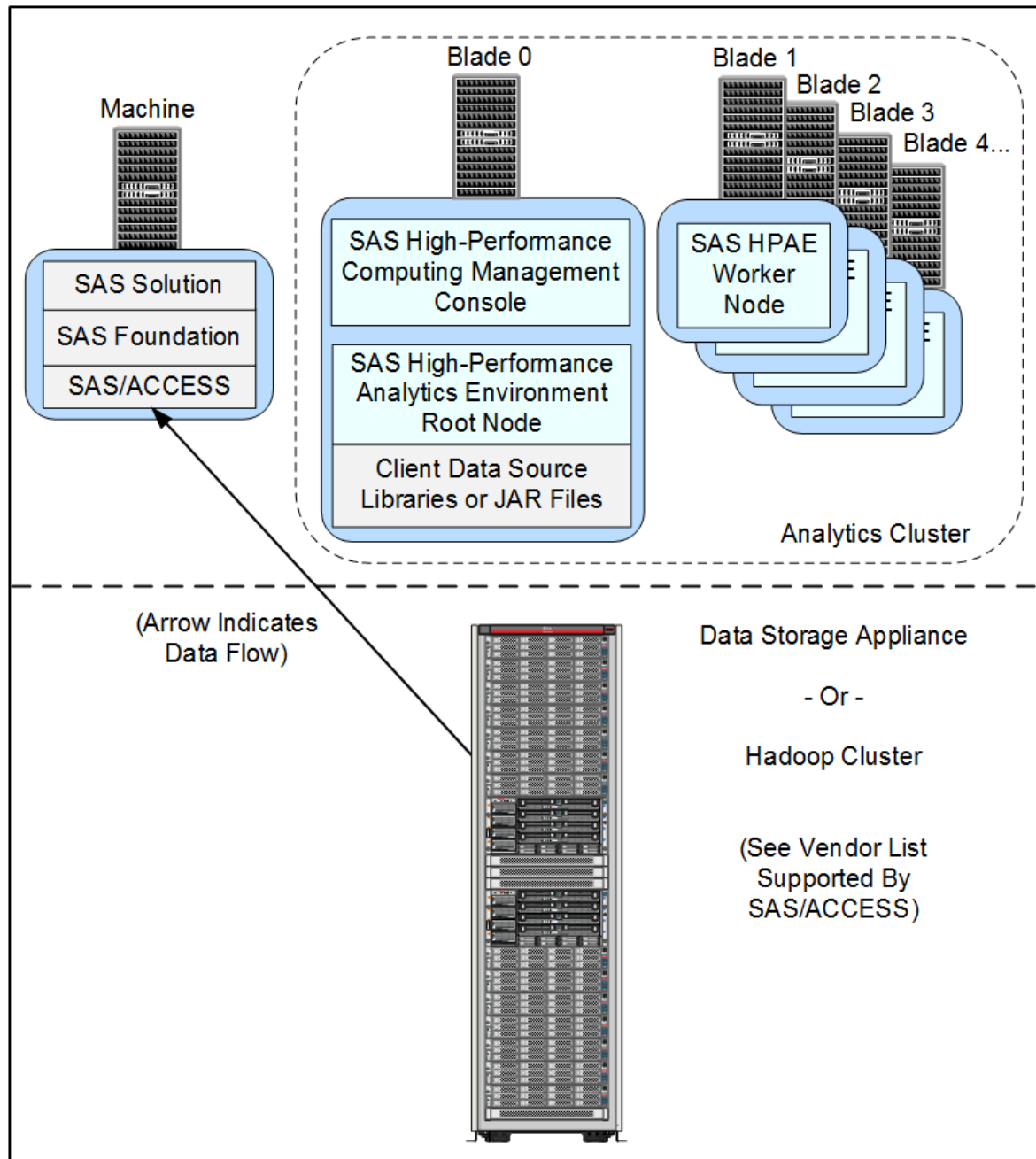


**Note:** For deployments that use Hadoop for the co-located data provider and access SASHDAT tables exclusively, SAS/ACCESS and the SAS Embedded Process are not needed.

## Analytics Cluster Remote from Your Data Store (Serial Connection)

The following figure shows the analytics cluster using a serial connection to your remote data store:

**Figure 1.3** Analytics Cluster Remote from Your Data Store (Serial Connection)



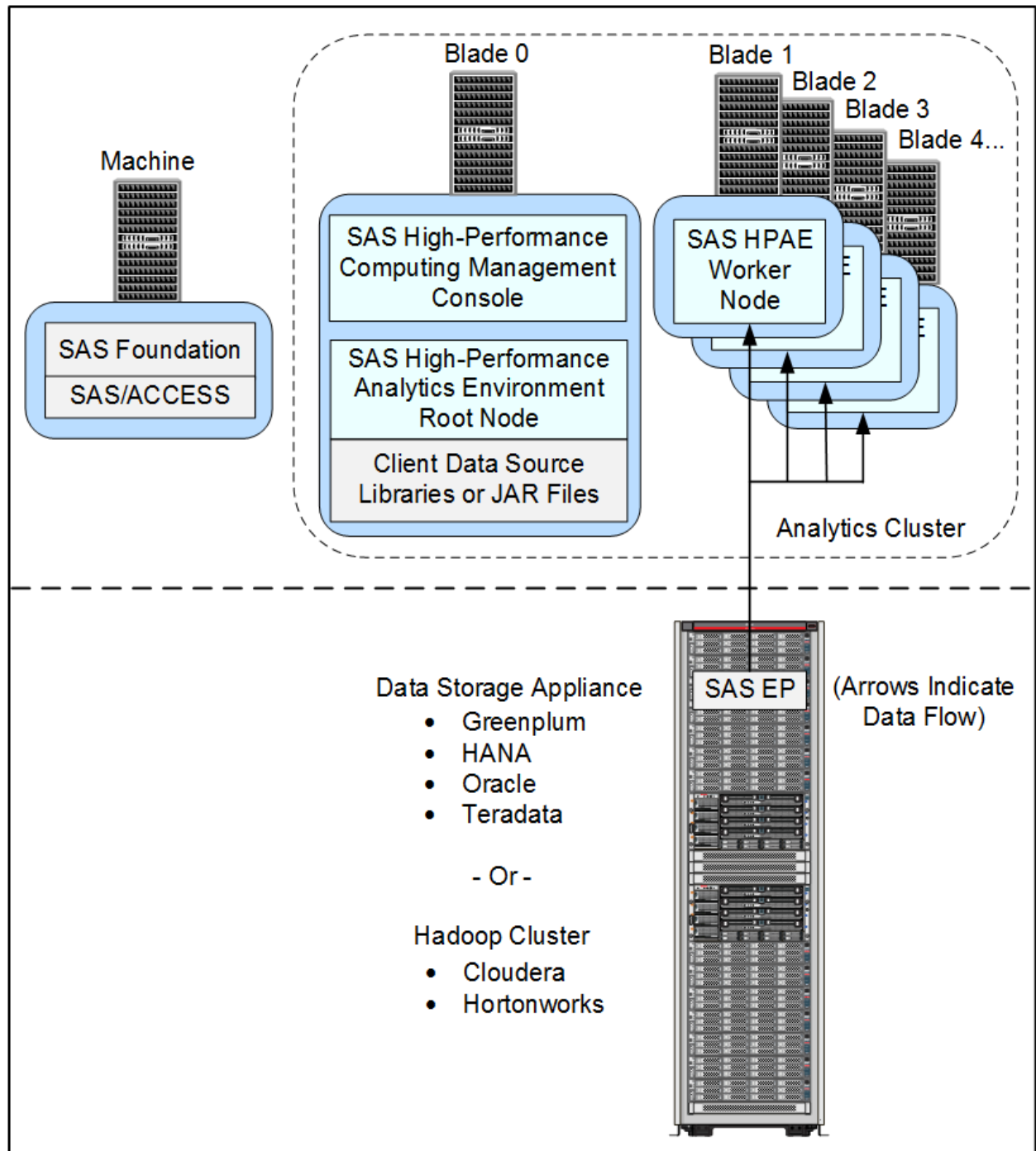
The serial connection between the analytics cluster and your data store is achieved by using the SAS/ACCESS Interface. SAS/ACCESS is orderable in a deployment package that is specific for your data source. For more information,

refer to the *SAS/ACCESS for Relational Databases: Reference*, available at <http://support.sas.com/documentation/onlinedoc/access/index.html>.

## Analytics Cluster Remote from Your Data Store (Parallel Connection)

The following figure shows the analytics cluster using a parallel connection to your remote data store:

**Figure 1.4** Analytics Cluster Remote from Your Data Store (Parallel Connection)





**Note:** In the second maintenance release of SAS 9.4, SAS Embedded Process supports the Cloudera and Hortonworks distributions of Hadoop. For more detailed information, see the SAS Foundation system requirements documentation for your operating environment, available at <http://support.sas.com/resources/sysreq/index.html>.

Together the SAS/ACCESS Interface and SAS Embedded Process provide a high-speed parallel connection that delivers data from your data source to the SAS High-Performance Analytics environment on the analytics cluster. These components are contained in a deployment package that is specific for your data source. For more information, refer to the *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.

---

## Deploying the Infrastructure

### Overview of Deploying the Infrastructure

The following list summarizes the steps required to install and configure the SAS High-Performance Analytics infrastructure:

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.
4. (Optional) Deploy SAS High-Performance Computing Management Console.
5. (Optional) Deploy co-located Hadoop.
6. Deploy the SAS High-Performance Analytics environment.
7. (Optional) Deploy the SAS Embedded Process for Hadoop.
8. (Optional) Configure the analytics environment for a remote parallel connection

The following sections provide a brief description of each of these tasks. Subsequent chapters in the guide provide the step-by-step instructions.

### Step 1: Create a SAS Software Depot

Create a SAS Software Depot, which is a special file system used to deploy your SAS software. The depot contains the SAS Deployment Wizard—the program used to install and initially configure most SAS software—one or more deployment plans, a SAS installation data file, order data, and product data.

**Note:** If you have chosen to receive SAS through Electronic Software Delivery, a SAS Software Depot is automatically created for you.

For more information, see “Creating a SAS Software Depot” in the *SAS Intelligence Platform: Installation and Configuration Guide*, available at <http://support.sas.com/documentation/cdl/en/biig/63852/HTML/default/p03intellplatform00installgd.htm>.

## Step 2: Check for Documentation Updates

It is very important to check for late-breaking installation information in SAS Notes and also to review the system requirements for your SAS software.

- SAS Notes

Go to this web page and click **Outstanding Alert Status Installation Problems**:

<http://support.sas.com/notes/index.html>.

- system requirements

Refer to the system requirements for your SAS solution, available at <http://support.sas.com/resources/sysreq/index.html>.

## Step 3: Prepare Your Analytics Cluster

Preparing your analytics cluster includes tasks such as creating a list of machine names in your grid hosts file. Setting up passwordless SSH is required, as well as considering system umask settings. You must determine which operating system is required to install, configure, and run the SAS High-Performance Analytics infrastructure. Also, you will need to designate ports for the various SAS components that you are deploying.

For more information, see [Chapter 2, “Preparing Your System to Deploy the SAS High-Performance Analytics Infrastructure,”](#) on page 13.

## Step 4: (Optional) Deploy SAS High-Performance Computing Management Console

SAS High-Performance Computing Management Console is an optional web application tool that eases the administrative burden on multiple machines in a distributed computing environment.

For example, when you are creating operating system accounts and passwordless SSH on all machines in the cluster or on blades across the appliance, the management console enables you to perform these tasks from one location.

For more information, see [Chapter 3, “Deploying SAS High-Performance Computing Management Console,”](#) on page 33.

## Step 5: (Optional) Deploy Hadoop

If your site wants to use Hadoop as the co-located data store, then you can install and configure SAS High-Performance Deployment of Hadoop or use one of the supported Hadoop distributions.

For more information, see [Chapter 4, “Deploying Co-Located Hadoop,”](#) on page 45.

## **Step 6: Deploy the SAS High-Performance Analytics Environment**

The SAS High-Performance Analytics environment consists of a root node and worker nodes. The product is installed by a self-extracting shell script.

Software for the root node is deployed on the first host. Software for a worker node is installed on each remaining machine in the cluster or database appliance.

For more information, see [Chapter 5, “Deploying the SAS High-Performance Analytics Environment,”](#) on page 71.

## **Step 7: (Optional) Deploy the SAS Embedded Process for Hadoop**

Together the SAS/ACCESS Interface and SAS Embedded Process provide a high-speed parallel connection that delivers data from the co-located SAS data source to the SAS High-Performance Analytics environment on the analytics cluster. These components are contained in a deployment package that is specific for your data source.

For more information, see [Chapter 6, “Deploying SAS Embedded Process for Hadoop,”](#) on page 85.

## **Step 8: (Optional) Configure the Analytics Environment for a Remote Parallel Connection**

You can optionally configure the SAS High-Performance Analytics Environment for a remote parallel connection.

For more information, see [Chapter 7, “Configuring the Analytics Environment for a Remote Parallel Connection,”](#) on page 107.



## 2

# Preparing Your System to Deploy the SAS High-Performance Analytics Infrastructure

<b>Infrastructure Deployment Process Overview</b>	<b>14</b>
<b>System Settings for the Infrastructure</b>	<b>14</b>
<b>List the Machines in the Cluster or Appliance</b>	<b>15</b>
<b>Review Passwordless Secure Shell Requirements</b>	<b>16</b>
<b>Preparing for Kerberos</b>	<b>16</b>
Kerberos Prerequisites	16
Generate and Test Host Principals	17
Configure Passwordless SSH to use Kerberos	18
Preparing the Analytics Environment for Kerberos	18
Preparing Hadoop for Kerberos	19
<b>Preparing to Install SAS High-Performance Computing Management Console</b>	<b>24</b>
User Account Considerations for the Management Console	24
Management Console Requirements	25
<b>Preparing to Deploy Hadoop</b>	<b>25</b>
Install Hadoop Using root	25
User Accounts for Hadoop	26
Preparing for YARN (Experimental)	27
Install a Java Runtime Environment	27
Plan for Hadoop Directories	28
<b>Preparing to Deploy the SAS High-Performance Analytics Environment</b>	<b>28</b>
User Accounts for the SAS High-Performance Analytics Environment	28
Consider Umask Settings	29
Additional Prerequisite for Greenplum Deployments	30
<b>Recommended Database Names</b>	<b>30</b>
<b>Pre-installation Ports Checklist for SAS</b>	<b>30</b>

---

## Infrastructure Deployment Process Overview

Preparing your analytics cluster is the third of eight steps required to install and configure the SAS High-Performance Analytics infrastructure.

1. Create a SAS Software Depot.
2. Check for documentation updates.
- ▶ **3. Prepare your analytics cluster.**
4. (Optional) Deploy SAS High-Performance Computing Management Console.
5. (Optional) Deploy co-located Hadoop.
6. Deploy the SAS High-Performance Analytics environment.
7. (Optional) Deploy the SAS Embedded Process for Hadoop.
8. (Optional) Configure the analytics environment for a remote parallel connection.

---

## System Settings for the Infrastructure

Understand the system requirements for a successful SAS High-Performance Analytics infrastructure deployment before you begin. The lists that follow offer recommended settings for the analytics infrastructure on every machine in the cluster or blade in the data appliance:

- Modify `/etc/ssh/sshd_config` with the following setting:  

```
MaxStartups 1000
```
- Modify `/etc/security/limits.conf` with the following settings:  

```
* soft nproc 65536
* hard nproc 65536
* soft nofile 350000
* hard nofile 350000
```
- Modify `/etc/security/limits.d/90-nproc.conf` with the following setting:  

```
* soft nproc 65536
```
- Modify `/etc/sysconfig/cpuspeed` with the following setting:  

```
GOVERNOR=performance
```
- The SAS High-Performance Analytics components require approximately 1.4 GB of disk space. SAS High-Performance Deployment of Hadoop requires approximately 300 MB of disk space for the software. This estimate does not include the disk space that is needed for storing data that is added to Hadoop Distributed File System (HDFS) for use by the SAS High-Performance Analytics environment.

For more information, refer to the system requirements for your SAS solution, available at <http://support.sas.com/resources/sysreq/index.html>.

## List the Machines in the Cluster or Appliance

Before the SAS High-Performance Analytics infrastructure can be installed on the machines in the cluster, you must create a file that lists all of the host names of the machines in the cluster.

On blade 0 (known as the Master Server on Greenplum), create an `/etc/gridhosts` file for use by SAS High-Performance Computing Management Console, SAS High-Performance Deployment of Hadoop, and the SAS High-Performance Analytic environment. (The grid hosts file is copied to the other machines in the cluster during the installation process.) If the management console is located on a machine that is not a member of the analytics cluster, then this machine must also contain a copy of `/etc/gridhosts` with its host name added to the list of machines. For more information, see [Chapter 3, “Deploying SAS High-Performance Computing Management Console,”](#) on page 33 before you start the installation.

You can use short names or fully qualified domain names so long as the host names in the file resolve to IP addresses. The long and short host names for each node must be resolvable from each node in the environment. The host names listed in the file must be in the same DNS domain and sub-domain. These host names are used for Message Passing Interface (MPI) communication and SAS High-Performance Deployment of Hadoop network communication.

The *root node* is listed first. This is also the machine that is configured as the following, depending on your data provider:

- SAS High-Performance Deployment of Hadoop or a supported Hadoop distribution: NameNode (blade 0)
- Greenplum Data Computing Appliance: Master Server

The following lines are an example of the file contents:

```
grid001
grid002
grid003
grid004
...
```

**TIP** You can use SAS High-Performance Computing Management Console to create and manage your grid hosts file. For more information, see *SAS High-Performance Computing Management Console: User's Guide*, available at <http://support.sas.com/documentation/solutions/hpainfrastructure/>.

---

## Review Passwordless Secure Shell Requirements

Secure Shell (SSH) has the following requirements:

- To support Kerberos, enable GSSAPI authentication methods in your implementation of Secure Shell (SSH).

**Note:** If you are using Kerberos, see [“Configure Passwordless SSH to use Kerberos” on page 18](#).

- Passwordless Secure Shell (SSH) is required on all machines in the cluster or on the data appliance for the following user accounts:

- ☐ root user account

The root account must run SAS High-Performance Computing Management Console and the simultaneous commands (for example, `simsh`, and `simcp`). For more information about management console user accounts, see [“Preparing to Install SAS High-Performance Computing Management Console” on page 24](#).

- ☐ Hadoop user account

For more information about Hadoop user accounts, see [“Preparing to Deploy Hadoop” on page 25](#).

- ☐ SAS High-Performance Analytics environment user account

For more information about the environment’s user accounts, see [“Preparing to Deploy the SAS High-Performance Analytics Environment” on page 28](#).

**TIP** Users’ home directories must be located in the same directory on each machine in the analytics cluster. For example, you will experience problems if user foo has a home directory at `/home/foo` on blade one and blade two, and a home directory at `/mnt/user/foo` on blade three.

---

## Preparing for Kerberos

### Kerberos Prerequisites

The SAS High-Performance Analytics infrastructure supports the Kerberos computer network authentication protocol. Throughout this document, we indicate the particular settings you need to perform in order to make parts of the infrastructure configurable for Kerberos. However, you must understand and be able to verify your security setup. If you are using Kerberos, you need the ability to get a Kerberos ticket.



**Note:** The SAS High-Performance Analytics environment using YARN is *not* supported with SAS High-Performance Deployment of Hadoop running in Secure Mode Hadoop (that is, configured to use Kerberos).

The list of Kerberos prerequisites are as follows:

- A Kerberos key distribution center (KDC)
- All machines configured as Kerberos clients
- Permissions to copy and secure Kerberos keytab files on all machines
- A user principal for the Hadoop user  
(This is used for setting up the cluster and performing administrative functions.)
- Encryption types supported on the Kerberos domain controller should be aes256-cts:normal and aes128-cts:normal

## Generate and Test Host Principals

Every machine in the analytics cluster must have a host principal and a Kerberos keytab in order to operate as Kerberos clients.

To generate and test host principals, follow these steps:

- 1 Execute `kadmin.local` on the KDC.
- 2 Run the following command for each machine in the cluster:  
`addprinc -randkey host/$machine-name`  
where *machine-name* is the host name of the particular machine.
- 3 Generate host keytab files in `kadmin.local` for each machine, by running the following command:  
`ktadd -norandkey -k $machine-name.keytab host/$machine-name`  
where *machine-name* is the name of the particular machine.

**TIP** When generating keytab files, it is a best practice to create files by machine. In the event a keytab file is compromised, the keytab will only contain the host principal associated with machine it resides on, instead of a single file that contains every machine in the environment.

- 4 Copy each generated keytab file to its respective machine under `/etc`, rename the file to `krb5.keytab`, and secure it with mode 600 and owned by root.

For example:

```
cp keytab /etc/krb5.keytab
chown root:root /etc/krb5.keytab
chmod 600 /etc/krb5.keytab
```

- 5 At this point, any user with a principal in Kerberos should be able to use `kinit` successfully to get a ticket granting ticket.

For example:

```
kinit
Password for hdfs@HOST.DOMAIN.NET:
```

As the Hadoop user, you can run the `klist` command to check the status of your Kerberos ticket. For example:

```
klist
Ticket cache: FILE:/tmp/krb5cc_493
Default principal: hdfs@HOST.DOMAIN.NET

Valid starting    Expires          Service principal
06/20/14 09:51:26 06/27/14 09:51:26 krbtgt/HOST.DOMAIN.NET@HOST.DOMAIN.NET
        renew until 06/22/14 09:51:26
```

**Note:** If you intend to deploy the SAS Embedded Process on the cluster for use with SAS/ACCESS Interface to Hadoop, then a user keytab file for the user ID that runs HDFS is required.

## Configure Passwordless SSH to use Kerberos

Passwordless access of some form is a requirement of the SAS High-Performance Analytics environment through its use of the Message Passing Interface (MPI). Traditionally, public key authentication in Secure Shell (SSH) is used to meet the passwordless access requirement. For Secure Mode Hadoop, GSSAPI with Kerberos is used as the passwordless SSH mechanism. GSSAPI with Kerberos not only meets the passwordless SSH requirements, but also supplies Hadoop with the credentials required for users to perform operations in HDFS with SAS LASR Analytic Server and SASHDAT files. Certain options must be set in the SSH daemon and SSH client configuration files. Those options are as follows and assume a default configuration of `sshd`.

To configure passwordless SSH to use Kerberos, follow these steps:

- 1 In the `sshd_config` file, set:

```
GSSAPIAuthentication yes
```

- 2 In the `ssh_config` file, set:

```
Host *.domain.net
```

```
GSSAPIAuthentication yes
```

```
GSSAPIDelegateCredentials yes
```

where *domain.net* is the domain name used by the machine in the cluster.

**TIP** Although you can specify `host *`, this is not recommended because it would allow GSSAPI Authentication with any host name.

## Preparing the Analytics Environment for Kerberos

During startup, the Message Passing Interface (MPI) sends a user's Kerberos credentials cache (KRB5CCNAME) which can cause an issue when Hadoop attempts to use Kerberos credentials to perform operations in HDFS.

Under Secure Shell (SSH), a random set of characters are appended to the credentials cache file, so the value of the KRB5CCNAME environment variable

is different for each machine. To set the correct value for KRB5CCNAME on each machine, you must use the option below when asked for additional options to MPIRUN during the analytics environment installation:

```
-genvlist `env | sed -e s/=.*// | sed /KRB5CCNAME/d | tr -d
'\n' `TKPATH,LD_LIBRARY_PATH
```

For more information, see [Table 5.2 on page 77](#).

You must use a launcher that supports GSSAPI authentication because the implementation of SSH that is included with SAS does not support it. Add the following to your SAS programs on the client:

```
option set=GRIDRSHCOMMAND="/path-to-file/ssh";
```

**TIP** Adding GRIDRSHCOMMAND to your sasv9\_usermods.cfg preserves the setting during SAS upgrades and avoids having to manually set that environment variable on the client before starting SAS.

## Preparing Hadoop for Kerberos

### Overview of Preparing Hadoop for Kerberos

Preparing SAS High-Performance Deployment of Hadoop for Kerberos, consists of the following steps:

- 1 [“Adding the Principals Required by Hadoop” on page 19](#)
- 2 [“Creating the Necessary Keytab Files” on page 20](#)
- 3 [“Download and Compile JSVC” on page 21](#)
- 4 [“Download and Use Unlimited Strength JCE Policy Files” on page 21](#)
- 5 [“Configure Self-Signed Certificates for Hadoop” on page 22](#)

### Adding the Principals Required by Hadoop

Secure Mode Hadoop requires a number of principals to work properly with Kerberos. Principals can be created using `addprinc` within `kadmin.local` on the KDC.

Your add principal command should resemble:

```
addprinc -randkey nn/$FQDN@$REALM
```

where `$FQDN` is a fully qualified domain name and `$REALM` is the name of the Kerberos Realm.

If you are using HDFS only, then only the HDFS-specific principals and keytab files are required. For HDFS, you need the following principals:

```
nn/$FQDN@$REALM
```

NameNode principal. Create this for the NameNode machine only.

```
sn/$FQDN@$REALM
```

Secondary NameNode principal. Create this for the NameNode machine only.

dn/\$FQDN@\$REALM

DataNode principal. Create this for every machine in the cluster except for the NameNode machine.

HTTP/\$FQDN@\$REALM

HTTP server principal, used by WebDFS. Create this for every machine in the cluster.

## Creating the Necessary Keytab Files

After creating the principals, keytab files must be created for each service and machine. Keytab files are created using `ktadd` within `kadmin.local` on the KDC. Your `ktadd` command for the NameNode service should resemble the following:

```
ktadd -k /path-to-file/service_name.keytab nn/$FQDN@$REALM
```

where *service\_name* is a value like `hdfs_nnservice`, `hdfs_dnservice`, or `http` as shown in the following examples.

For example, your keytab files should be similar to the following:

- an `hdfs_nnservice.keytab` file containing three principals, with two encryption types per principal.

The NameNode principal starts with `nn`, the Secondary NameNode principal starts with `sn`, and the host principal is included in the keytab file as described in the Apache Hadoop documentation. Your KVNO value might differ. `hdfs_nnservice.keytab` is copied to the NameNode only and owned by the Hadoop user with a mode of 600.

```
Keytab name: FILE:hdfs_nnservice.keytab
KVNO Principal
-----
2 nn/node1.domain.net@DOMAIN.NET
2 nn/node1.domain.net@DOMAIN.NET
2 sn/node1.domain.net@DOMAIN.NET
2 sn/node1.domain.net@DOMAIN.NET
3 host/node1.domain.net@DOMAIN.NET
3 host/node1.domain.net@DOMAIN.NET
```

- an `hdfs_dnservice.keytab` file that contains a principal for each DataNode host, with two encryption types per principal.

The DataNode principle starts with `dn` and the host principals are included in the keytab file as described in the Apache Hadoop documentation. Your KVNO value might differ. `hdfs_dnservice.keytab` is copied only to the DataNodes and owned by the Hadoop user with a mode of 600.

```
Keytab name: FILE:hdfs_dnservice.keytab
KVNO Principal
-----
2 dn/node2.domain.net@DOMAIN.NET
2 dn/node2.domain.net@DOMAIN.NET
2 dn/node3.domain.net@DOMAIN.NET
2 dn/node3.domain.net@DOMAIN.NET
...
2 dn/noden.domain.net@DOMAIN.NET
2 dn/noden.domain.net@DOMAIN.NET
3 host/node2.domain.net@DOMAIN.NET
3 host/node2.domain.net@DOMAIN.NET
3 host/node3.domain.net@DOMAIN.NET
```

```

3 host/node3.domain.net@DOMAIN.NET
...
3 host/noden.domain.net@DOMAIN.NET
3 host/noden.domain.net@DOMAIN.NET

```

- an `http.keytab` file that contains a principal for each machine in the cluster, with two encryption types per principal.

Your KVNO value might differ. `http.keytab` is copied to all machines and is owned by the Hadoop user with a mode of 600.

```

Keytab name: FILE:http.keytab
KVNO Principal
-----
2 HTTP/node1.domain.net@DOMAIN.NET
2 HTTP/node1.domain.net@DOMAIN.NET
2 HTTP/node2.domain.net@DOMAIN.NET
2 HTTP/node2.domain.net@DOMAIN.NET
...
2 HTTP/noden.domain.net@DOMAIN.NET
2 HTTP/noden.domain.net@DOMAIN.NET

```

## Download and Compile JSVC

The JSVC binary is required to start secure DataNodes on a privileged port. A server is started on the privileged port by root and the process is then switched to the secure DataNode user. The JSVC binary is not currently included with Apache Hadoop. This section details where to get the source, how to compile it on a machine, and where to copy it.

To download and compile JSVC, follow these steps:

- 1 Download the JSVC source from apache at <http://archive.apache.org/dist/commons/daemon/source/commons-daemon-1.0.15-src.tar.gz>.
- 2 Extract the file into a directory you have write access to.
- 3 Change directory to `commons-daemon-1.0.15-src/src/native/unix`.  
**Note:** This directory contains the `INSTALL.txt` file that describes the installation process.
- 4 Execute `./configure` and correct any issues found during the pre-make.
- 5 After a successful configure, compile the binary by running `make`. This generates a file called `jsvc` in the directory.
- 6 Copy the `jsvc` file to `$HADOOP_HOME/sbin` on every DataNode in the cluster.

Note the path to `jsvc` because this path is used later in `hadoop-env.sh`.

## Download and Use Unlimited Strength JCE Policy Files

For encryption strengths above 128 bit, you must download the latest JCE (Java Cryptography Extension) for the JRE you are using. For this document, the JCE was used to provide 256-bit, AES encryption. Keep export and import laws in mind when dealing with encryption. Check with your site's legal department if you have any questions.

## Configure Self-Signed Certificates for Hadoop

In order to secure Hadoop communications between cluster machines, you must setup the cluster to use HTTPS. This section goes through the process of generating the necessary files and the configuration options required to enable HTTPS using self-signed certificates.

To configure self-signed certificates for Hadoop, follow these steps:

- 1 Create the client and server key directories by running the following commands on each machine:

```
mkdir -p /etc/security/serverKeys
mkdir -p /etc/security/clientKeys
```

- 2 Create the key store on each machine:

```
cd /etc/security/serverKeys
keytool -genkey -alias $shortname -keyalg RSA -keysize 1024
-dname "CN=
$shortname.domain.net,OU=unit,O=company,L=location,ST=state,
C=country" -keypass $somepass -keystore keystore -storepass
$somepass
```

- 3 Create the certificate on each machine:

```
cd /etc/security/serverKeys
keytool -export -alias $shortname -keystore keystore -rfc -
file $shortname.cert -storepass $somepass
```

- 4 Import the certificate into the key store on each machine:

```
cd /etc/security/serverKeys
keytool -import -noprompt -alias $shortname -file
$shortname.cert -keystore truststore -storepass $somepass
```

- 5 Import the certificate for each machine into the all store file. After you complete this on the first machine, copy the generated allstore file to the next machine until you have copied the file and run the following command on each machine. The end result is an allstore file containing each machine's certificate.

```
cd /etc/security/serverKeys
keytool -import -noprompt -alias $shortname -file
$shortname.cert -keystore allstore -storepass $somepass
```

- 6 Use the command below to verify the allstore file has a certificate from every node.

```
keytool -list -v -keystore allstore -storepass $somepass
```

- 7 Move the allstore file to the /etc/security/clientKeys directory on each machine.

- 8 Refer to [Table 2.1](#) and make sure the generated files on each machine are in their respective locations, with appropriate ownership and mode.

The allstore file to be used is the one containing all the certificates, which was verified in the previous step. The directories /etc/security/

`serverKeys` and `/etc/security/clientKeys` should have a mode of 755 and owned by `hdfs:hadoop`.

**Table 2.1** Summary of Certificates

Filename	Location	Ownership	Mode
keystore	<code>/etc/security/serverKeys</code>	<code>hdfs:hadoop</code>	<code>r--r-----</code>
truststore	<code>/etc/security/serverKeys</code>	<code>hdfs:hadoop</code>	<code>r--r-----</code>
allstore	<code>/etc/security/clientKeys</code>	<code>hdfs:hadoop</code>	<code>r--r--r--</code>
<code>\$shortname.cer</code>	<code>/etc/security/serverKeys</code>	<code>hdfs:hadoop</code>	<code>r--r-----</code>

**9** Make the following SSL related additions or changes to each respective file:

- `core-site.xml`
- `ssl-server.xml`
- `ssl-client.xml`

`core-site.xml`:

```
<property>
  <name>hadoop.ssl.require.client.cert</name>
  <value>>false</value>
</property>
<property>
  <name>hadoop.ssl.hostname.verifier</name>
  <value>DEFAULT</value>
</property>
<property>
  <name>hadoop.ssl.keystores.factory.class</name>
  <value>org.apache.hadoop.security.ssl.FileBasedKeyStoresFactory</value>
</property>
<property>
  <name>hadoop.ssl.server.conf</name>
  <value>ssl-server.xml</value>
</property>
<property>
  <name>hadoop.ssl.client.conf</name>
  <value>ssl-client.xml</value>
</property>
```

`ssl-server.xml`:

**Note:** The `ssl-server.xml` file should be owned by `hdfs:hadoop` with a mode of 440. Replace `$somepass` with the keystore password. Because this file has the keystore password in clear-text, make sure that only Hadoop service accounts are added to the `hadoop` group.

```
<property>
  <name>ssl.server.truststore.location</name>
  <value>/etc/security/serverKeys/truststore</value>
</property>
<property>
```

```

        <name>ssl.server.truststore.password</name>
        <value>$somepass</value>
    </property>
    <property>
        <name>ssl.server.truststore.type</name>
        <value>jks</value>
    </property>
    <property>
        <name>ssl.server.keystore.location</name>
        <value>/etc/security/serverKeys/keystore</value>
    </property>
    <property>
        <name>ssl.server.keystore.password</name>
        <value>$somepass</value>
    </property>
    <property>
        <name>ssl.server.keystore.type</name>
        <value>jks</value>
    </property>
    <property>
        <name>ssl.server.keystore.keypassword</name>
        <value>$somepass</value>
    </property>

```

ssl-client.xml:

**Note:** The ssl-client.xml file should be owned by `hdfs:hadoop` with a mode of 440. The same information from the preceding note applies to this file.

```

    <property>
        <name>ssl.client.truststore.location</name>
        <value>/etc/security/clientKeys/allstore</value>
    </property>
    <property>
        <name>ssl.client.truststore.password</name>
        <value>$somepass</value>
    </property>
    <property>
        <name>ssl.client.truststore.type</name>
        <value>jks</value>
    </property>

```

---

## Preparing to Install SAS High-Performance Computing Management Console

### User Account Considerations for the Management Console

SAS High-Performance Computing Management Console is installed from either an RPM or a tarball package and must be installed and configured with the root user ID. The root user account must have passwordless secure shell (SSH) access between all the machines in the cluster. The console includes a web server. The web server is started with the root user ID, and it runs as the root user ID.



The reason that the web server for the console must run as the root user ID is that the console can be used to add, modify, and delete operating system user accounts from the local passwords database (`/etc/passwd` and `/etc/shadow`). Only the root user ID has Read and Write access to these files.

Be aware that you do not need to log on to the console with the root user ID. In fact, the console is typically configured to use console user accounts. Administrators can log on to the console with a console user account that is managed by the console itself and does not have any representation in the local passwords database or whatever security provider the operating system is configured to use.

## Management Console Requirements

Before you install SAS High-Performance Computing Management Console, make sure that you have performed the following tasks:

- Make sure that the Perl extension `perl-Net-SSLeay` is installed.
- For PAM authentication, make sure that the `Authen::PAM PERL` module is installed.

**Note:** The management console can manage operating system user accounts if the machines are configured to use the `/etc/passwd` local database only.

- Create the list of all the cluster machines in the `/etc/gridhosts` file. You can use short names or fully qualified domain names so long as the host names in the file resolve to IP addresses. These host names are used for Message Passing Interface (MPI) communication and Hadoop network communication. For more information, see “List the Machines in the Cluster or Appliance” on page 15.
- Locate the software.

Make sure that your SAS Software Depot has been created. (For more information, see “Creating a SAS Software Depot” in the *SAS Intelligence Platform: Installation and Configuration Guide*, available at <http://support.sas.com/documentation/cdl/en/biig/63852/HTML/default/p03intellplatform00installgd.htm>.)

---

## Preparing to Deploy Hadoop

If you are using Kerberos, see also “Preparing for Kerberos” on page 16.

### Install Hadoop Using root

As is the case with most enterprise Hadoop distributions such as Cloudera or Hortonworks, root privileges are needed when installing SAS High-Performance Deployment of Hadoop.

The installer must be root in order to `chown` and `chmod` files appropriately. Unlike earlier releases, there is a new user (`yarn`), and the Hadoop user (`hdfs`) cannot change file ownership to another user. Also, installing Hadoop using root facilitates implementation of Kerberos and Secure Mode Hadoop. For more

information, refer to the Apache document available at <http://hadoop.apache.org/docs/r2.4.0/hadoop-project-dist/hadoop-common/SecureMode.html>.

## User Accounts for Hadoop

Apache recommends that the HDFS and YARN daemons and the MapReduce JobHistory server run as different Linux users. It is also recommended that these users share the same primary Linux group. The following table summarizes Hadoop user and group information:

**Table 2.2** *Hadoop Users and Their Primary Group*

User:Group	Daemons
hdfs:hadoop	NameNode, Secondary NameNode, JournalNode, DataNode
yarn:hadoop	ResourceManager, NodeManager
mapred:hadoop	MapReduce JobHistory Server

The accounts with which you deploy Hadoop, MapReduce, and YARN must have passwordless secure shell (SSH) access between all the machines in the cluster.

**TIP** Although the Hadoop installation program can run as any user, you might find it easier to run `hadoopInstall` as root so that it can set permissions and ownership of the Hadoop data directories for the user account that runs Hadoop.

As a convention, this document uses an account and group named `hdfs` when describing how to deploy and run SAS High-Performance Deployment of Hadoop. `mapred` and `yarn` are used for the MapReduce JobHistory Server user and the YARN user, respectively. If you do not already have an account that meets the requirements, you can use SAS High-Performance Computing Management Console to add the appropriate user ID.

If your site has a requirement for a reserved UID and GID for the `hdfs` user account, then create the user and group on each machine before continuing with the installation.

**Note:** We recommend that you install SAS High-Performance Computing Management Console before setting up the user accounts that you will need for the rest of the SAS High-Performance Analytics infrastructure. The console enables you to easily manage user accounts across the machines of a cluster. For more information, see “[Create the First User Account and Propagate the SSH Key](#)” on page 41.

SAS High-Performance Deployment of Hadoop is installed from a TAR.GZ file. An installation and configuration program, `hadoopInstall`, is available after the archive is extracted.

## Preparing for YARN (Experimental)

When deploying the SAS High-Performance Deployment for Hadoop, you must decide whether to use YARN. YARN stands for “Yet Another Resource Negotiator.” It consists of a framework that manages execution and schedules resource requests for distributed applications. For information about how to configure the analytics environment with YARN, see [“Resource Management for the Analytics Environment” on page 81](#).

**Note:** The SAS High-Performance Analytics environment using YARN is *not* supported with SAS High-Performance Deployment of Hadoop running in Secure Mode Hadoop (that is, configured to use Kerberos).

If you decide to use YARN with the SAS High-Performance Deployment for Hadoop, you must do the following:

- Create Linux user accounts for YARN and MapReduce to run YARN and MapReduce jobs on the machines in the cluster.

These user accounts must exist on all the machines in the cluster and must be configured for passwordless SSH. For more information, see [“User Accounts for Hadoop” on page 26](#).

- Create a Linux group and make it the primary group for the Hadoop, YARN, and MapReduce users.
- Provide YARN-related input when prompted during the SAS High-Performance Deployment for Hadoop installation.

For more information, see [“Install SAS High-Performance Deployment of Hadoop” on page 47](#).

## Install a Java Runtime Environment

Hadoop requires a Java Runtime Environment (JRE) or Java Development Kit (JDK) on every machine in the cluster. The path to the Java executable must be the same on all of the machines in the cluster. If this requirement is already met, make a note of the path and proceed to installing SAS High-Performance Deployment of Hadoop.

If the requirement is not met, then install a JRE or JDK on the machine that is used as the grid host. You can use the `simsh` and `simcp` commands to copy the files to the other machines in the cluster.

**Note:** For information about the simultaneous commands, see [“Simultaneous Utilities Commands” on page 129](#).

**Example Code 2.1** *Sample simsh and simcp Commands*

```
/opt/TKGrid/bin/simsh mkdir /opt/java
/opt/TKGrid/bin/simcp /opt/java/jdk1.6.0_31 /opt/java
```

For information about the supported Java version, see <http://wiki.apache.org/hadoop/HadoopJavaVersions>. SAS High-Performance Deployment of Hadoop uses the Apache Hadoop 2.4 version.

## Plan for Hadoop Directories

The following table lists the default directories where the SAS High-Performance Deployment of Hadoop stores content:

**Table 2.3** Default SAS High-Performance Deployment of Hadoop Directory Locations

Default Directory Location	Description
<code>hadoop-name</code>	The <code>hadoop-name</code> directory is the location on the file system where the NameNode stores the namespace and transactions logs persistently. This location is formatted by Hadoop during the configuration stage.
<code>hadoop-data</code>	The <code>hadoop-data</code> directory is the location on the file system where the DataNodes store data in blocks.
<code>hadoop-local</code>	The <code>hadoop-local</code> directory is the location on the file system where temporary MapReduce data is written.
<code>hadoop-system</code>	The <code>hadoop-system</code> directory is the location on the file system where the MapReduce framework writes system files.

**Note:** These Hadoop directories must reside on local storage. The exception is the `hadoop-data` directory, which can be on a storage area network (SAN). Network attached storage (NAS) devices are not supported.

You create the Hadoop installation directory on the NameNode machine. The installation script prompts you for this Hadoop installation directory and the names for each of the subdirectories (listed in [Table 2.3](#)) which it creates for you on every machine in the cluster.

Especially in the case of the data directory, it is important to designate a location that is large enough to contain all of your data. If you want to use more than one data device, see [“\(Optional\) Deploy with Multiple Data Devices”](#) on page 53.

---

## Preparing to Deploy the SAS High-Performance Analytics Environment

If you are using Kerberos, see also [“Preparing for Kerberos”](#) on page 16.

### User Accounts for the SAS High-Performance Analytics Environment

This topic describes the user account requirements for deploying and running the SAS High-Performance Analytics environment:

- Installation and configuration must be run with the same user account.

- The installer account must have passwordless secure shell (SSH) access between all the machines in the cluster.

**TIP** We recommend that you install SAS High-Performance Computing Management Console before setting up the user accounts that you will need for the rest of the SAS High-Performance Analytics infrastructure. The console enables you to easily manage user accounts across the machines of a cluster. For more information, see “[User Account Considerations for the Management Console](#)” on page 24.

The SAS High-Performance Analytics environment uses a shell script installer. You can use a SAS installer account to install this software if the user account meets the following requirements:

- The SAS installer account has Write access to the directory that you want to use and Write permission to the same directory path on every machine in the cluster.
- The SAS installer account is configured for passwordless SSH on all the machines in the cluster.

The root user ID can be used to install the SAS High-Performance Analytics environment, but it is not a requirement. When users start a process on the machines in the cluster with SAS software, the process runs under the user ID that starts the process. Any user accounts running analytic environment processes must also be configured with passwordless SSH.

## Consider Umask Settings

The SAS High-Performance Analytics environment installation script (described in a later section) prompts you for a umask setting. Its default is no setting.

If you do not enter any umask setting, then jobs, servers, and so on, that use the analytics environment create files with the user’s pre-existing umask set on the operating system. If you set a value for umask, then that umask is used and overrides each user’s system umask setting.

Entering a value of 027 ensures that only users in the same operating system group can read these files.

**Note:** Remember that the account used to run the LASRMonitor process (by default, sas) must be able to read the table and server files in `/opt/VADP/var` and any other related subdirectories.

**Note:** Remember that the LASRMonitor process that is part of SAS Visual Analytics must be run with an account (by default, sas) that can read the server signature file. (This signature file is created when you start a SAS LASR Analytic Server and the file is specified in SAS metadata. For more information, see “Establishing Connectivity to a SAS LASR Analytic Server” in Chapter 4 of *SAS Intelligence Platform: Data Administration Guide*, available at <http://support.sas.com/documentation/cdl/en/bidsag/67493/HTML/default/viewer.htm#n1y0g0l4bgiduzn1o6jdy4l8c61d.htm>.

You can also add umask settings to the resource settings file for the SAS Analytics environment. For more information, see “[Resource Management for the Analytics Environment](#)” on page 81.

For more information about using umask, refer to your Linux documentation.

## Additional Prerequisite for Greenplum Deployments

For deployments that rely on Greenplum data appliances, the SAS High-Performance Analytics environment requires that you also deploy the appropriate SAS/ACCESS interface and SAS Embedded Process that SAS supplies with SAS In-Database products. For more information, see *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.

---

## Recommended Database Names

SAS solutions, such as SAS Visual Analytics, that rely on a co-located data provider can make use of two database instances.

The first instance often already exists and is expected to have your operational or transactional data that you want to explore and analyze.

A second database instance is used to support the self-service data access features of SAS Visual Analytics. This database is commonly named “vapublic,” but you can specify a different name if you prefer. Keep these names handy, as the SAS Deployment Wizard prompts you for them when deploying your SAS solution.

---

## Pre-installation Ports Checklist for SAS

While you are creating operating system user accounts and groups, you need to review the set of ports that SAS will use by default. If any of these ports is unavailable, select an alternate port, and record the new port on the ports pre-installation checklist that follows.

The following checklist indicates what ports are used for SAS by default and gives you a place to enter the port numbers that you will actually use.

We recommend that you document each SAS port that you reserve in the following standard location on each machine: `/etc/services`. This practice will help avoid port conflicts on the affected machines.

**Note:** These checklists are superseded by more complete and up-to-date checklists that can be found at <http://support.sas.com/installcenter/plans>. This website also contains a corresponding deployment plan and an architectural diagram. If you are a SAS solutions customer, consult the pre-installation checklist provided by your SAS representative for a complete list of ports that you must designate.

**Table 2.4** Pre-installation Checklist for SAS Ports

SAS Component	Default Port	Data Direction	Actual Port
YARN ResourceManager Scheduler	8030	Inbound	

SAS Component	Default Port	Data Direction	Actual Port
YARN ResourceManager Resource Tracker	8031	Inbound	
YARN ResourceManager	8032	Inbound	
YARN ResourceManager Admin	8033	Inbound	
YARN Node Manager Localizer	8040	Inbound	
YARN Node Manager Web Application	8042	Inbound	
YARN ResourceManager Web Application	8088	Inbound	
SAS High-Performance Computing Management Console server	10020	Inbound	
MapReduce Job History	10021	Inbound	
YARN Web Proxy	10022	Inbound	
MapReduce Job History Admin	10033	Inbound	
MapReduce Job History Web Application	19888	Inbound	
Hadoop Service on the NameNode	15452	Inbound	
Hadoop Service on the DataNode	15453	Inbound	
Hadoop DataNode Address	50010	Inbound	
Hadoop DataNode IPC Address	50020	Inbound	
Hadoop JobTracker	50030	Inbound	
Hadoop TaskTracker	50060	Inbound	
Hadoop Name Node web interface	50070	Inbound	
Hadoop DataNode HTTP Address	50075	Inbound	
Hadoop Secondary NameNode	50090	Inbound	
Hadoop Name Node Backup Address	50100	Inbound	
Hadoop Name Node Backup HTTP Address	50105	Inbound	
Hadoop Name Node HTTPS Address	50470	Inbound	

SAS Component	Default Port	Data Direction	Actual Port
Hadoop DataNode HTTPS Address	50475	Inbound	
SAS High-Performance Deployment of Hadoop	54310	Inbound	
SAS High-Performance Deployment of Hadoop	54311	Inbound	



## 3

## Deploying SAS High-Performance Computing Management Console

<i>Infrastructure Deployment Process Overview</i> .....	33
<i>Benefits of the Management Console</i> .....	34
<i>Overview of Deploying the Management Console</i> .....	34
<i>Installing the Management Console</i> .....	35
Install SAS High-Performance Computing Management Console Using RPM .....	35
Install the Management Console Using tar .....	36
<i>Configure the Management Console</i> .....	36
<i>Create the Installer Account and Propagate the SSH Key</i> .....	38
<i>Create the First User Account and Propagate the SSH Key</i> .....	41

### Infrastructure Deployment Process Overview

Installing and configuring SAS High-Performance Computing Management Console is an optional fourth of eight steps required to install and configure the SAS High-Performance Analytics infrastructure.

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.
- **4. (Optional) Deploy SAS High-Performance Computing Management Console.**
5. (Optional) Deploy co-located Hadoop.
6. Deploy the SAS High-Performance Analytics environment.
7. (Optional) Deploy the SAS Embedded Process for Hadoop.
8. (Optional) Configure the analytics environment for a remote parallel connection.

---

## Benefits of the Management Console

Passwordless SSH is required to start and stop SAS LASR Analytic Servers and to load tables. For some SAS solutions, such as SAS High-Performance Risk and SAS High-Performance Analytic Server, passwordless SSH is required to run jobs on the machines in the cluster.

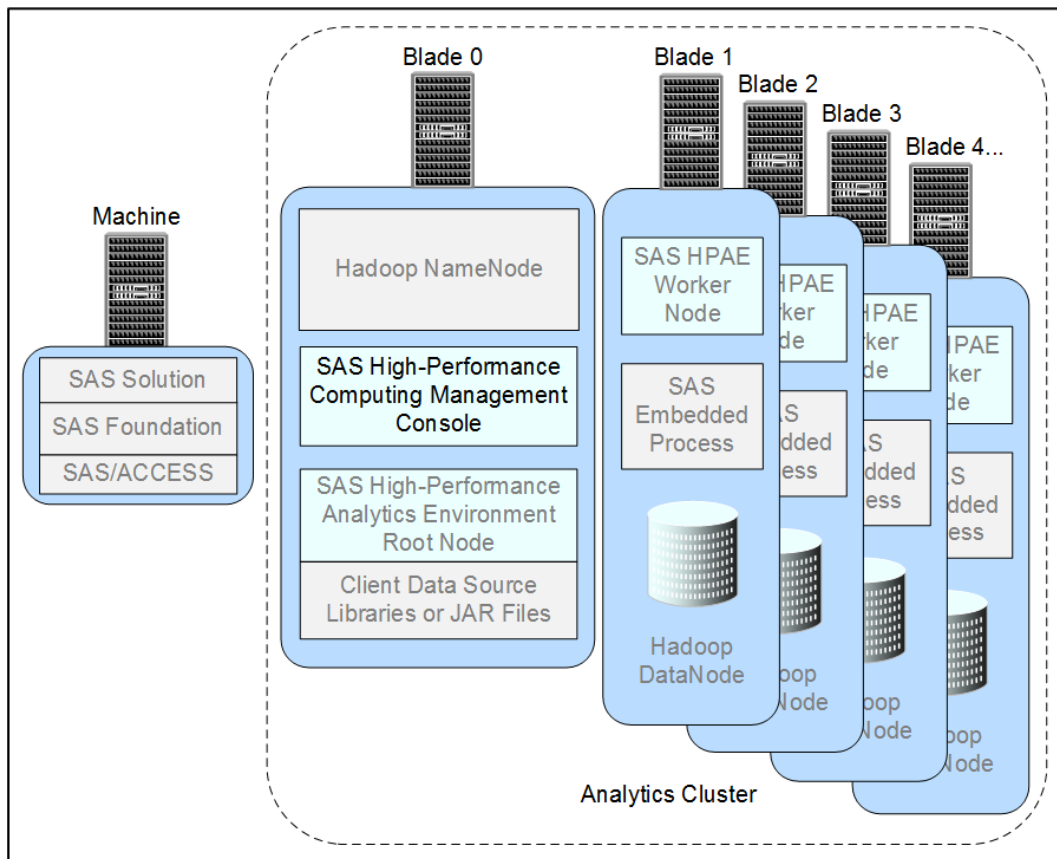
Also, users of some SAS solutions must have an operating system (external) account on all the machines in the cluster and must have the key distributed across the cluster. For more information, see [“Create the First User Account and Propagate the SSH Key” on page 41](#).

SAS High-Performance Computing Management Console enables you to perform these tasks from one location. When you create *new* user accounts using SAS High-Performance Computing Management Console, the console propagates the public key across all the machines in the cluster in a single operation. For more information, see *SAS High-Performance Computing Management Console: User's Guide*, available at <http://support.sas.com/documentation/solutions/hpainfrastructure/>.

---

## Overview of Deploying the Management Console

The SAS High-Performance Computing Management Console is deployed on the machine where the SAS High-Performance Analytics environment is deployed. In this document, that machine is blade 0.

**Figure 3.1** Management Console Deployed with a Data Appliance

## Installing the Management Console

There are two ways to install SAS High-Performance Computing Management Console.

### Install SAS High-Performance Computing Management Console Using RPM

To install SAS High-Performance Computing Management Console using RPM, follow these steps:

**Note:** For information about updating the console, see [Appendix 1, “Updating the SAS High-Performance Analytics Infrastructure,”](#) on page 117.

- 1 Make sure that you have reviewed all of the information contained in the section [“Preparing to Install SAS High-Performance Computing Management Console”](#) on page 24.
- 2 Log on to the target machine as root.
- 3 In your SAS Software Depot, locate the `standalone_installs/SAS_High-Performance_Computing_Management_Console/2_6/Linux_for_x64` directory.

- 4 Enter one of the following commands:
  - To install in the default location of `/opt`:
 

```
rpm -ivh sashpcmc*.rpm
```
  - To install in a location of your choice:
 

```
rpm -ivh --prefix=directory sashpcmc*.rpm
```

where *directory* is an absolute path where you want to install the console.
- 5 Proceed to the topic [“Configure the Management Console” on page 36.](#)

## Install the Management Console Using tar

Some versions of Linux use different RPM libraries and require an alternative means to install SAS High-Performance Computing Management Console. Follow these steps to install the management console using tar:

- 1 Make sure that you have reviewed all of the information contained in the section [“Preparing to Install SAS High-Performance Computing Management Console” on page 24.](#)
- 2 Log on to the target machine as root.
- 3 In your SAS Software Depot, locate the `standalone_installs/SAS_High-Performance_Computing_Management_Console/2_6/Linux_for_x64` directory.
- 4 Copy `sashpcmc-2.6.tar.gz` to the location where you want to install the management console.
- 5 Change to the directory where you copied the tar file, and run the following command:
 

```
tar -xzf sashpcmc-2.6.tar.gz
```

tar extracts the contents into a directory called `sashpcmc`.
- 6 Proceed to the topic [“Configure the Management Console” on page 36.](#)

---

## Configure the Management Console

After installing SAS High-Performance Computing Management Console, you must configure it. This is done with the setup script.

- 1 Log on to the SAS Visual Analytics server and middle tier machine (blade 0) as root.
- 2 Run the setup script by entering the following command:

```
management-console-installation-directory/opt/webmin/utilbin/setup
```

Answer the prompts that follow.

Enter the username for initial login to SAS HPC MC below.

This user will have rights to everything in the SAS HPC MC and

can either be an OS account or new console user. If an OS account exists for the user, then system authentication will be used. If an OS account does not exist, you will be prompted for a password.

**3** Enter the user name for the initial login.

```
Creating using system authentication
Use SSL\HTTPS (yes|no)
```

**4** If you want to use Secure Sockets Layer (SSL) when running the console, enter **yes**. Otherwise, enter **no**.

**5** If you chose not to use SSL, then skip to [Step 7 on page 37](#). Otherwise, the script prompts you to use a pre-existing certificate and key file or to create a new one.

```
Use existing combined certificate and key file or create a new one (file|create)?
```

**6** Make one of two choices:

- Enter **create** for the script to generate the combined private key and SSL certificate file for you.

The script displays output of the `openssl` command that it uses to create the private key pair for you.

- Enter **file** to supply the path to a valid private key pair.

When prompted, enter the absolute path for the combined certificate and key file.

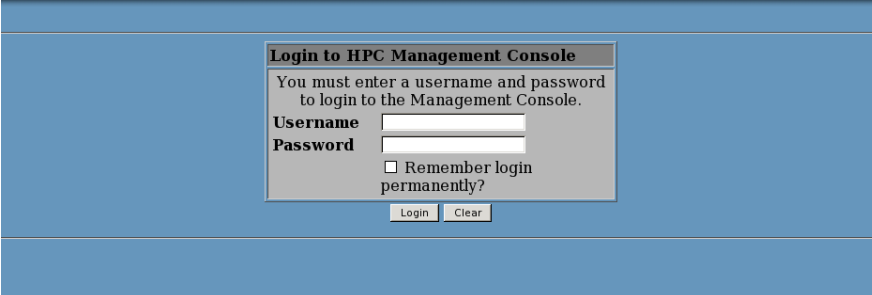
**7** To start the SAS High-Performance Computing Management Console server, enter the following command from any directory:

```
service sashpcmc start
```

**8** Open a web browser and, in the address field, enter the fully qualified domain name for the blade 0 host followed by port 10020.

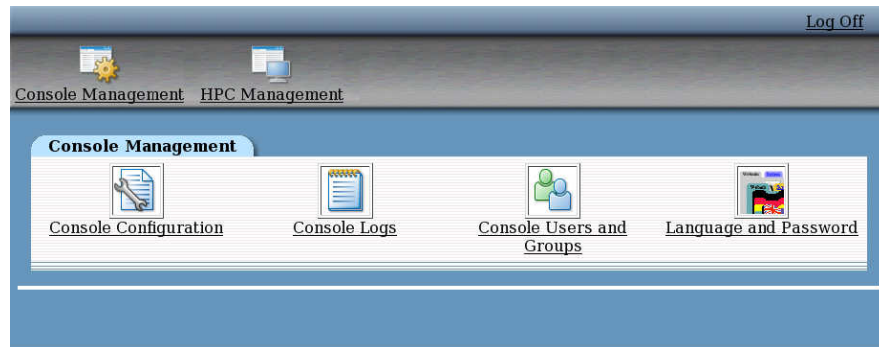
For example: `https://myserver.example.com:10020`

The Login page appears.



**9** Log on to SAS High-Performance Computing Management Console using the credentials that you specified in [Step 2](#).

The Console Management page appears.



## Create the Installer Account and Propagate the SSH Key

The user account needed to start and stop server instances and to load and unload tables to those servers must be configured with passwordless secure shell (SSH).

To reduce the number of operating system (external) accounts, it can be convenient to use the SAS Installer account for both of these purposes.

Implementing passwordless SSH requires that the public key be added to the `authorized_keys` file across all machines in the cluster. When you create user accounts using SAS High-Performance Computing Management Console, the console propagates the public key across all the machines in the cluster in a single operation.

To create an operating system account and propagate the public key, follow these steps:

- 1 Make sure that the SAS High-Performance Computing Management Console server is running. While logged on as the root user, enter the following command from any directory:

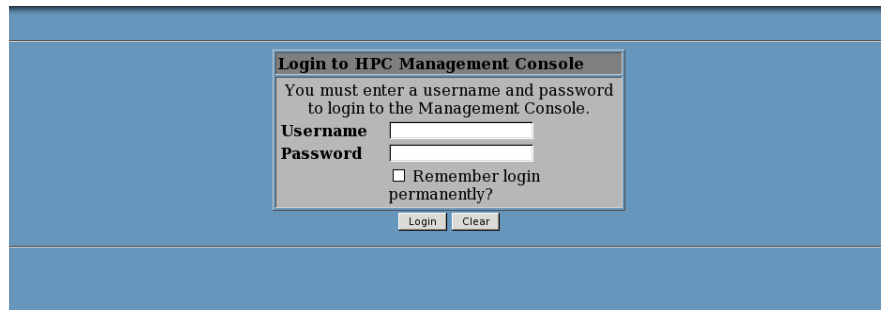
```
service sashpcmc status
```

(If you are logged on as a user other than the root user, the script returns the message `sashpcmc is stopped`.) For more information, see [To start the SAS High-Performance Computing Management Console server on page 37](#).

- 2 Open a web browser and, in the address field, enter the fully qualified domain name for the blade 0 host followed by port 10020.

For example: `http://myserver.example.com:10020`

The Login page appears.



**Login to HPC Management Console**

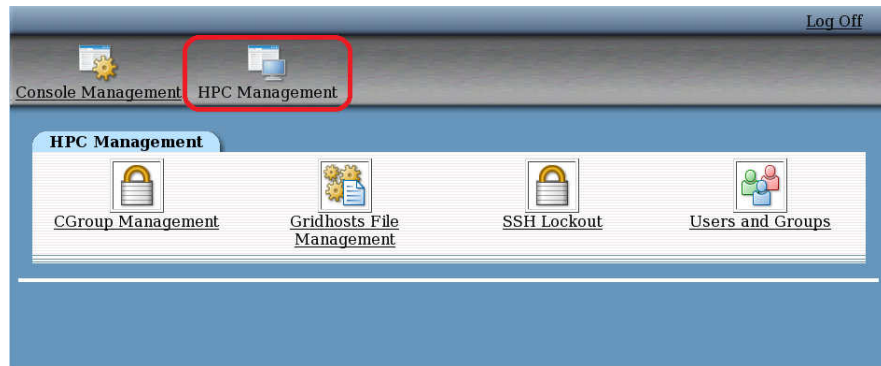
You must enter a username and password to login to the Management Console.

**Username**

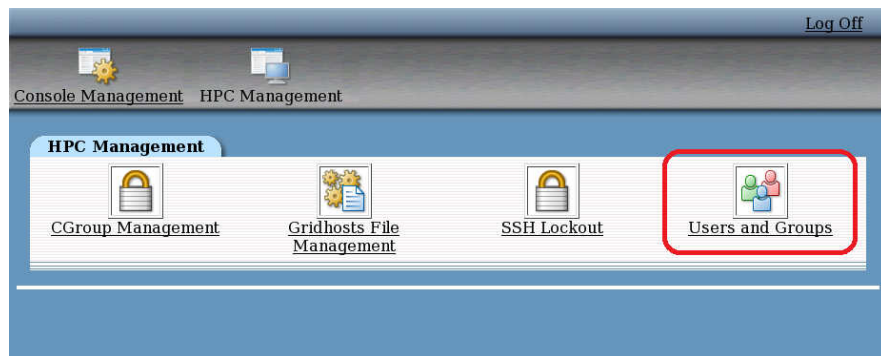
**Password**

☐ Remember login permanently?

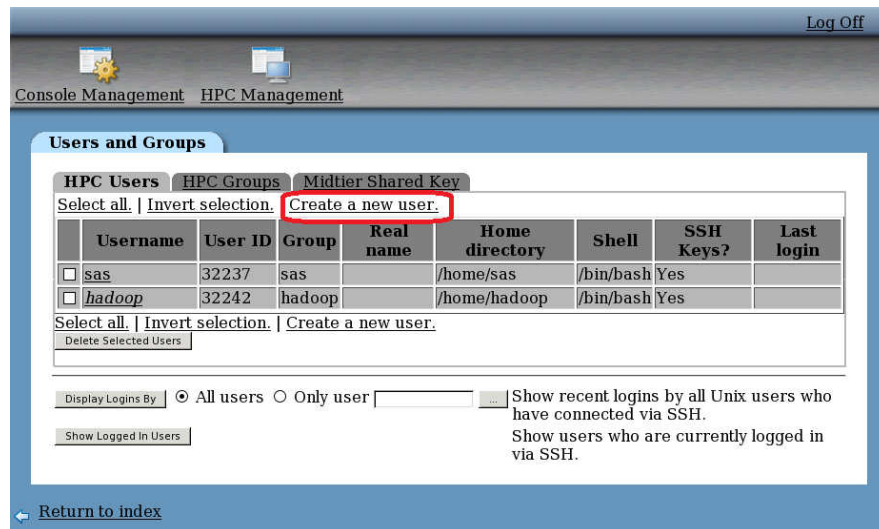
- 3 Log on to SAS High-Performance Computing Management Console.  
The Console Management page appears.



- 4 Click **HPC Management**.  
The HPC Management page appears.



- 5 Click **Users and Groups**.  
The Users and Groups page appears.



6 Click **Create a new user**.

The Create User page appears.

7 Enter information for the new user, using the security policies in place at your site.

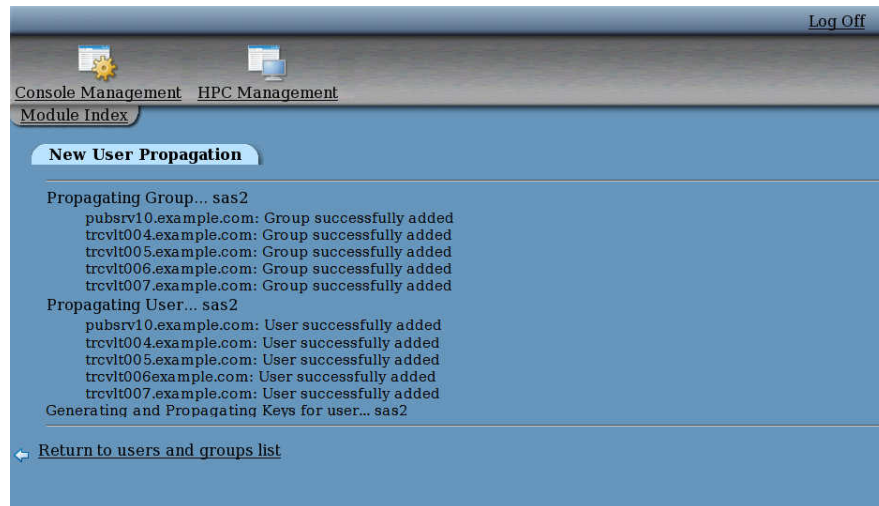
Be sure to choose **Yes** for the following:



- **Propagate User**
- **Generate and Propagate SSH Keys**

When you are finished making your selections, click **Create**.

The New User Propagation page appears and lists the status of the create user command. Your task is successful if you see output similar to the following figure.



## Create the First User Account and Propagate the SSH Key

Depending on their configuration, some SAS solution users must have an operating system (external) account on all the machines in the cluster. Furthermore, the public key might be distributed on each cluster machine in order for their secure shell (SSH) access to operate properly. SAS High-Performance Computing Management Console enables you to perform these two tasks from one location.

To create an operating system account and propagate the public key for SSH, follow these steps:

- 1 Make sure that the SAS High-Performance Computing Management Console server is running. Enter the following command from any directory:

```
service sashpcmc status
```

For more information, see [To start the SAS High-Performance Computing Management Console server on page 37](#).

- 2 Open a web browser and, in the address field, enter the fully qualified domain name for the blade 0 host followed by port 10020.

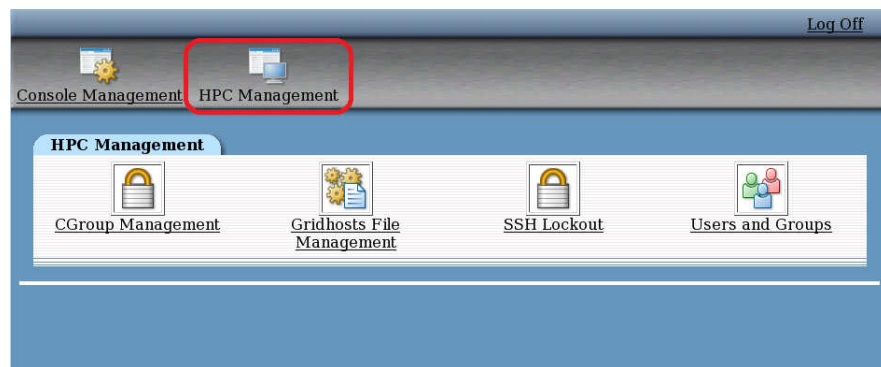
For example: `http://myserver.example.com:10020`

The Login page appears.



A login dialog box titled "Login to HPC Management Console". It contains the text "You must enter a username and password to login to the Management Console." Below this are two input fields: "Username" and "Password". There is a checkbox labeled "Remember login permanently?". At the bottom are two buttons: "Login" and "Clear".

- 3 Log on to SAS High-Performance Computing Management Console.  
The Console Management page appears.



- 4 Click **HPC Management**.  
The Console Management page appears.



- 5 Click **Users and Groups**.  
The Users and Groups page appears.

Log Off

Console Management HPC Management

**Users and Groups**

HPC Users HPC Groups Midtier Shared Key

Select all. | Invert selection. **Create a new user.**

	Username	User ID	Group	Real name	Home directory	Shell	SSH Keys?	Last login
<input type="checkbox"/>	sas	32237	sas		/home/sas	/bin/bash	Yes	
<input type="checkbox"/>	hadoop	32242	hadoop		/home/hadoop	/bin/bash	Yes	

Select all. | Invert selection. | Create a new user.

Delete Selected Users

Display Logins By ☒ All users ☐ Only user  Show recent logins by all Unix users who have connected via SSH.

Show Logged in Users Show users who are currently logged in via SSH.

[Return to index](#)

6 Click **Create a new user**.

The Create User page appears.

**Create User**

**User Details**

Username

User ID ☒ Automatic ☐ Calculated

Real name

Home directory ☒ Automatic ☐ Directory

Shell

Password ☒ No password required ☐ Normal password

☐ Pre-encrypted password

**Password Options**

Password changed  Expiry date

Minimum days  Maximum days

Warning days  Inactive days

**Group Membership**

Primary group ☒ New group with same name as user ☐ New group  ☐ Existing group

Secondary groups

All groups

In groups

**HPC Actions and Settings**

Propagate User ☒ Yes ☐ No

Generate and Propagate SSH Keys ☒ Yes ☐ No

Add Shared Midtier Key ☐ Yes ☒ No

**Upon Creation..**

Create home directory? ☒ Yes ☐ No

Copy template files to home directory? ☒ Yes ☐ No

Create

[Return to users and groups list](#)

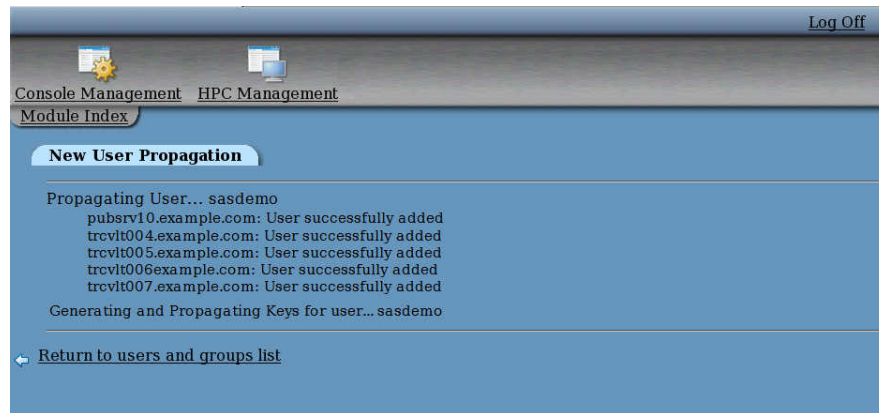
7 Enter information for the new user, using the security policies in place at your site.

Be sure to choose **Yes** for the following:

- **Propagate User**
- **Generate and Propagate SSH Keys**

When you are finished making your selections, click **Create**.

The New User Propagation page appears and lists the status of the create user command. Your task is successful if you see output similar to the following figure.



## 4

## Deploying Co-Located Hadoop

<i>Infrastructure Deployment Process Overview</i> .....	<b>45</b>
<i>Overview of Deploying Hadoop</i> .....	<b>46</b>
<i>Deploying SAS High-Performance Deployment of Hadoop</i> .....	<b>47</b>
What Is SAS High-Performance Deployment of Hadoop? .....	47
Overview of Deploying SAS High-Performance Deployment of Hadoop .....	47
Install SAS High-Performance Deployment of Hadoop .....	47
Post-Installation Steps for Hadoop .....	51
<i>Configuring Existing Hadoop Clusters</i> .....	<b>60</b>
Overview of Configuring Existing Hadoop Clusters .....	60
Prerequisites for Existing Hadoop Clusters .....	60
Configuring the Existing Cloudera Hadoop Cluster .....	61
Configuring the Existing Hortonworks Data Platform Hadoop Cluster .....	65
Configuring the Existing IBM BigInsights Hadoop Cluster .....	67
Configuring the Existing MapR Hadoop Cluster .....	68
Configuring the Existing Pivotal HD Hadoop Cluster .....	68

---

## Infrastructure Deployment Process Overview

Deploying a co-located Hadoop is an optional fifth of eight steps required to install and configure the SAS High-Performance Analytics infrastructure.

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.
4. (Optional) Deploy SAS High-Performance Computing Management Console.
- **5. (Optional) Deploy co-located Hadoop.**
6. Deploy the SAS High-Performance Analytics environment.
7. (Optional) Deploy the SAS Embedded Process for Hadoop.
8. (Optional) Configure the analytics environment for a remote parallel connection.

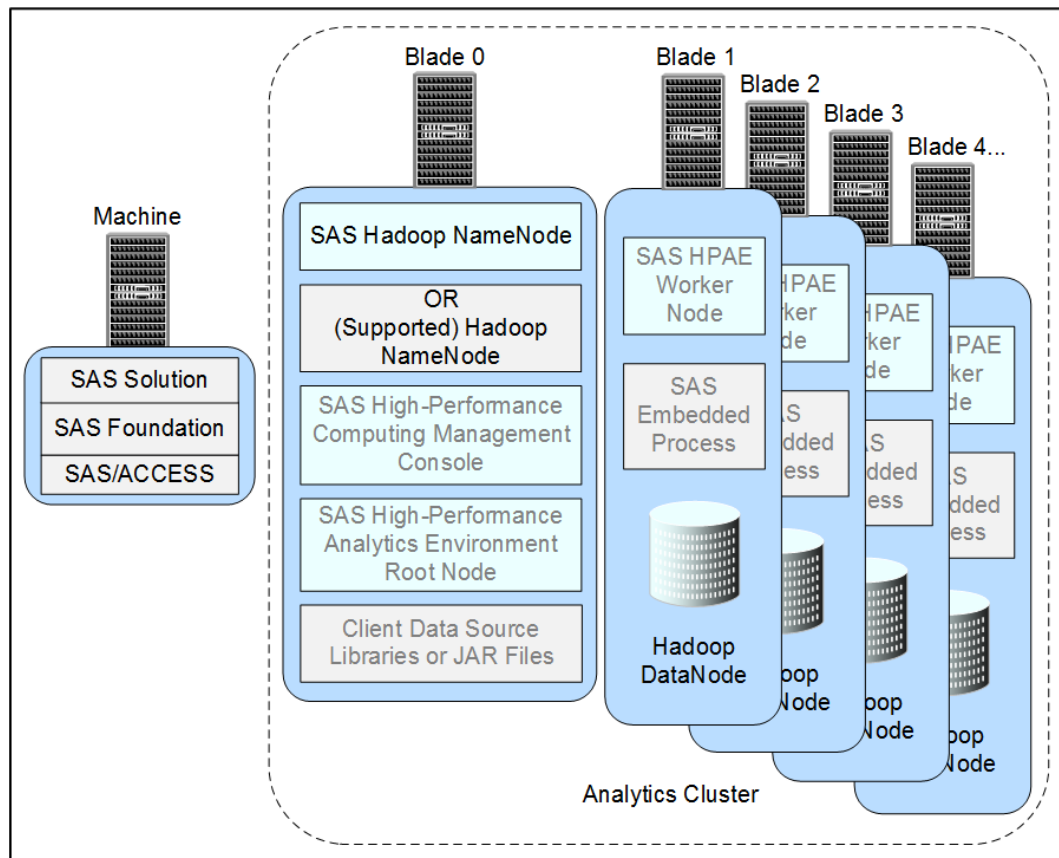
## Overview of Deploying Hadoop

The SAS High-Performance Analytics environment relies on a Hadoop Distributed File System. You have the option of using a Hadoop supplied by SAS, or using another supported Hadoop:

- “Deploying SAS High-Performance Deployment of Hadoop” on page 47.
- “Configuring Existing Hadoop Clusters” on page 60.

Deploying Hadoop requires installing and configuring components on the NameNode machine and DataNodes on the remaining machines in the cluster. In this document, the NameNode is deployed on blade 0.

**Figure 4.1** Analytics Cluster Co-located on the Hadoop Cluster



## Deploying SAS High-Performance Deployment of Hadoop

### What Is SAS High-Performance Deployment of Hadoop?

Some solutions, such as SAS Visual Analytics, rely on a SAS data store that is co-located with the SAS High-Performance Analytic environment on the analytic cluster. One option for this co-located data store is the SAS High-Performance Deployment for Hadoop. This is an Apache Hadoop distribution that is easily configured for use with the SAS High-Performance Analytics environment. It adds services to Apache Hadoop to write SASHDAT file blocks evenly across the HDFS filesystem. This even distribution provides a balanced workload across the machines in the cluster and enables SAS analytic processes to read SASHDAT tables at very impressive rates.

Alternatively, these SAS high-performance analytic solutions can use a pre-existing, supported Hadoop deployment.

### Overview of Deploying SAS High-Performance Deployment of Hadoop

The following steps are required to deploy the SAS High-Performance Deployment of Hadoop:

**Note:** If you want to upgrade a pre-existing SAS High-Performance Deployment of Hadoop system, then see [“Updating SAS High-Performance Deployment of Hadoop” on page 118](#).

- 1 [Prepare for Hadoop on page 25](#)
- 2 [Install Hadoop on page 47](#)
- 3 [Perform post-installation steps on page 51](#)

### Install SAS High-Performance Deployment of Hadoop

The software that is needed for SAS High-Performance Deployment of Hadoop is available from within the SAS Software Depot that was created by the site depot administrator:

```
depot-installation-location/standalone_installs/  
SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_  
x64/sashadoop.tar.gz
```

- 1 Make sure that you have reviewed all of the information contained in the section [“Preparing to Deploy Hadoop” on page 25](#).
- 2 Log on to the Hadoop NameNode machine (blade 0) as root.

For more information, see [“Install Hadoop Using root” on page 25](#).

- 3 Decide where to install Hadoop, and create that directory if it does not exist.

```
mkdir hadoop
```

- 4 Record the name of this directory, as you will need it later in the install process.

- 5 Copy the `sashadoop.tar.gz` file to a temporary location and extract it:

```
cp sashadoop.tar.gz /tmp
cd /tmp
tar xzf sashadoop.tar.gz
```

A directory that is named `sashadoop` is created.

- 6 Change directory to the `sashadoop` directory and run the `hadoopInstall` command:

```
cd sashadoop
./hadoopInstall
```

- 7 Respond to the prompts from the configuration program:

**Table 4.1** SAS High-Performance Deployment of Hadoop Configuration Parameters

Parameter	Description
Choose the type of installation to perform: 1) New installation of SAS Apache Hadoop 2.4.0 with new HDFS. 2) Add the latest LASR support to an existing SAS Apache Hadoop. Leave existing HDFS unmodified. 3) New installation of SAS Apache Hadoop 2.4.0 with upgrade of your existing HDFS directory structure. 4) Quit. [This utility is not used with 3rd-party Hadoop distributions.] Enter choice (1-4). Default is 4: (1/2/3/4)?	Specify 1 and press Enter to perform a new installation.  If you want to upgrade Hadoop (options 2 or 3), see <a href="#">“Overview of Updating SAS High-Performance Deployment of Hadoop” on page 118</a> .
Enter path to install Hadoop. The directory 'hadoop-2.4.0' will be created in the path specified.	Specify the directory that you created in <a href="#">Step 3 on page 48</a> and press Enter.
Do you wish to use Yarn and MR Jobhistory Server? (y/N)	Enter either <code>y</code> or <code>n</code> and press Enter. If you are using YARN, be sure to review <a href="#">“Preparing for YARN (Experimental)” on page 27</a> before proceeding.
Enter replication factor. Default 2	To accept the default, press Enter. Or specify a preferred number of replications for blocks (0 - 10) and press Enter. This prompt corresponds to the <code>dfs.replication</code> property for HDFS.



Parameter	Description
Enter port number for fs.defaultFS. Default 54310	To accept the default port numbers, press Enter for each prompt. Or specify a different port and press Enter. These ports are listed in <a href="#">“Pre-installation Ports Checklist for SAS”</a> on page 30.
Enter port number for dfs.namenode.https-address. Default 50470	
Enter port number for dfs.datanode.https.address. Default 50475	
Enter port number for dfs.datanode.address. Default 50010	
Enter port number for dfs.datanode.ipc.address. Default 50020	
Enter port number for dfs.namenode.http-address. Default 50070	
Enter port number for dfs.datanode.http.address. Default 50075	
Enter port number for dfs.secondary.http.address. Default 50090	
Enter port number for dfs.namenode.backup.address. Default 50100	
Enter port number for dfs.namenode.backup.http-address. Default 50105	
Enter port number for com.sas.lasr.hadoop.service.namenode.port. Default 15452	
Enter port number for com.sas.lasr.hadoop.service.datanode.port. Default 15453	
<hr/>	
[The following port prompts are displayed when you choose to deploy YARN:]	To accept the default port numbers, press Enter for each prompt. Or specify a different port and press Enter. These ports are listed in <a href="#">“Pre-installation Ports Checklist for SAS”</a> on page 30.
Enter port number for mapreduce.jobhistory.admin.address. Default 10033	
Enter port number for mapreduce.jobhistory.webapp.address. Default 19888	
Enter port number for mapreduce.jobhistory.address. Default 10021	
Enter port number for yarn.resourcemanager.scheduler.address. Default 8030	
Enter port number for yarn.resourcemanager.resource-tracker.address. Default 8031	
Enter port number for yarn.resourcemanager.address. Default 8032	
Enter port number for yarn.resourcemanager.admin.address. Default 8033	
Enter port number for yarn.resourcemanager.webapp.address. Default 8088	
Enter port number for yarn.nodemanager.localizer.address. Default 8040	
Enter port number for yarn.nodemanager.webapp.address. Default 8042	
Enter port number for yarn.web-proxy.address. Default 10022	

Parameter	Description
Enter maximum memory allocation per Yarn container. Default 5905	This is the maximum amount of memory (in MB) that YARN can allocate on a particular machine in the cluster. To accept the default, press Enter. Or specify a different value and press Enter.
Enter user that will be running the HDFS server process.	Specify the user name (for example, hdfs) and press Enter.  For more information, see <a href="#">“User Accounts for Hadoop” on page 26</a> .
Enter user that will be running Yarn services	Specify the user name (for example, yarn) and press Enter.  For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter user that will be running the Map Reduce Job History Server.	Specify the user name (for example, mapred) and press Enter.  For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter common primary group for users running Hadoop services.	Apache recommends that the hdfs, mapred, and yarn user accounts share the same primary Linux group (for example, hadoop). Enter a group name and press Enter. For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter path for JAVA_HOME directory. (Default: /usr/lib/jvm/jre)	To accept the default, press Enter. Or specify a different path to the JRE or JDK and press Enter.  <b>Note:</b> The configuration program does not verify that a JRE is installed at <code>/usr/lib/jvm/jre</code> , which is the default path for some Linux vendors.
Enter path for Hadoop data directory. This should be on a large drive. Default is '/hadoop/hadoop-data'.  Enter path for Hadoop name directory. Default is '/hadoop/hadoop-name'.	To accept the default, press Enter. Or specify different paths and press Enter.  <b>Note:</b> The data directory cannot be the root directory of a partition or mount.  <b>Note:</b> If you have more than one data device, enter one of the data directories now, and after the installation, refer to <a href="#">“(Optional) Deploy with Multiple Data Devices” on page 53</a> .
Enter full path to machine list. The NameNode <code>'host'</code> should be listed first.	Enter <code>/etc/gridhosts</code> and press Enter.

- 8** The installation program installs SAS High-Performance Deployment of Hadoop on the local host, configures several files, and then provides a prompt:

The installer can now copy '/hadoop/hadoop-2.4.0' to all the slave machines using scp, skipping the first entry. Perform copy? (YES/no)

Enter **Yes** and press Enter to install SAS High-Performance Deployment of Hadoop on the other machines in the cluster.

The installation program installs Hadoop. When you see output similar to the following, the installation is finished:

Installation complete. (HADOOP\_HOME=/opt/hadoop/hadoop-2.4.0)

-->Follow the remaining instructions in your installation guide.

- 9 Proceed to [“Overview of Post-Installation Steps for Hadoop”](#).

## Post-Installation Steps for Hadoop

### Overview of Post-Installation Steps for Hadoop

You must perform these manual steps after installing SAS High-Performance Deployment of Hadoop:

- 1 Use the appropriate user ID when invoking these processes:
  - Run HDFS commands as user ID **hdfs**.
  - Run as YARN as user ID **yarn**.
  - Run the Map Reduce Jobhistory Server as user ID **mapred**.
- 2 Define the environment variable, HADOOP\_HOME:
 

```
export HADOOP_HOME=/hadoop-installation-directory/hadoop-2.4.0
```
- 3 [Format the NameNode](#).
- 4 [If you are using more than one data device, update hdfs-site.xml and push it to each machine in the cluster.](#)
- 5 [If you are implementing Kerberos, see “Post-Installation Configuration Changes to Hadoop for Kerberos” on page 53.](#)
- 6 Start Hadoop.
  - With Kerberos:
 

[See “Start HDFS \(with Kerberos\)” on page 58](#)
  - Without Kerberos:
 

[See “Start HDFS \(without Kerberos\)” on page 59](#)
- 7 [Check the HDFS filesystem and create HDFS directories.](#)
- 8 [Validate your Hadoop deployment.](#)
- 9 [If your deployment includes SAS/ACCESS Interface to Hadoop, install the SAS Embedded Process on your Hadoop machine cluster. For more information, see Chapter 6, “Deploying SAS Embedded Process for Hadoop,” on page 85.](#)

### Format the Hadoop NameNode

To format the SAS High-Performance Deployment of Hadoop NameNode, follow these steps:

- 1 Change to the **hdfs** user account:

```
su - hdfs
```

**2** Export the HADOOP\_HOME environment variable.

For example:

```
export "HADOOP_HOME=/hadoop/hadoop-2.4.0"
```

**3** Format the NameNode:

```
$HADOOP_HOME/bin/hadoop namenode -format
```

**4** At the Re-format filesystem in `/hadoop-install-dir/hadoop-name ? (Y or N)` prompt, enter **y**. A line similar to the following highlighted output indicates that the format is successful:

```
Formatting using clusterid: CID-5b96061a-79f4-4264-87e0-99f351b749af
14/06/12 17:17:02 INFO util.HostsFileReader:
Refreshing hosts (include/exclude) list
14/06/12 17:17:03 INFO blockmanagement.DatanodeManager:
dfs.block.invalidate.limit=1000
14/06/12 17:17:03 INFO util.GSet: VM type           = 64-bit
14/06/12 17:17:03 INFO util.GSet: 2% max memory = 19.33375 MB
14/06/12 17:17:03 INFO util.GSet: capacity       = 2^21 = 2097152 entries
14/06/12 17:17:03 INFO util.GSet: recommended=2097152, actual=2097152
14/06/12 17:17:03 INFO blockmanagement.BlockManager:
dfs.block.access.token.enable=false
14/06/12 17:17:03 INFO blockmanagement.BlockManager: defaultReplication = 2
14/06/12 17:17:03 INFO blockmanagement.BlockManager: maxReplication      = 512
14/06/12 17:17:03 INFO blockmanagement.BlockManager: minReplication      = 1
14/06/12 17:17:03 INFO blockmanagement.BlockManager:
maxReplicationStreams      = 2
14/06/12 17:17:03 INFO blockmanagement.BlockManager:
shouldCheckForEnoughRacks = false
14/06/12 17:17:03 INFO blockmanagement.BlockManager:
replicationRecheckInterval = 3000
14/06/12 17:17:03 INFO namenode.FSNamesystem: fsOwner=nn/node1.domain.net@DOMAIN.NET (auth:KERBEROS)
14/06/12 17:17:03 INFO namenode.FSNamesystem: supergroup=supergroup
14/06/12 17:17:03 INFO namenode.FSNamesystem: isPermissionEnabled=true
14/06/12 17:17:03 INFO namenode.NameNode:
Caching file names occurring more than 10 times
14/06/12 17:17:04 INFO namenode.NNStorage: Storage directory
/hadoop/hadoop-name has been successfully formatted.
14/06/12 17:17:04 INFO namenode.FSImage: Saving image file
/hadoop/hadoop-name/current/fsimage.ckpt_000000000000000000 using no compression
14/06/12 17:17:04 INFO namenode.FSImage: Image file of size 119 saved in 0 seconds.
14/06/12 17:17:04 INFO namenode.NNStorageRetentionManager:
Going to retain 1 images with txid >= 0
14/06/12 17:17:04 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at node1.domain.net/192.0.0.0
*****/
```

**Note:** Without Kerberos, the log record for fsOwner is similar to the following:

```
14/06/12 17:17:03 INFO namenode.FSNamesystem: fsOwner=hdfs (auth:SIMPLE)
```

**5** Return to “[Overview of Post-Installation Steps for Hadoop](#)” on page 51 for instructions on creating well-known HDFS directories.

### (Optional) Deploy with Multiple Data Devices

If you plan to use more than one data device with the SAS High-Performance Deployment of Hadoop, then you must manually declare each device's Hadoop data directory in `hdfs-site.xml` and push it out to all of your DataNodes.

To deploy SAS High-Performance Deployment for Hadoop with more than one data device, follow these steps:

- 1 Log on to the Hadoop NameNode using the account with which you plan to run Hadoop.
- 2 In a text editor, open `hadoop-installation-directory/etc/hadoop/hdfs-site.xml`.
- 3 Locate the `dfs.data.dir` property, specify the location of your additional data devices' data directories, and save the file.

Separate multiple data directories with a comma.

For example:

```
<property>
  <name>dfs.data.dir</name>
  <value>/hadoop/hadoop-data,/data/dn</value>
</property>
```

- 4 Copy `hdfs-site.xml` to all of your Hadoop DataNodes using the `simcp` command.

For information about `simcp`, see [Appendix 2, "SAS High-Performance Analytics Infrastructure Command Reference,"](#) on page 129.

- 5 If you are using Kerberos, proceed to ["Post-Installation Configuration Changes to Hadoop for Kerberos"](#) on page 53.

Otherwise, proceed to ["Start HDFS \(without Kerberos\)"](#) on page 59.

### Post-Installation Configuration Changes to Hadoop for Kerberos

There are additional HDFS options not covered by the SAS Hadoop installer that need to be specified in order for Secure Mode Hadoop to work properly. Those additional options are defined in the various Hadoop configuration files. Your configuration files should match the ones below. You could copy and paste the files below and make environment-specific changes for the following items:

- hostnames
- JAVA\_HOME
- HADOOP\_HOME
- DOMAIN.NET is used as the example Kerberos realm

**Note:** Do not replace `_HOST`, as shown in the example files, with Kerberos principal names.

Be aware that you need to check and correct line breaks. Additions and changes relative to Secure Mode Hadoop are highlighted.

`hadoop-env.sh`:

```
export JAVA_HOME=/usr/java/latest
```

```

export HADOOP_HOME=/hadoop/hadoop-2.4.0
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export HADOOP_LOG_DIR=$HADOOP_HOME/logs/hdfs
for f in $HADOOP_HOME/contrib/capacity-scheduler/*.jar; do
    if [ "$HADOOP_CLASSPATH" ]; then
        export HADOOP_CLASSPATH=$HADOOP_CLASSPATH:$f
    else
        export HADOOP_CLASSPATH=$f
    fi
done

for f in $HADOOP_HOME/share/hadoop/sas/*.jar; do
    if [ "$HADOOP_CLASSPATH" ]; then
        export HADOOP_CLASSPATH=$HADOOP_CLASSPATH:$f
    else
        export HADOOP_CLASSPATH=$f
    fi
done

export HADOOP_COMMON_LIB_NATIVE_DIR=${HADOOP_PREFIX}/lib/native
export HADOOP_OPTS="$HADOOP_OPTS -Djava.net.preferIPv4Stack=true
-Djava.library.path=$HADOOP_PREFIX/lib"
export HADOOP_NAMENODE_OPTS="-Dhadoop.security.logger=
${HADOOP_SECURITY_LOGGER:
- INFO,RFAS} -Dhdfs.audit.logger=${HDFS_AUDIT_LOGGER:
-INFO,NullAppender} $HADOOP_NAMENODE_OPTS"
export HADOOP_DATANODE_OPTS="-Dhadoop.security.logger=
ERROR,RFAS $HADOOP_DATANODE_OPTS"
export HADOOP_SECONDARYNAMENODE_OPTS="-Dhadoop.security.logger=
${HADOOP_SECURITY_LOGGER:-INFO,RFAS} -Dhdfs.audit.logger=
${HDFS_AUDIT_LOGGER:-INFO,NullAppender} $HADOOP_SECONDARYNAMENODE_OPTS"
export HADOOP_CLIENT_OPTS="-Xmx512m $HADOOP_CLIENT_OPTS"
export HADOOP_SECURE_DN_USER=hdfs
export JSVC_HOME=/hadoop/hadoop-2.4.0/sbin
export HADOOP_SECURE_DN_LOG_DIR=${HADOOP_LOG_DIR}/${HADOOP_HDFS_USER}
export HADOOP_PID_DIR=/hadoop/hadoop-2.4.0/tmp
export HADOOP_SECURE_DN_PID_DIR=${HADOOP_PID_DIR}
export HADOOP_IDENT_STRING=$USER

```

#### core-site.xml:

```

<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!-- Put site-specific property overrides in this file. -->
<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://node1.domain.net:54310</value>
  </property>
  <property>
    <name>io.file.buffer.size</name>
    <value>102400</value>
  </property>
  <property>
    <name>hadoop.tmp.dir</name>
    <value>/hadoop/hadoop-2.4.0/tmp</value>
  </property>

```

```

<property>
  <name>hadoop.security.authentication</name>
  <value>kerberos</value>
</property>
<property>
  <name>hadoop.security.authorization</name>
  <value>true</value>
</property>
<property>
  <name>hadoop.rpc.protection</name>
  <value>privacy</value>
</property>
<property>
  <name>hadoop.ssl.require.client.cert</name>
  <value>false</value>
</property>
<property>
  <name>hadoop.ssl.hostname.verifier</name>
  <value>DEFAULT</value>
</property>
<property>
  <name>hadoop.ssl.keystores.factory.class</name>
  <value>org.apache.hadoop.security.ssl.FileBasedKeyStoresFactory</value>
</property>
<property>
  <name>hadoop.ssl.server.conf</name>
  <value>ssl-server.xml</value>
</property>
<property>
  <name>hadoop.ssl.client.conf</name>
  <value>ssl-client.xml</value>
</property>
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[2:$1;$2] (^dn;.*$) s/^.*$/hdfs/
    RULE:[2:$1;$2] (^sn;.*$) s/^.*$/hdfs/
    RULE:[2:$1;$2] (^nn;.*$) s/^.*$/hdfs/
    RULE:[1:$1@$0] (.@DOMAIN.NET) s/@.//
    DEFAULT
  </value>
</property>
</configuration>

```

#### hdfs-site.xml:

```

<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!-- Put site-specific property overrides in this file. -->
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>2</value>
    <description>Default block replication.
      The actual number of replications can be specified when the file is created.
      The default is used if replication is not specified in create time.
    </description>
  </property>

```

```

</property>
<property>
  <name>dfs.namenode.name.dir</name>
  <value>file:///hadoop/hadoop-name</value>
</property>
<property>
  <name>dfs.datanode.data.dir</name>
  <value>/hadoop/hadoop-data</value>
</property>
<property>
  <name>dfs.permissions.enabled</name>
  <value>true</value>
</property>
<property>
  <name>dfs.namenode.plugins</name>
  <value>com.sas.lasr.hadoop.NameNodeService</value>
</property>
<property>
  <name>dfs.datanode.plugins</name>
  <value>com.sas.lasr.hadoop.DataNodeService</value>
</property>
<property>
  <name>com.sas.lasr.hadoop.fileinfo</name>
  <value>ls -l {0}</value>
  <description>The command used to get the user, group, and permission
  information for a file.
  </description>
</property>
<property>
  <name>com.sas.lasr.service.allow.put</name>
  <value>true</value>
  <description>Flag indicating whether the PUT command is enabled when
  running as a service. The default is false.
  </description>
</property>
<property>
  <name>dfs.namenode.https-address</name>
  <value>0.0.0.0:50470</value>
</property>
<property>
  <name>dfs.datanode.https.address</name>
  <value>0.0.0.0:50475</value>
</property>
<property>
  <name>dfs.datanode.ipc.address</name>
  <value>0.0.0.0:50020</value>
</property>
<property>
  <name>dfs.namenode.http-address</name>
  <value>0.0.0.0:50070</value>
</property>
<property>
  <name>dfs.secondary.http.address</name>
  <value>0.0.0.0:50090</value>
</property>
<property>

```



```

    <name>dfs.namenode.backup.address</name>
    <value>0.0.0.0:50100</value>
  </property>
  <property>
    <name>dfs.namenode.backup.http-address</name>
    <value>0.0.0.0:50105</value>
  </property>
  <property>
    <name>com.sas.lasr.hadoop.service.namenode.port</name>
    <value>15452</value>
  </property>
  <property>
    <name>com.sas.lasr.hadoop.service.datanode.port</name>
    <value>15453</value>
  </property>
  <property>
    <name>dfs.namenode.fs-limits.min-block-size</name>
    <value>0</value>
  </property>
  <property>
    <name>dfs.block.access.token.enable</name>
    <value>true</value>
  </property>
  <property>
    <name>dfs.http.policy</name>
    <value>HTTPS_ONLY</value>
  </property>
  <property>
    <name>dfs.namenode.keytab.file</name>
    <value>/hadoop/hadoop-2.4.0/etc/hadoop/hdfs_nnservice.keytab</value>
  </property>
  <property>
    <name>dfs.namenode.kerberos.principal</name>
    <value>nn/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.namenode.kerberos.https.principal</name>
    <value>host/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.secondary.namenode.https-port</name>
    <value>50471</value>
  </property>
  <property>
    <name>dfs.secondary.namenode.keytab.file</name>
    <value>/hadoop/hadoop-2.4.0/etc/hadoop/hdfs_nnservice.keytab</value>
  </property>
  <property>
    <name>dfs.secondary.namenode.kerberos.principal</name>
    <value>sn/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.secondary.namenode.kerberos.https.principal</name>
    <value>host/_HOST@DOMAIN.NET</value>
  </property>
  <property>

```

```

    <name>dfs.datanode.data.dir.perm</name>
    <value>700</value>
  </property>
  <property>
    <name>dfs.datanode.address</name>
    <value>0.0.0.0:1004</value>
  </property>
  <property>
    <name>dfs.datanode.http.address</name>
    <value>0.0.0.0:1006</value>
  </property>
  <property>
    <name>dfs.datanode.keytab.file</name>
    <value>/hadoop/hadoop-2.4.0/etc/hadoop/hdfs_dnservice.keytab</value>
  </property>
  <property>
    <name>dfs.datanode.kerberos.principal</name>
    <value>dn/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.datanode.kerberos.https.principal</name>
    <value>host/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.encrypt.data.transfer</name>
    <value>true</value>
  </property>
  <property>
    <name>dfs.webhdfs.enabled</name>
    <value>true</value>
  </property>
  <property>
    <name>dfs.web.authentication.kerberos.principal</name>
    <value>HTTP/_HOST@DOMAIN.NET</value>
  </property>
  <property>
    <name>dfs.web.authentication.kerberos.keytab</name>
    <value>/hadoop/hadoop-2.4.0/etc/hadoop/http.keytab</value>
  </property>
</configuration>

```

Proceed to “[Start HDFS \(with Kerberos\)](#)” on page 58 .

## Start HDFS (with Kerberos)

To start HDFS using Kerberos, follow these steps:

- 1 Connect to all machines using SSH as the `hdfs` user.
- 2 Log on to the NameNode machine as the `hdfs` user and start the NameNode. For example:

```

export HADOOP_HOME=/hadoop/hadoop-2.4.0
$HADOOP_HOME/sbin/hadoop-daemon.sh start namenode

```

This command starts the NameNode and the Secondary NameNode.

- 3 To start the DataNodes, log on to the first DataNode as the `root` user.

- 4 Run the following commands on each DataNode in the cluster:

```
export HADOOP_HOME=/hadoop/hadoop-2.4.0
$HADOOP_HOME/sbin/hadoop-daemon.sh start datanode
```

You start the process with the `root` user, but it switches to the user ID specified in the `HADOOP_SECURE_DN_USER` variable from the `hadoop-env.sh` file.

All secure DataNodes should be running.

- 5 Proceed to [“Check the HDFS Filesystem and Create HDFS Directories”](#).

### Start HDFS (without Kerberos)

Log on to the NameNode machine as the `hdfs` user and start HDFS. For example:

```
export HADOOP_HOME=/hadoop/hadoop-2.4.0
$HADOOP_HOME/sbin/start-dfs.sh
```

This command starts the NameNode, Secondary NameNode, and the DataNodes in the cluster.

Proceed to [“Check the HDFS Filesystem and Create HDFS Directories”](#).

### Check the HDFS Filesystem and Create HDFS Directories

To perform a filesystem check and create the initial HDFS directories, follow these steps:

- 1 Log on to the NameNode as the `hdfs` user.

**Note:** If you are using Kerberos, then use `kinit` to get a ticket. For example:  
`kinit hdfs@DOMAIN.NET.`

- 2 Run the following commands to check the filesystem and display the number of DataNodes:

```
export HADOOP_HOME=/hadoop/hadoop-2.4.0
$HADOOP_HOME/bin/hadoop fsck /
```

- 3 Run the following command to create the directories in HDFS:

```
$HADOOP_HOME/sbin/initial-sas-hdfs-setup.sh
```

- 4 Run the following command to verify the directories have been created:

```
$HADOOP_HOME/bin/hadoop fs -ls /
```

You should output similar to the following:

```
drwxrwxrwx - hdfs supergroup 0 2014-07-24 13:29 /hps
drwxrwxrwx - hdfs supergroup 0 2014-07-23 13:59 /test
drwxrwxrwt - hdfs supergroup 0 2014-07-24 13:29 /tmp
drwxr-xr-x - hdfs supergroup 0 2014-07-24 13:29 /user
drwxrwxrwt - hdfs supergroup 0 2014-07-24 13:29 /vapublic
```

- 5 Proceed to [“Validate Your Hadoop Deployment”](#).

## Validate Your Hadoop Deployment

You can confirm that Hadoop is running successfully by opening a browser to `http://NameNode:50070/dfshealth.jsp`. Review the information in the cluster summary section of the page. Confirm that the number of live nodes equals the number of DataNodes and that the number of dead nodes is zero.

**Note:** It can take a few seconds for each node to start. If you do not see every node, then refresh the connection in the web interface.

---

## Configuring Existing Hadoop Clusters

### Overview of Configuring Existing Hadoop Clusters

If your site uses a Hadoop implementation that is supported, then you can configure your Hadoop cluster for use with the SAS High-Performance Analytics environment.

The following steps are needed to configure your existing Hadoop cluster:

- 1 Make sure that your Hadoop deployment meets the analytic environment prerequisites. For more information, see [“Prerequisites for Existing Hadoop Clusters” on page 60](#)
- 2 Follow steps specific to your implementation of Hadoop:
  - [“Configuring the Existing Cloudera Hadoop Cluster” on page 61](#)
  - [“Configuring the Existing Hortonworks Data Platform Hadoop Cluster” on page 65](#)
  - [“Configuring the Existing IBM BigInsights Hadoop Cluster” on page 67](#)
  - [“Configuring the Existing MapR Hadoop Cluster” on page 68](#)
  - [“Configuring the Existing Pivotal HD Hadoop Cluster” on page 68](#)

### Prerequisites for Existing Hadoop Clusters

The following is required for existing Hadoop clusters that will be configured for use with the SAS High-Performance Analytics environment:

- Each machine machine in the cluster must be able to resolve the host name of all the other machines.
- The NameNode and secondary NameNode are not defined as the same host.
- The NameNode host does not also have a DataNode configured on it.
- For Kerberos, in the SAS High-Performance Analytics environment, `/etc/hosts` must contain the machine names in the cluster in this order: short name, fully qualified domain name.
- Time must be synchronized across all machines in the cluster.
- (Cloudera 5 only) Make sure that all machines configured for the SAS High-Performance Analytics environment are in the same role group.

## Configuring the Existing Cloudera Hadoop Cluster

### Managing Cloudera Configuration Priorities

Cloudera uses the Linux `alternatives` command for client configuration files. Therefore, make sure that the client configuration path has the highest priority for all machines in the cluster. (Often, the mapreduce client configuration has a higher priority over the hdfs configuration.)

If the output of the command `alternatives --display hadoop-conf` returns the Cloudera server configuration, or if mapreduce client configuration has priority over the client configuration, you will experience problems because SAS makes additions to the client configuration. For more information about `alternatives`, refer to its man page.

### Configure the Existing Cloudera Hadoop Cluster, Version 5

Use the Cloudera Manager to configure your existing Cloudera 5 Hadoop deployment to interoperate with the SAS High-Performance Analytics environment.

- 1 Untar the SAS High-Performance Deployment for Hadoop tarball, and propagate three files (identified in the following steps) on every machine in your Cloudera Hadoop cluster:
  - a Navigate to the SAS High-Performance Deployment for Hadoop tarball in your SAS Software depot:

```
cd depot-installation-location/standalone_installs/
SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_x64/
```

- b Copy `sashadoop.tar.gz` to a temporary location where you have Write access.

- c Untar `sashadoop.tar.gz`:

```
tar xzf sashadoop.tar.gz
```

- d If not already done, set the following environment variables before running the Hadoop commands.

```
export JAVA_HOME=/path-to-java
export HADOOP_HOME=/opt/cloudera/parcels/CDH-5.0.0-0.cdh5b1.p0.57/lib/hadoop
```

- e Locate `sas.lasr.jar` and `sas.lasr.hadoop.jar` and propagate these two JAR files to every machine in the Cloudera Hadoop cluster into the CDH library path.

**TIP** If you have already installed the SAS High-Performance Computing Management Console or the SAS High-Performance Analytics environment, you can issue a single `simcp` command to propagate JAR files across all machines in the cluster. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.jar
/opt/cloudera/parcels/CDH-5.0.0-0.cdh5b1.p0.57/lib/hadoop/lib/
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.hadoop.jar
```

```
/opt/cloudera/parcels/CDH-5.0.0-0.cdh5b1.p0.57/lib/hadoop/lib/
```

For more information, see [Appendix 2, “SAS High-Performance Analytics Infrastructure Command Reference,”](#) on page 129.

- f** Locate `saslasrfd` and propagate this file to every machine in the Cloudera Hadoop cluster into the CDH `bin` directory. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp saslasrfd
/opt/cloudera/parcels/CDH-5.0.0-0.cdh5b1.p0.57/lib/hadoop/bin/
```

- 2** Log on to the Cloudera Manager as an administrator.
- 3** Add the following to the plug-in configuration for the NameNode:

```
com.sas.lasr.hadoop.NameNodeService
```

- 4** Add the following to the plug-in configuration for DataNodes:

```
com.sas.lasr.hadoop.DataNodeService
```

- 5** Add the following lines to the advanced configuration for service-wide. These lines are placed in the HDFS Service Advanced Configuration Snippet (Safety Valve) for `hdfs-site.xml`:

```
<property>
<name>com.sas.lasr.service.allow.put</name>
<value>true</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name> dfs.namenode.fs-limits.min-block-size</name>
<value>0</value>
</property>
```

- 6** Add the following property to the **HDFS Client Configuration Safety Valve** under **Advanced** within the **Gateway Default Group**. Make sure that you change *path-to-data-dir* to the data directory location for your site (for example, `<value>file://dfs/dn</value>`):

```
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name> dfs.namenode.fs-limits.min-block-size</name>
<value>0</value>
</property>
```

```
<property>
<name>dfs.datanode.data.dir</name>
<value>file://path-to-data-dir</value>
</property>
```

- 7 Add the location of JAVA\_HOME to the **HDFS Client Environment Advanced Configuration Snippet for `hadoop-env.sh` (Safety Valve)**, located under **Advanced** in the **Gateway Default Group**. For example:
 

```
JAVA_HOME=/usr/lib/java/jdk1.7.0_07
```
- 8 Save your changes and deploy the client configuration to each host in the cluster.
- 9 Restart the HDFS service and any dependencies in Cloudera Manager.
- 10 Create and mode the `/test` directory in HDFS for testing the cluster with SAS test jobs. You might need to set HADOOP\_HOME first, and you must run the following commands as the user running HDFS (typically, `hdfs`).

```
$HADOOP_HOME/bin/hadoop fs -mkdir /test
```

```
$HADOOP_HOME/bin/hadoop fs -chmod 777 /test
```

- 11 Make sure that the client configuration path has the highest priority for all machines in the cluster. For more information, see [“Managing Cloudera Configuration Priorities” on page 61](#).

## Configure the Existing Cloudera Hadoop Cluster, Version 4

Use the Cloudera Manager to configure your existing Cloudera 4 Hadoop deployment to interoperate with the SAS High-Performance Analytics environment.

**TIP** In Cloudera 4.2 and earlier, you must install the enterprise license, even if you are below the stated limit of 50 nodes in the Hadoop cluster for requiring a license.

- 1 Untar the SAS High-Performance Deployment for Hadoop tarball, and propagate three files (identified in the following steps) on every machine in your Cloudera Hadoop cluster:
  - a Navigate to the SAS High-Performance Deployment for Hadoop tarball in your SAS Software depot:
 

```
cd depot-installation-location/standalone_installs/
   SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_x64/
```
  - b Copy `sashadoop.tar.gz` to a temporary location where you have Write access.
  - c Untar `sashadoop.tar.gz`:
 

```
tar xzf sashadoop.tar.gz
```
  - d Locate `sas.lasr.jar` and `sas.lasr.hadoop.jar` and propagate these two JAR files to every machine in the Cloudera Hadoop cluster into the CDH library path.

**TIP** If you have already installed the SAS High-Performance Computing Management Console or the SAS High-Performance Analytics environment, you can issue a single `simcp` command to propagate JAR files across all machines in the cluster. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.jar
/opt/cloudera/parcels/CDH-4.4.0-1.cdh4.4.0.p0.39/lib/hadoop/lib/
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.hadoop.jar
/opt/cloudera/parcels/CDH-4.4.0-1.cdh4.4.0.p0.39/lib/hadoop/lib/
```

For more information, see [Appendix 2, “SAS High-Performance Analytics Infrastructure Command Reference,”](#) on page 129.

- e Locate `saslasrfd` and propagate this file to every machine in the Cloudera Hadoop cluster into the CDH `bin` directory. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp saslasrfd
/opt/cloudera/parcels/CDH-4.4.0-1.cdh4.4.0.p0.39/lib/hadoop/bin/
```

- 2 Log on to the Cloudera Manager as an administrator.
- 3 Add the following to the plug-in configuration for the NameNode:

```
com.sas.lasr.hadoop.NameNodeService
```

- 4 Add the following to the plug-in configuration for DataNodes:

```
com.sas.lasr.hadoop.DataNodeService
```

- 5 Add the following lines to the advanced configuration for service-wide. These lines are placed in the **HDFS Service Configuration Safety Valve** property for `hdfs-site.xml`:

```
<property>
<name>com.sas.lasr.service.allow.put</name>
<value>true</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name>dfs.namenode.fs-limits.min-block-size</name>
<value>0</value>
</property>
```

- 6 Restart all Cloudera Manager services.
- 7 Create and set the mode for the `/test` directory in HDFS for testing. You might need to set `HADOOP_HOME` first, and you must run the following commands as the user running HDFS (normally, the `hdfs` user).
- 8 If needed, set the following environment variables before running the Hadoop commands.



```
export HADOOP_HOME=/opt/cloudera/parcels/CDH-4.4.0-1.cdh4.4.0.p0.39/lib/hadoop
```

- 9 Run the following commands to create the `/test` directory in HDFS. This directory is to be used for testing the cluster with SAS test jobs.

```
$HADOOP_HOME/bin/hadoop fs -mkdir /test
```

```
$HADOOP_HOME/bin/hadoop fs -chmod 777 /test
```

- 10 Add the following to the **HDFS Client Configuration Safety Valve**:

```
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name>dfs.datanode.data.dir</name>
<value>file:///hadoop/hadoop-data</value>
</property>
```

- 11 Add the location of `JAVA_HOME` to the **Client Environment Safety Valve** for `hadoop-env.sh`. For example:

```
JAVA_HOME=/usr/lib/java/jdk1.7.0_07
```

- 12 Save your changes and deploy the client configuration to each host in the cluster.
- 13 Make sure that the client configuration path has the highest priority for all machines in the cluster. For more information, see [“Managing Cloudera Configuration Priorities” on page 61](#).

**TIP** Remember the value of `HADOOP_HOME` as the SAS High-Performance Analytics environment prompts for this during its install. By default, these are the values for Cloudera:

- Cloudera 4.5:

```
/opt/cloudera/parcels/CDH-4.5.0-1.cdh4.5.0.p0.30
```

- Cloudera 4.2 and earlier:

```
/opt/cloudra/parcels/CDH-4.2.0-1.cdh4.2.0.po.10/lib/Hadoop
```

## Configuring the Existing Hortonworks Data Platform Hadoop Cluster

Use the Ambari interface to configure your existing Hortonworks Data Platform deployment to interoperate with the SAS High-Performance Analytics environment.

- 1 Log on to Ambari as an administrator, and stop all HDP services.

- 2 Untar the SAS High-Performance Deployment for Hadoop tarball, and propagate three files (identified in the following steps) on every machine in your Hortonworks Hadoop cluster:

- a Navigate to the SAS High-Performance Deployment for Hadoop tarball in your SAS Software depot:

```
cd depot-installation-location/standalone_installs/
SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_x64/
```

- b Copy sashadoop.tar.gz to a temporary location where you have Write access.

- c Untar sashadoop.tar.gz:

```
tar xzf sashadoop.tar.gz
```

- d Locate sas.lasr.jar and sas.lasr.hadoop.jar and propagate these two JAR files to every machine in the HDP cluster into the HDP library path.

**TIP** If you have already installed the SAS High-Performance Computing Management Console or the SAS High-Performance Analytics environment, you can issue a single `simcp` command to propagate JAR files across all machines in the cluster. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.jar
/usr/lib/hadoop/lib/
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.hadoop.jar
/usr/lib/hadoop/lib/
```

For more information, see [Appendix 2, “SAS High-Performance Analytics Infrastructure Command Reference,”](#) on page 129.

- e Locate saslasrfd and propagate this file to every machine in the HDP cluster into the HDP `bin` directory. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp saslasrfd /usr/lib/hadoop/bin/
```

- 3 In the Ambari interface, create a custom `hdfs-site.xml` and add the following properties:

**dfs.namenode.plugins**

```
com.sas.lasr.hadoop.NameNodeService
```

**dfs.datanode.plugins**

```
com.sas.lasr.hadoop.DataNodeService
```

**com.sas.lasr.hadoop.fileinfo**

```
ls -l {0}
```

**com.sas.lasr.service.allow.put**

```
true
```

**com.sas.lasr.hadoop.service.namenode.port**

```
15452
```

**com.sas.lasr.hadoop.service.datanode.port**

```
15453
```

**dfs.namenode.fs-limits.min-block-size**

```
0
```

- 4 Save the properties and start the HDFS service.

- 5 Run the following commands as the `hdfs` user to create the `/test` directory in HDFS. This directory is used for testing your cluster with SAS test jobs.

```
hadoop fs -mkdir /test
hadoop fs -chmod 777 /test
```

## Configuring the Existing IBM BigInsights Hadoop Cluster

To configure your existing IBM BigInsights Hadoop deployment to interoperate with the SAS High-Performance Analytics environment, follow these steps:

- 1 Untar the SAS High-Performance Deployment for Hadoop tarball, and propagate three files (identified in the following steps) on every machine in your BigInsights Hadoop cluster:

- a Navigate to the SAS High-Performance Deployment for Hadoop tarball in your SAS Software depot:

```
cd depot-installation-location/standalone_installs/
SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_x64/
```

- b Copy `sashadoop.tar.gz` to a temporary location where you have Write access.

- c Untar `sashadoop.tar.gz`:

```
tar xzf sashadoop.tar.gz
```

- d Locate `sas.lasr.jar` and `sas.lasr.hadoop.jar` and propagate these two JAR files to every machine in the BigInsights cluster into the library path.

**Note:** `HADOOP_HOME`: default location: `/opt/ibm/biginsights/INC.BIGINSIGHT_HOME`; default location: `/opt/ibm/biginsights`.

**TIP** If you have already installed the SAS High-Performance Computing Management Console or the SAS High-Performance Analytics environment, you can issue a single `simcp` command to propagate JAR files across all machines in the cluster. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.jar
$HADOOP_HOME/share/hadoop/hdfs/libs
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.hadoop.jar
$HADOOP_HOME/share/hadoop/hdfs/libs
```

For more information, see [Appendix 2, “SAS High-Performance Analytics Infrastructure Command Reference,”](#) on page 129.

- e Locate `saslasrfd` and propagate this file to every machine in the BigInsights cluster into the `$HADOOP_HOME/bin` directory. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp saslasrfd $HADOOP_HOME/bin
```

- 2 On the machine where you initially installed BigInsights, add the following properties for SAS for the HDFS configuration to the file `$BIGINSIGHT_HOME/hdm/hadoop-conf-staging/hdfs-site.xml`. Adjust values appropriately for your deployment:

```
<property>
```

```

<name>dfs.datanode.plugins</name>
<value>com.sas.lasr.hadoop.DataNodeService</value>
</property>
<property>
<name>com.sas.lasr.service.allow.put</name>
<value>true</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name>dfs.namenode.fs-limits.min-block-size</name>
<value>0</value>
</property>

```

- 3 Synchronize this new configuration by running the following command on the machine where you initially deployed BigInsights:

```
$BIGINSIGHT_HOME/bin/synconf.sh.
```

- 4 On the machine where you initially deployed BigInsights, log on as the biadmin user and run the following commands to restart the cluster with the new configuration:

```
stop-all.sh
```

```
start-all.sh
```

- 5 Note the location of HADOOP\_HOME. You will need to refer to this value when installing the SAS High-Performance Analytics environment.
- 6 Run the following commands as the hdfs user to create the /test directory in HDFS. This directory is used for testing your cluster with SAS test jobs.

```
hadoop fs -mkdir /test
```

```
hadoop fs -chmod 777 /test
```

## Configuring the Existing MapR Hadoop Cluster

Configuring your existing MapR deployment to interoperate with the SAS High-Performance Analytics environment consists of setting up a proper NFS mount. For more information, see <http://doc.mapr.com/display/MapR/Setting+Up+MapR+NFS>.

## Configuring the Existing Pivotal HD Hadoop Cluster

Use the Pivotal Command Center (PCC) to configure your existing Pivotal HD deployment to interoperate with the SAS High-Performance Analytics environment.

- 1 Log on to PCC as gpadmin. (The default password is gpadmin.)

- 2 Untar the SAS High-Performance Deployment for Hadoop tarball, and propagate three files (identified in the following steps) on every machine in your Cloudera Hadoop cluster:

- a Navigate to the SAS High-Performance Deployment for Hadoop tarball in your SAS Software depot:

```
cd depot-installation-location/standalone_installs/
SAS_High-Performance_Deployment_for_Hadoop/2_8/Linux_for_x64/
```

- b Copy sashadoop.tar.gz to a temporary location where you have Write access.

- c Untar sashadoop.tar.gz:

```
tar xzf sashadoop.tar.gz
```

- d Locate sas.lasr.jar and sas.lasr.hadoop.jar and propagate these two JAR files to every machine in the Pivotal HD cluster into the library path.

**TIP** If you have already installed the SAS High-Performance Computing Management Console or the SAS High-Performance Analytics environment, you can issue a single `simcp` command to propagate JAR files across all machines in the cluster. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.jar
/usr/lib/gphd/hadoop/lib/
/opt/sashpcmc/opt/webmin/utilbin/simcp sas.lasr.hadoop.jar
/usr/lib/gphd/hadoop/lib/
```

For more information, see [Appendix 2, “SAS High-Performance Analytics Infrastructure Command Reference,”](#) on page 129.

- e Locate saslasrfd and propagate this file to every machine in the Pivotal HD cluster into the Pivotal HD `bin` directory. For example:

```
/opt/sashpcmc/opt/webmin/utilbin/simcp saslasrfd /usr/lib/gphd/hadoop/bin/
```

- 3 In the PCC, for YARN, make sure that Resource Manager, History Server, and Node Managers have unique host names.
- 4 In the PCC, make sure that the Zookeeper Server contains a unique host name.
- 5 Add the following properties for SAS for the HDFS configuration to the file `hdfs-site.xml`:

```
<property>
<name>dfs.datanode.plugins</name>
<value>com.sas.lasr.hadoop.DataNodeService</value>
</property>
<property>
<name>com.sas.lasr.service.allow.put</name>
<value>true</value>
</property>
<property>
<name>com.sas.lasr.hadoop.service.namenode.port</name>
<value>15452</value>
</property>
<property>
```

```
<name>com.sas.lasr.hadoop.service.datanode.port</name>
<value>15453</value>
</property>
<property>
<name> dfs.namenode.fs-limits.min-block-size</name>
<value>0</value>
</property>
```

- 6** Save your changes and deploy.
- 7** Restart your cluster using PCC and verify that HDFS is running in the dashboard.
- 8** Run the following commands as the `gpadmin` user to create the `/test` directory in HDFS. This directory is used for testing your cluster with SAS test jobs.

```
hadoop fs -mkdir /test
hadoop fs -chmod 777 /test
```

## 5

## Deploying the SAS High-Performance Analytics Environment

<i>Infrastructure Deployment Process Overview</i> .....	71
<i>Overview of Deploying the Analytics Environment</i> .....	72
<i>Encrypting SASHDAT Files</i> .....	75
<i>Install the Analytics Environment</i> .....	76
<i>Configuring the Analytics Environment for SASHDAT Encryption</i> .....	79
<i>Validating the Analytics Environment Deployment</i> .....	80
Overview of Validating .....	80
Use simsh to Validate .....	80
Use MPI to Validate .....	81
<i>Resource Management for the Analytics Environment</i> .....	81
Resource Settings File .....	81
Request Memory with TKMPI_INFO .....	83

---

### Infrastructure Deployment Process Overview

Installing and configuring the SAS High-Performance Analytics environment is the sixth of eight steps.

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.
4. (Optional) Deploy SAS High-Performance Computing Management Console.
5. (Optional) Deploy co-located Hadoop.
- ▶ **6. Deploy the SAS High-Performance Analytics environment.**
7. (Optional) Deploy the SAS Embedded Process for Hadoop.
8. (Optional) Configure the analytics environment for a remote parallel connection.

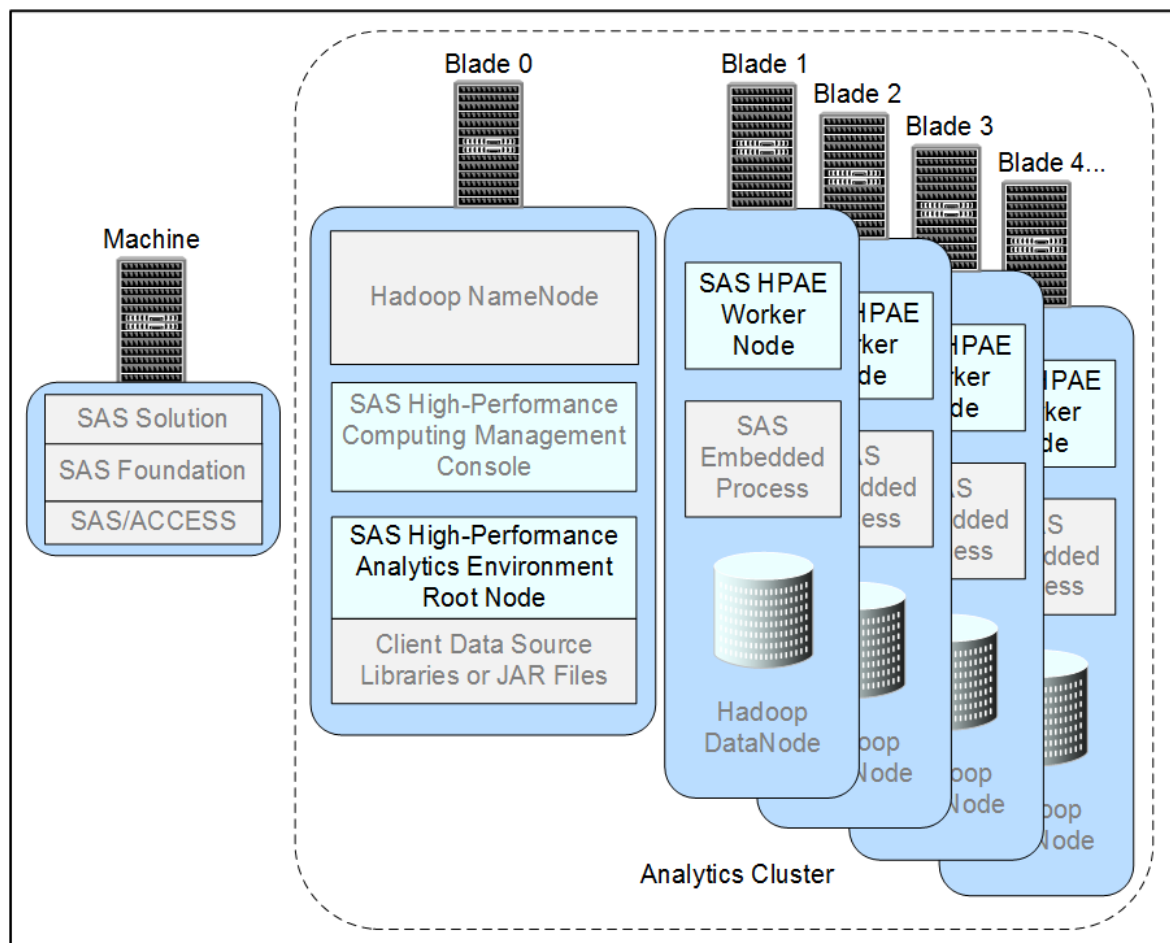
This chapter describes how to install and configure all of the components for the SAS High-Performance Analytics environment on the machines in the cluster.

## Overview of Deploying the Analytics Environment

Deploying the SAS High-Performance Analytics environment requires installing and configuring components on the root node machine and on the remaining machines in the cluster. In this document, the root node is deployed on blade 0.

The following figure shows the SAS High-Performance Analytics environment co-located on your Hadoop cluster:

**Figure 5.1** Analytics Environment Co-Located on the Hadoop Cluster

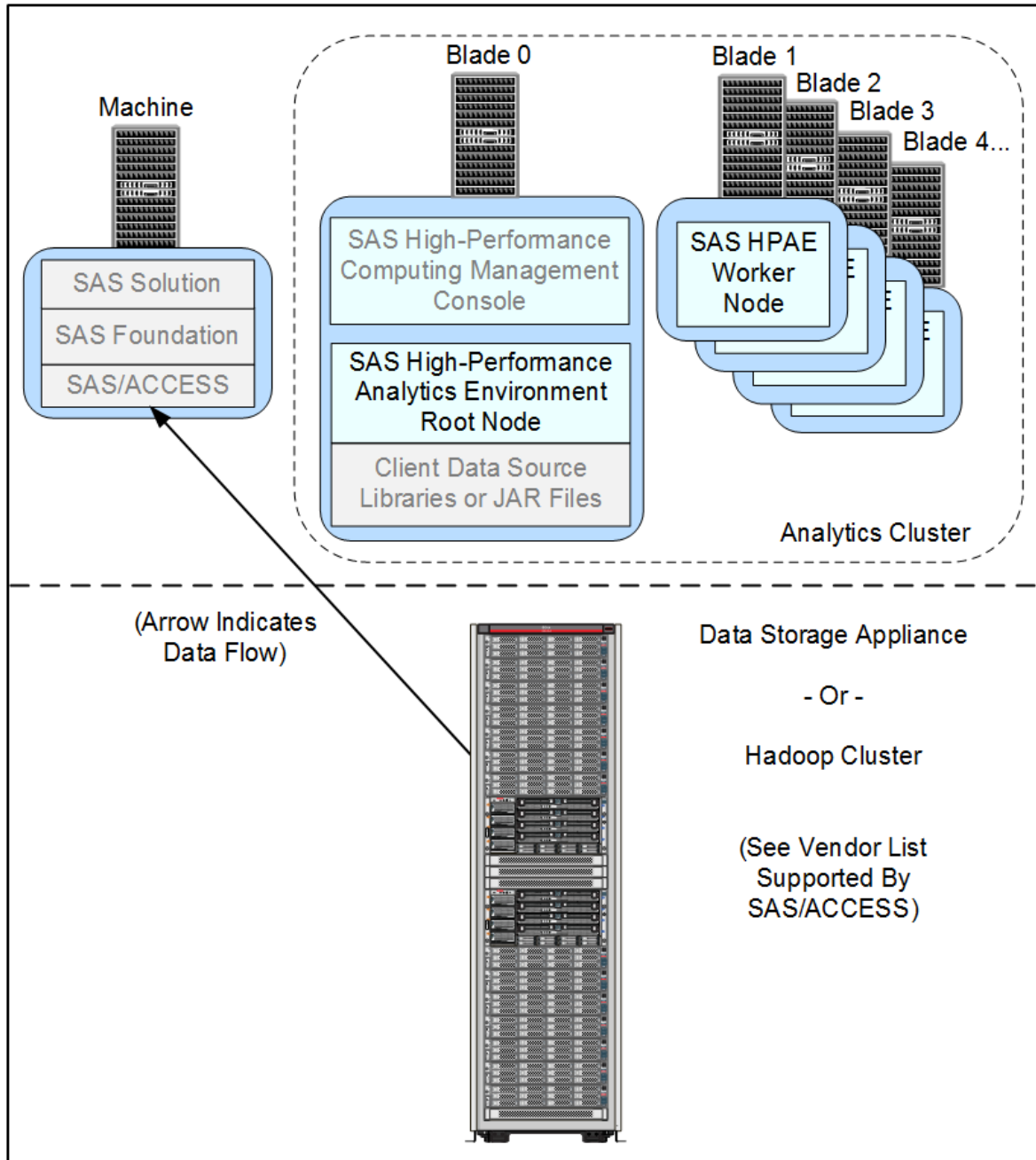


**Note:** For deployments that use Hadoop for the co-located data provider and access SASHDAT tables exclusively, SAS/ACCESS and SAS Embedded Process are not needed.



The following figure shows the SAS High-Performance Analytics environment using a serial connection through the SAS/ACCESS Interface to your remote data store:

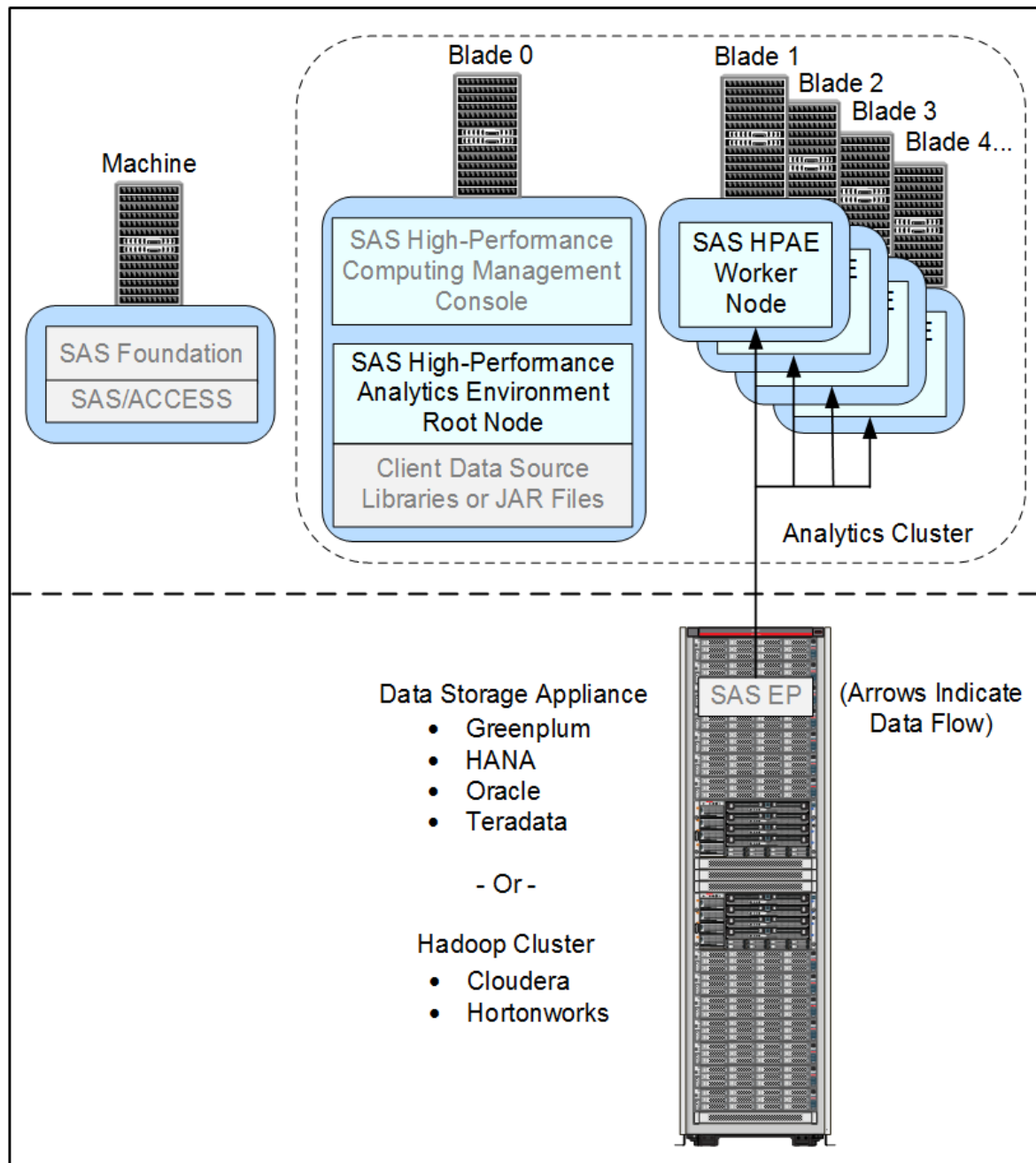
**Figure 5.2** Analytics Environment Remote from Your Data Store (Serial Connection)



**TIP** There might be solution-specific criteria that you should consider when determining your analytics cluster location. For more information, see the installation or administration guide for your specific SAS solution.

The following figure shows the SAS High-Performance Analytics environment using a parallel connection through the SAS Embedded Process to your remote data store:

**Figure 5.3** Analytics Environment Remote from Your Data Store (Parallel Connection)



The SAS High-Performance Analytics environment is packaged in separate executables. Refer to the following table for more information:

**Table 5.1** *Installation and Configuration Packages for the SAS High-Performance Analytics Environment*

Order to install	Filename	Purpose
1	<a href="#">TKGrid_Linux_x86_64.sh</a>	Analytics environment installation script for Red Hat Linux 6 and other equivalent, kernel-level Linux systems.
	<a href="#">TKGrid_Linux_x86_64_rhel5.sh</a>	Analytics environment installation script for Red Hat Linux (pre-version 6) and SUSE Linux 10 systems.
2	<a href="#">TKTGDat.sh</a>	SAS linguistic binary files required to perform text analysis in SAS LASR Analytic Server with SAS Visual Analytics and to run PROC HPTMINE and HPTMSCORE with SAS Text Miner.
3 (optional)	<a href="#">TKGrid_SEC_x86_64.sh</a>	Installation script for enabling the analytics environment to read and write encrypted SASHDAT files.
4 (optional)	<a href="#">TKGrid_REP_x86_64.sh</a>	Script for configuring the SAS High-Performance Analytics environment with a SAS Embedded Process for Red Hat Linux 6 and other equivalent, kernel-level Linux systems.
	<a href="#">TKGrid_REP_x86_64_rhel5.sh</a>	Script for configuring the SAS High-Performance Analytics environment with a SAS Embedded Process for Red Hat Linux (pre-version 6) and SUSE Linux 10 systems.

## Encrypting SASHDAT Files

In release 2.94, the SAS High-Performance Analytics environment supports reading and writing files using AES encryption with 256-bit keys. (This feature is very similar to the AES encryption provided by the SAS BASE Engine.) SASHDAT encryption is designed to bolster privacy protection for data at rest—that is, data stored in SASHDAT for analytic purposes.

Remember that SASHDAT data is typically not the system of record, but rather a copy of operational data that has been arranged for the purposes of analytics. In addition to encrypting data, many SAS users also anonymize their data when preparing it for analytics.

To enable the SAS High-Performance Analytics environment to read and write SASHDAT using encryption, you must install the TKGrid\_SEC package. For more information, see [“Configuring the Analytics Environment for SASHDAT Encryption” on page 79](#).

---

## Install the Analytics Environment

The SAS High-Performance Analytics environment components are installed with two shell scripts. Follow these steps to install:

- 1 Make sure that you have reviewed all of the information contained in the section [“Preparing to Deploy the SAS High-Performance Analytics Environment”](#) on page 28.
- 2 The software that is needed for the SAS High-Performance Analytics environment is available from within the SAS Software Depot that was created by the site depot administrator: *depot-installation-location/standalone\_installs/SAS\_High-Performance\_Node\_Installation/2\_94/Linux\_for\_x64*.
- 3 Copy the file that is appropriate for your operating system to the `/tmp` directory of the root node of the cluster:
  - Red Hat Linux (pre-version 6) and SUSE Linux 10:  
`TKGrid_Linux_x86_64_rhel5.sh`
  - Red Hat Linux 6 and other equivalent, kernel-level Linux systems:  
`TKGrid_Linux_x86_64.sh`
- 4 Copy `TKTGDat.sh` to the `/tmp` directory of the root node of the cluster.

**Note:** `TKTGDat.sh` contains the SAS linguistic binary files required to perform text analysis in SAS LASR Analytic Server with SAS Visual Analytics and to run PROC HPTMINE and HPTMSCORE with SAS Text Miner.
- 5 Log on to the machine that will serve as the root node of the cluster or the data appliance with a user account that has the necessary permissions.

For more information, see [“User Accounts for the SAS High-Performance Analytics Environment”](#) on page 28.
- 6 Change directories to the desired installation location, such as `/opt`.

Record the location of where you installed the analytics environment, as other configuration programs will prompt you for this path later in the deployment process.
- 7 Run the `TKGrid` shell script in this directory.

The shell script creates the `TKGrid` subdirectory and places all files under that directory.
- 8 Respond to the prompts from the shell script:

**Table 5.2** Configuration Parameters for the TKGrid Shell Script

Parameter	Description
Shared install or replicate to each node? (Y=SHARED/n=replicated)	If you are installing to a local drive on each node, then specify <b>n</b> and press Enter to indicate that this is a replicated installation. If you are installing to a drive that is shared across all the nodes (for example, NFS), then specify <b>y</b> and press Enter.
Enter additional paths to include in LD_LIBRARY_PATH, separated by colons (:)	If you have any external library paths that you want to be accessible to the SAS High-Performance Analytics environment, enter the paths here and press Enter. Otherwise, press Enter.
Enter NFS mount to MAPR directory (ie: /mapr/my.cluster.com, default is none).	<p>If you want the analytics environment to be able to read and write MapR data directly, enter the NFS mount here (for example, /mapr/my.cluster.com).</p> <p>The mount point must exist on all nodes, including the head node.</p> <p>For more information, see <a href="http://doc.mapr.com/display/MapR/Accessing+Data+with+NFS">http://doc.mapr.com/display/MapR/Accessing+Data+with+NFS</a>.</p>
Enter additional options to mpirun.	<p>If you have any mpirun options to add, specify them and press Enter.</p> <p>If you are using Kerberos, specify the following option and press Enter:</p> <pre>-genvlist `env   sed -e s/=.*//   sed /KRB5CCNAME/d   tr -d '\n' `TKPATH,LD_LIBRARY_PATH</pre> <p>If you have no additional options, press Enter.</p>
Enter path to use for Utility files. (default is /tmp).	<p>SAS High-Performance Analytics applications might write scratch files. By default, these files are created in the /tmp directory. To accept the default, press Enter. Or, to redirect the files to a different location, specify the path and press Enter.</p> <p><b>Note:</b> If the directory that you specified does not exist, you must create it manually.</p>
Enter path to Hadoop. (default is Hadoop not installed).	<p>If your site uses Hadoop, enter the installation directory (the value of the variable, HADOOP_HOME) and press Enter. If your site does not use Hadoop, press Enter.</p> <p>If you are using SAS High-Performance Deployment of Hadoop, use the directory that you specified earlier in <a href="#">Step 3 on page 48</a>.</p>
Force Root Rank to run on headnode? (y/N)	If the appliance resides behind a firewall and only the root node can connect back to the client machines, specify <b>y</b> and press Enter. Otherwise, specify <b>n</b> and press Enter.
Enter full path to machine list. The head node 'head-node-machine-name' should be listed first.	Specify the name of the file that you created in the section <a href="#">“List the Machines in the Cluster or Appliance”</a> (for example, /etc/gridhosts) and press Enter.

Parameter	Description
Enter maximum runtime for grid jobs (in seconds). Default 7200 (2 hours).	<p>If a SAS High-Performance Analytics application executes for more than the maximum allowable run time, it is automatically terminated. You can adjust that run-time limit here.</p> <p>To accept the default, press Enter. Or, specify a different maximum run time (in seconds) and press Enter.</p>
Enter value for UMASK. (default is unset.)	<p>To set <i>no</i> umask value, press Enter. Or, specify a umask value and press Enter.</p> <p>For more information, see <a href="#">“Consider Umask Settings” on page 29</a>.</p>

- 9** If you selected a replicated installation at the first prompt, you are now prompted to choose the technique for distributing the contents to the appliance nodes:

```
The install can now copy this directory to all the machines
listed in 'filename' using scp, skipping the first entry.
Perform copy?
(YES/no)
```

Press Enter if you want the installation program to perform the replication. Enter **no** if you are distributing the contents of the installation directory by some other technique.

- 10** Next, in the same directory from which you ran the TKGGrid shell script, run TKTGDat.sh.

The shell script creates the **TKTGDat** subdirectory and places all files in that directory.

- 11** Respond to the prompts from the shell script:

**Table 5.3** Configuration Prompts for the TKTG Dat Shell Script

Shared install or replicate to each node? (Y=SHARED/n=replicated)	If you are installing to a local drive on each node, then specify <b>n</b> and press Enter to indicate that this is a replicated installation. If you are installing to a drive that is shared across all the nodes (for example, NFS), then specify <b>y</b> and press Enter.
Enter full path to machine list.	Specify the name of the file that you created in the section <a href="#">“List the Machines in the Cluster or Appliance”</a> (for example, <code>/etc/gridhosts</code> ) and press Enter.

- 12** If you selected a replicated installation at the first prompt, you are now prompted to choose the technique for distributing the contents to the appliance nodes:

```
The install can now copy this directory to all the machines
listed in 'filename' using scp, skipping the first entry.
Perform copy? (YES/no)
```

If you want the installation program to perform the replication, specify **yes** and press Enter. If you are distributing the contents of the installation directory by some other technique, specify **no** and press Enter.

**13** Make one of the following choices:

- To enable the SAS High-Performance Analytics environment to read and write SASHDAT using encryption, proceed to [“Configuring the Analytics Environment for SASHDAT Encryption” on page 79](#).
- To configure the analytics environment for a SAS Embedded Process, proceed to [“Configuring for Access to a Data Store with a SAS Embedded Process” on page 113](#).
- To validate your analytics environment, proceed to [“Validating the Analytics Environment Deployment” on page 80](#).

---

## Configuring the Analytics Environment for SASHDAT Encryption

In release 2.94, the SAS High-Performance Analytics environment supports reading and writing files using AES encryption with 256-bit keys. (This feature is very similar to the AES encryption provided by the SAS BASE Engine.)

**Note:** For U.S. export purposes, SAS designates each product based on the encryption algorithms and the product’s functional capability. The ability to encrypt SASHDAT files is available to most commercial and government users inside and outside the U.S. However, some countries (for example, Russia, China, and France) have import restrictions on products that contain encryption, and the U.S. prohibits the export of encryption software to specific embargoed or restricted destinations.

To enable the SAS High-Performance Analytics environment to read and write SASHDAT using encryption, follow these steps:

- 1** The software that is needed for the SAS High-Performance Analytics environment is available from within the SAS Software Depot that was created by the site depot administrator: `depot-installation-location/standalone_installs/SAS_High-Performance_Encryption_Installation/2_94/Linux_for_x64`.
- 2** Copy `TKGrid_SEC_x86_64.sh` to the `/tmp` directory of the root node of the cluster.
- 3** Log on to the machine that will serve as the root node of the cluster or the data appliance with a user account that has the necessary permissions.  
For more information, see [“User Accounts for the SAS High-Performance Analytics Environment” on page 28](#).
- 4** Change directories to the desired installation location, such as `/opt`.
- 5** Run the `TKGrid_SEC_x86_64` shell script in this directory.
- 6** Respond to the prompts from the shell script:

**Table 5.4** Configuration Prompts for the TKGrid\_SEC\_x86\_64 Shell Script

Shared install or replicate to each node? (Y=SHARED/n=replicated)	If you are installing to a local drive on each node, then specify <code>n</code> and press Enter to indicate that this is a replicated installation. If you are installing to a drive that is shared across all the nodes (for example, NFS), then specify <code>Y</code> and press Enter.
--	--

- 7** If you selected a replicated installation at the first prompt, you are now prompted to choose the technique for distributing the contents to the appliance nodes:

```
The install can now copy this directory to all the machines
listed in 'filename' using scp, skipping the first entry.
Perform copy?
(YES/no)
```

Press Enter if you want the installation program to perform the replication. Enter `no` if you are distributing the contents of the installation directory by some other technique.

**Note:** The contents of TKGrid\_SEC must be distributed to every machine in the analytics cluster.

The shell script creates a `lib2` subdirectory and a file named `VERSION2`.

**TIP** If you are using Hadoop as your data provider, make sure that you follow the steps described for your distribution of Hadoop in [“Configuring Existing Hadoop Clusters” on page 60](#).

- 8** To validate your analytics environment, proceed to [“Validating the Analytics Environment Deployment” on page 80](#).

## Validating the Analytics Environment Deployment

### Overview of Validating

You have at least two methods to validate your SAS High-Performance Analytics environment deployment:

- [“Use simsh to Validate” on page 80](#).
- [“Use MPI to Validate” on page 81](#).

### Use simsh to Validate

To validate your SAS High-Performance Analytics environment deployment by issuing a `simsh` command, follow these steps:

- 1** Log on to one of the machines in the analytics cluster.



- 2 Enter the following command:

```
/HPA-environment-installation-directory/bin/simsh hostname
```

This command invokes the `hostname` command on each machine in the cluster. The host name for each machine is printed to the screen.

You should see a list of known hosts similar to the following:

```
myblade006.example.com: myblade006.example.com
myblade007.example.com: myblade007.example.com
myblade004.example.com: myblade004.example.com
myblade005.example.com: myblade005.example.com
```

- 3 Proceed to [Chapter 7, “Configuring the Analytics Environment for a Remote Parallel Connection,”](#) on page 107.

## Use MPI to Validate

To validate your SAS High-Performance Analytics environment deployment by issuing a Message Passing Interface (MPI) command, follow these steps:

- 1 Log on to the root node using the SAS High-Performance Analytics environment installation account.

- 2 Enter the following command:

```
/HPA-environment-installation-directory/TKGrid/mpich2-install/bin/mpirun
-f /etc/gridhosts hostname
```

You should see a list of known hosts similar to the following:

```
myblade006.example.com
myblade007.example.com
myblade004.example.com
myblade005.example.com
```

- 3 Proceed to [Chapter 7, “Configuring the Analytics Environment for a Remote Parallel Connection,”](#) on page 107.

---

# Resource Management for the Analytics Environment

## Resource Settings File

You can set limits on any TKGrid process running across the SAS High-Performance Analytics environment with a resource settings file supplied by SAS. Located in `/opt/TKGrid/`, `resource.settings` is in the format of a shell script. When the analytics environment starts, the environment variables contained in the file are set and last for the duration of the run.

Initially, all of the settings in `resource.settings` are commented. Uncomment the variables and add values that make sense for your site. For more information,

see Appendix 5, “Using CGroups and Memory Limits,” in *SAS LASR Analytic Server: Reference Guide* .

When you are finished editing, copy resource.settings to every machine in the analytics environment:

```
/opt/TKGrid/bin/simcp /opt/TKGrid/resource.settings /opt/TKGrid
```

If YARN is used on the cluster, then you can configure the analytic environment to participate in the resource accounting that YARN performs. For more information, see Appendix 5, “Managing Resources,” in *SAS LASR Analytic Server: Reference Guide*.

resource.settings consists of the following:

```
# if [ "$USER" = "lasradm" ]; then
# Custom settings for any process running under the lasradm account.
#   export TKMPI_ULIMIT="-v 50000000"
#   export TKMPI_MEMSIZE=50000
#   export TKMPI_CGROUP="cgexec -g cpu:75"
# fi

# if [ "$TKMPI_APPNAME" = "lasr" ]; then
# Custom settings for a lasr process running under any account.
#   export TKMPI_ULIMIT="-v 50000000"
#   export TKMPI_MEMSIZE=50000
#   export TKMPI_CGROUP="cgexec -g cpu:75"

# Allow other users to read server and tables, but not add or term.
#   export TKMPI_UMASK=0033

# Allow no access by other users to lasr server.
#   export TKMPI_UMASK=0077

# To exclude from YARN resource manager.
#   unset TKMPI_RESOURCEMANAGER

# Use default nice for LASR
# unset TKMPI_NICE
# fi

# if [ "$TKMPI_APPNAME" = "tklogis" ]; then
# Custom settings for a tklogis process running under any account.
#   export TKMPI_ULIMIT="-v 25000000"
#   export TKMPI_MEMSIZE=25000
#   export TKMPI_CGROUP="cgexec -g cpu:25"
#   export TKMPI_MAXRUNTIME=7200
# fi

# fi

# if [ "$TKMPI_INFO" = "LASRLOAD" ]; then
#   TKMPI_INFO is an environment variable that will be passed from
#   MVA SAS to the grid. It can be used to distinguish a
#   proc lasr create from a proc lasr add, by including
#   this line before the proc lasr add:
#   options set=TKMPI_INFO="LASRLOAD";
#   To exclude from YARN resource manager.
#   unset TKMPI_RESOURCEMANAGER
```

```
# fi
```

## Request Memory with TKMPI\_INFO

When programmers use TKMPI\_INFO in their SAS code, the SAS High-Performance Analytics environment can better decide how much memory to request.

Consider this example: the \$TKMPI\_APPNAME variable is set to `lasr` for both a SAS Analytic LASR Server (PROC LASR CREATE) and for a SAS Analytic LASR Server Proxy used when loading a table (PROC LASR ADD). This makes it impossible to set a YARN memory limit differently for these two cases. Most likely, a SAS Analytic LASR Server would want a large amount of memory and the proxy server would require a smaller amount.

Here is an example of how you might use TKMPI\_INFO in a SAS program to solve the memory issue:

```
options set=TKMPI_INFO="LASRSTART";
proc lasr create port=17761;
performance nodes=2; run;

options set=TKMPI_INFO="LASRLOAD";
proc lasr add data=sashelp.cars port=17761; run
```

In `resource.settings`, you might add an entry similar to the following:

```
if [ "$TKMPI_APPNAME" = "lasr" ]; then
  if [ "$TKMPI_INFO" = "LASRSTART" ];
    export TKMPI_MEMSIZE=60000
  fi
  if [ "$TKMPI_INFO" = "LASRLOAD" ];
    export TKMPI_MEMSIZE=4000
  fi
fi
```

Note that TKMPI\_INFO is not limited to SAS Analytic LASR Server. TKMPI\_INFO can also be used for any other HPA PROC. You could use the variable to pass any kind of information you need to `resource.settings` (for example SMALL, MEDIUM, LARGE classes).



## 6

## Deploying SAS Embedded Process for Hadoop

<i>Infrastructure Deployment Process Overview</i> .....	<b>85</b>
<i>Important Note</i> .....	<b>86</b>
<i>In-Database Deployment Package for Hadoop</i> .....	<b>86</b>
Prerequisites .....	86
Overview of the In-Database Deployment Package for Hadoop .....	88
<i>Hadoop Installation and Configuration</i> .....	<b>89</b>
Hadoop Installation and Configuration Steps .....	89
Upgrading from or Reinstalling a Previous Version .....	89
Moving the SAS Embedded Process and SAS Hadoop	
MapReduce JAR File Install Scripts .....	91
Installing the SAS Embedded Process and SAS Hadoop	
MapReduce JAR Files .....	92
How to Merge Configuration File Properties .....	96
Copying Hadoop JAR Files to the Client Machine .....	97
<i>SASEP-SERVERS.SH Script</i> .....	<b>98</b>
Overview of the SASEP-SERVERS.SH Script .....	98
SASEP-SERVERS.SH Syntax .....	99
Starting the SAS Embedded Process .....	103
Stopping the SAS Embedded Process .....	103
Determining the Status of the SAS Embedded Process .....	104
<i>Hadoop Permissions</i> .....	<b>104</b>
<i>Documentation for Using In-Database Processing in Hadoop</i> .....	<b>105</b>

---

### Infrastructure Deployment Process Overview

Installing and configuring the SAS High-Performance Analytics environment is an optional seventh of eight steps.

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.

4. (Optional) Deploy SAS High-Performance Computing Management Console.
5. (Optional) Deploy co-located Hadoop.
6. Deploy the SAS High-Performance Analytics environment.
- **7. (Optional) Deploy the SAS Embedded Process for Hadoop.**
8. (Optional) Configure the analytics environment for a remote parallel connection.

This chapter describes how to install and configure SAS Embedded Process for Hadoop.

---

## Important Note

This chapter contains information about the SAS Embedded Process for Hadoop. There are other third-party data providers that are supported. For more information, see the *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.

---

## In-Database Deployment Package for Hadoop

### Prerequisites

The following prerequisites are required before you install and configure the in-database deployment package for Hadoop:

- SAS Foundation and the SAS/ACCESS Interface to Hadoop are installed.
- You have working knowledge of the Hadoop vendor distribution that you are using (for example, Cloudera or Hortonworks).

You also need working knowledge of the Hadoop Distributed File System (HDFS), MapReduce 1, MapReduce 2, YARN, Hive, and HiveServer2 services. For more information, see the [Apache website](#) or the vendor's website.

- The SAS Scoring Accelerator for Hadoop requires a specific configuration file. For more information, see [“How to Merge Configuration File Properties” on page 96](#).
- The HDFS, MapReduce, YARN, and Hive services must be running on the Hadoop cluster.
- The SAS Scoring Accelerator for Hadoop requires a specific version of the Hadoop distribution. For more information, see the SAS Foundation system requirements documentation for your operating environment.
- You have root or sudo access. Your user has Write permission to the root of HDFS.
- You know the location of the MapReduce home.

- You know the host name of the Hive server and the NameNode.
  - You understand and can verify your Hadoop user authentication.
  - You understand and can verify your security setup.
- If you are using Kerberos, you need the ability to get a Kerberos ticket.
- You have permission to restart the Hadoop MapReduce service.
  - In order to avoid SSH key mismatches during installation, add the following two options to the SSH `config` file, under the user's home `.ssh` folder. An example of a home `.ssh` folder is `/root/.ssh/`. `nodes` is a list of nodes separated by a space.

```
host nodes
    StrictHostKeyChecking no
    UserKnownHostsFile /dev/null
```

For more details about the SSH `config` file, see the SSH documentation.

- All machines in the cluster are set up to communicate with passwordless SSH. Verify that the nodes can access the node that you chose to be the master node by using SSH.

Traditionally, public key authentication in Secure Shell (SSH) is used to meet the passwordless access requirement. SSH keys can be generated with the following example.

```
[root@raincloud1 .ssh]# ssh-keygen -t rsa
Generating public/private rsa key pair.
Enter file in which to save the key (/root/.ssh/id_rsa):
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /root/.ssh/id_rsa.
Your public key has been saved in /root/.ssh/id_rsa.pub.
The key fingerprint is:
09:f3:d7:15:57:8a:dd:9c:df:e5:e8:1d:e7:ab:67:86 root@raincloud1
```

```
add id_rsa.pub public key from each node to the master node authorized
key file under /root/.ssh/authorized_keys
```

For Secure Mode Hadoop, GSSAPI with Kerberos is used as the passwordless SSH mechanism. GSSAPI with Kerberos not only meets the passwordless SSH requirements, but also supplies Hadoop with the credentials required for users to perform operations in HDFS with SAS LASR Analytic Server and SASHDAT files. Certain options must be set in the SSH daemon and SSH client configuration files. Those options are as follows and assume a default configuration of `sshd`.

To configure passwordless SSH to use Kerberos, follow these steps:

- 1 In the `sshd_config` file, set:

```
GSSAPIAuthentication yes
```

- 2 In the `ssh_config` file, set:

```
Host *.domain.net
GSSAPIAuthentication yes
GSSAPIDelegateCredentials yes
```

where *domain.net* is the domain name used by the machine in the cluster.

**TIP** Although you can specify `host *`, this is not recommended because it allows GSSAPI Authentication with any host name.

## Overview of the In-Database Deployment Package for Hadoop

This section describes how to install and configure the in-database deployment package for Hadoop (SAS Embedded Process).

The in-database deployment package for Hadoop must be installed and configured before you can perform the following tasks:

- Run a scoring model in Hadoop Distributed File System (HDFS) using the SAS Scoring Accelerator for Hadoop.
- Run DATA step scoring programs in Hadoop.
- Run DS2 threaded programs in Hadoop using the SAS In-Database Code Accelerator for Hadoop.
- Read and write data to HDFS in parallel for SAS High-Performance Analytics.

**Note:** For deployments that use SAS High-Performance Deployment of Hadoop for the co-located data provider, and access SASHDAT tables exclusively, SAS/ACCESS and SAS Embedded Process are not needed.

**Note:** If you are installing the SAS High-Performance Analytics environment, you must perform additional steps after you install the SAS Embedded Process. For more information, see *SAS High-Performance Analytics Infrastructure: Installation and Configuration Guide*.

- Transform data in Hadoop and extract transformed data out of Hadoop for analysis in SAS with the SAS Data Loader for Hadoop. For more information, see *SAS Data Loader for Hadoop: User's Guide*.
- Perform data quality operations in Hadoop using the SAS Data Loader for Hadoop. For more information, see *SAS Data Loader for Hadoop: User's Guide*.

The in-database deployment package for Hadoop includes the SAS Embedded Process and two SAS Hadoop MapReduce JAR files. The SAS Embedded Process is a SAS server process that runs within Hadoop to read and write data. The SAS Embedded Process contains macros, run-time libraries, and other software that is installed on your Hadoop system.

The SAS Embedded Process must be installed on all nodes capable of executing either MapReduce 1 or MapReduce 2 tasks. For MapReduce 1, this would be nodes where a TaskTracker is running. For MapReduce 2, this would be nodes where a NodeManager is running. Usually, every DataNode node has a YARN NodeManager or a MapReduce 1 TaskTracker running. By default, the SAS Embedded Process install script (`sasep-servers.sh`) discovers the cluster topology and installs the SAS Embedded Process on all DataNode nodes, including the host node from where you run the script (the Hadoop master NameNode). This occurs even if a DataNode is not present. If you want to limit the list of nodes on which you want the SAS Embedded Process installed, you should run the `sasep-servers.sh` script with the `-host <hosts>` option. The SAS Hadoop MapReduce JAR files must be installed on all nodes of a Hadoop cluster.



---

# Hadoop Installation and Configuration

## Hadoop Installation and Configuration Steps

Before you begin the Hadoop installation and configuration, review [“Prerequisites” on page 86](#).

- 1 If you are upgrading from or reinstalling a previous release, follow the instructions in [“Upgrading from or Reinstalling a Previous Version” on page 89](#) before installing the in-database deployment package.

- 2 Move the SAS Embedded Process and SAS Hadoop MapReduce JAR file install scripts to the Hadoop master node (the NameNode).

For more information, see [“Moving the SAS Embedded Process and SAS Hadoop MapReduce JAR File Install Scripts” on page 91](#).

**Note:** The location where you transfer the install scripts becomes the SAS Embedded Process home and is referred to as *SASEPHome* throughout this chapter.

**Note:** Both the SAS Embedded Process install script and the SAS Hadoop MapReduce JAR file install script must be transferred to the *SASEPHome* directory.

- 3 Install the SAS Embedded Process and the SAS Hadoop MapReduce JAR files.

For more information, see [“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92](#).

- 4 If you want to use the SAS Scoring Accelerator for Hadoop, merge the properties of several configuration files into one configuration file.

For more information, see [“How to Merge Configuration File Properties” on page 96](#).

- 5 Copy the Hadoop core and common Hadoop JAR files to the client machine.

For more information, see [“Copying Hadoop JAR Files to the Client Machine” on page 97](#).

**Note:** If you are installing the SAS High-Performance Analytics environment, you must perform additional steps after you install the SAS Embedded Process. For more information, see *SAS High-Performance Analytics Infrastructure: Installation and Configuration Guide*.

## Upgrading from or Reinstalling a Previous Version

To upgrade or reinstall a previous version, follow these steps.

- 1 If you are upgrading from SAS 9.3, follow these steps. If you are upgrading from SAS 9.4, start with Step 2.
  - a Stop the Hadoop SAS Embedded Process.

```
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.35/bin/sasep-stop.all.sh
```

**SASEPHome** is the master node where you installed the SAS Embedded Process.

- b** Delete the Hadoop SAS Embedded Process from all nodes.

```
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.35/bin/sasep-delete.all.sh
```

- c** Verify that the `sas.hadoop.ep.distribution-name.jar` files have been deleted.

The JAR files are located at **HadoopHome/lib**.

For Cloudera, the JAR files are typically located here:

```
/opt/cloudera/parcels/CDH/lib/hadoop/lib
```

For Hortonworks, the JAR files are typically located here:

```
/usr/lib/hadoop/lib
```

- d** Continue with Step 3.

## 2 If you are upgrading from SAS 9.4, follow these steps.

- a** Stop the Hadoop SAS Embedded Process.

```
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.*/bin/sasep-servers.sh
-stop -hostfile host-list-filename | -host <">host-list<">
```

**SASEPHome** is the master node where you installed the SAS Embedded Process.

For more information, see [“SASEP-SERVERS.SH Script” on page 98](#).

- b** Remove the SAS Embedded Process from all nodes.

```
SASEPHome/SAS/SASTKInDatabaseForServerHadoop/9.*/bin/sasep-servers.sh
-remove -hostfile host-list-filename | -host <">host-list<">
-mrhome dir
```

**Note:** This step ensures that all old SAS Hadoop MapReduce JAR files are removed.

For more information, see [“SASEP-SERVERS.SH Script” on page 98](#).

- c** Verify that the `sas.hadoop.ep.apache*.jar` files have been deleted.

The JAR files are located at **HadoopHome/lib**.

For Cloudera, the JAR files are typically located here:

```
/opt/cloudera/parcels/CDH/lib/hadoop/lib
```

For Hortonworks, the JAR files are typically located here:

```
/usr/lib/hadoop/lib
```

- d** Manually remove the SAS Embedded Process directories and files on the node from which you ran the script.

The `sasep-servers.sh -remove` script removes the file everywhere except on the node from which you ran the script. The `sasep-servers.sh -remove` script displays instructions that are similar to the following example.

```
localhost WARN: Apparently, you are trying to uninstall SAS Embedded Process
for Hadoop from the local node.
```

```

The binary files located at
  local_node/SAS/SASTKInDatabaseServerForHadoop/local_node/
  SAS/SASACCESSToHadoopMapReduceJARFiles will not be removed.
localhost WARN: The init script will be removed from /etc/init.d and the
  SAS Map Reduce JAR files will be removed from /usr/lib/hadoop-mapreduce/lib.
localhost WARN: The binary files located at local_node/SAS
  should be removed manually.

```

### 3 Continue the installation process.

For more information, see [“Moving the SAS Embedded Process and SAS Hadoop MapReduce JAR File Install Scripts” on page 91.](#)

## Moving the SAS Embedded Process and SAS Hadoop MapReduce JAR File Install Scripts

### Creating the SAS Embedded Process Directory

Before you can install the SAS Embedded Process and the SAS Hadoop MapReduce JAR files, you must move the SAS Embedded Process and SAS Hadoop MapReduce JAR file install scripts to a directory on the Hadoop master node (the NameNode).

Create a new directory that is not part of an existing directory structure, such as */sasep*.

This path will be created on each node in the Hadoop cluster during the SAS Embedded Process installation. Do not use existing system directories such as */opt* or */usr*. This new directory becomes the SAS Embedded Process home and is referred to as *SASEPHome* throughout this chapter.

### Moving the SAS Embedded Process Install Script

The SAS Embedded Process install script is contained in a self-extracting archive file named *tkindbsrv-9.42-n\_lax.sh*. *n* is a number that indicates the latest version of the file. If this is the initial installation, *n* has a value of 1. Each time you reinstall or upgrade, *n* is incremented by 1. The self-extracting archive file is located in the *SAS-installation-directory/SASTKInDatabaseServer/9.4/HadooponLinuxx64/* directory.

Using a method of your choice, transfer the SAS Embedded Process install script to your Hadoop master node.

This example uses secure copy, and *SASEPHome* is the location where you want to install the SAS Embedded Process.

```
scp tkindbsrv-9.42-n_lax.sh username@hadoop:/SASEPHome
```

**Note:** Both the SAS Embedded Process install script and the SAS Hadoop MapReduce JAR file install script must be transferred to the *SASEPHome* directory.

### Moving the SAS Hadoop MapReduce JAR File Install Script

The SAS Hadoop MapReduce JAR file install script is contained in a self-extracting archive file named *hadoopmrjars-9.42-n\_lax.sh*. *n* is a number that indicates the latest version of the file. If this is the initial installation, *n* has a value of 1. Each time you reinstall or upgrade, *n* is incremented by 1. The self-

extracting archive file is located in the *SAS-installation-directory/SASACCESStoHadoopMapReduceJARFiles/9.41* directory.

Using a method of your choice, transfer the SAS Hadoop MapReduce JAR file install script to your Hadoop master node.

This example uses Secure Copy, and *SASEPHome* is the location where you want to install the SAS Hadoop MapReduce JAR files.

```
scp hadoopmrjars-9.42-n_lax.sh username@hadoop:/SASEPHome
```

**Note:** Both the SAS Embedded Process install script and the SAS Hadoop MapReduce JAR file install script must be transferred to the *SASEPHome* directory.

## Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files

To install the SAS Embedded Process, follow these steps.

**Note:** Permissions are needed to install the SAS Embedded Process and SAS Hadoop MapReduce JAR files. For more information, see [“Hadoop Permissions” on page 104](#).

- 1 Log on to the server using SSH as root with sudo access.

```
ssh username@serverhostname
sudo su - root
```

- 2 Move to your Hadoop master node where you want the SAS Embedded Process installed.

```
cd /SASEPHome
```

*SASEPHome* is the same location to which you copied the self-extracting archive file. For more information, see [“Moving the SAS Embedded Process Install Script” on page 91](#).

**Note:** Before continuing with the next step, ensure that each self-extracting archive file has Execute permission.

- 3 Use the following script to unpack the *tkindbsrv-9.42-n\_lax.sh* file.

```
./tkindbsrv-9.42-n_lax.sh
```

*n* is a number that indicates the latest version of the file. If this is the initial installation, *n* has a value of 1. Each time you reinstall or upgrade, *n* is incremented by 1.

**Note:** If you unpack in the wrong directory, you can move it after the unpack.

After this script is run and the files are unpacked, the script creates the following directory structure where *SASEPHome* is the master node from Step 1.

```
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/misc
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/sasexe
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/utilities
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/build
```

The content of the

***SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin*** directory should look similar to this.

```
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sas.ep4hadoop.template
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sasep-servers.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sasep-common.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sasep-server-start.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sasep-server-status.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sasep-server-stop.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/InstallTKIndbsrv.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/MANIFEST.MF
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/qkbpush.sh
SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/sas.tools.qkb.hadoop.jar
```

#### 4 Use this command to unpack the SAS Hadoop MapReduce JAR files.

```
./hadoopmrjars-9.42-n_lax.sh
```

After the script is run, the script creates the following directory and unpacks these files to that directory.

```
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/ep-config.xml
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache023.jar
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache023.nls.jar
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache121.jar
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache121.nls.jar
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache205.jar
SASEPHome/SAS/SASACCESSToHadoopMapReduceJARFiles/9.42/lib/
sas.hadoop.ep.apache205.nls.jar
```

#### 5 Use the `sasep-servers.sh -add` script to deploy the SAS Embedded Process installation across all nodes. The SAS Embedded Process is installed as a Linux service.

**TIP** There are many options available when installing the SAS Embedded Process. We recommend that you review the script syntax before running it. For more information, see [“SASEP-SERVERS.SH Script” on page 98](#).

**Note:** If you are running on a cluster with Kerberos, complete both steps a and b. If you are not running with Kerberos, complete only step b.

##### a If you are running on a cluster with Kerberos, you must kinit the HDFS user.

```
sudo su - root
su - hdfs | hdfs-userid
kinit -kt location of keytab file
      user for which you are requesting a ticket
exit
```

Here is an example:

```
sudo su - root
su - hdfs
```

```
kinit -kt hdfs.keytab hdfs
exit
```

**Note:** The default HDFS user is `hdfs`. You can specify a different user ID with the `-hdfsuser` argument when you run the `sasep-servers.sh -add script`.

**Note:** If you are running on a cluster with Kerberos, a keytab is required for the `-hdfsuser` running the `sasep-servers.sh -add script`.

**Note:** You can run `klist` while you are running as the `-hdfsuser` user to check the status of your Kerberos ticket on the server. Here is an example:

```
klist
Ticket cache: FILE/tmp/krb5cc_493
Default principal: hdfs@HOST.COMPANY.COM

Valid starting    Expires          Service principal
06/20/14 09:51:26 06/27/14 09:51:26 krbtgt/HOST.COMPANY.COM@HOST.COMPANY.COM
        renew until 06/22/14 09:51:26
```

- b** Run the `sasep-servers.sh` script. Review all of the information in this step before running the script.

```
cd $HADOOP_HOME/SAS/SASTKInDatabaseServerForHadoop/9.42/bin
./sasep-servers.sh -add
```

During the install process, the script asks whether you want to start the SAS Embedded Process. If you choose `y` or `Y`, the SAS Embedded Process is started on all nodes after the install is complete. If you choose `n` or `N`, you can start the SAS Embedded Process later by running `./sasep-servers.sh -start`.

**Note:** When you run the `sasep-servers.sh -add script`, a user and group named `sasep` is created. You can specify a different user and group name with the `-epuser` and `-epgroup` arguments when you run the `sasep-servers.sh -add script`.

**Note:** The `sasep-servers.sh` script can be run from any location. You can also add its location to the `PATH` environment variable.

**TIP** Although you can install the SAS Embedded Process in multiple locations, the best practice is to install only one instance. Only one version of the SASEP JAR files will be installed in your `HadoopHome/lib` directory.

**Note:** The SAS Embedded Process must be installed on all nodes capable of executing either MapReduce 1 or MapReduce 2 tasks. For MapReduce 1, this would be nodes where a TaskTracker is running. For MapReduce 2, this would be nodes where a NodeManager is running. Usually, every DataNode node has a YARN NodeManager or a MapReduce 1 TaskTracker running. By default, the SAS Embedded Process install script (`sasep-servers.sh`) discovers the cluster topology and installs the SAS Embedded Process on all DataNode nodes, including the host node from where you run the script (the Hadoop master NameNode). This occurs even if a DataNode is not present. If you want to limit the list of nodes on which you want the SAS Embedded Process installed, you should run the `sasep-servers.sh` script with the `-host <hosts>` option.

**Note:** If you install the SAS Embedded Process on a large cluster, the SSHD daemon might reach the maximum number of concurrent connections. The `ssh_exchange_identification: Connection closed by remote host` SSHD error might occur. To work around the problem, edit the `/etc/ssh/sshd_config` file, change the `MaxStartups` option to the number that accommodates your cluster, and save the file. Then, reload the SSHD daemon by running the `/etc/init.d/sshd reload` command.

- 6 Verify that the SAS Embedded Process is installed and running. Change directories and then run the `sasep-servers.sh` script with the `-status` option.

```
cd $SASEPHOME/SAS/SASTKInDatabaseServerForHadoop/9.42/bin
./sasep-servers.sh -status
```

This command returns the status of the SAS Embedded Process running on each node of the Hadoop cluster. Verify that the SAS Embedded Process home directory is correct on all the nodes.

**Note:** The `sasep-servers.sh -status` script will not run successfully if the SAS Embedded Process is not installed.

- 7 Verify that the `sas.hadoop.ep.apache*.jar` files are now in place on all nodes.

The JAR files are located at `HadoopHome/lib`.

For Cloudera, the JAR files are typically located here:

```
/opt/cloudera/parcels/CDH/lib/hadoop/lib
```

For Hortonworks, the JAR files are typically located here:

```
/usr/lib/hadoop/lib
```

- 8 Restart the Hadoop YARN or MapReduce service.

This enables the cluster to reload the SAS Hadoop JAR files (`sas.hadoop.ep.*.jar`).

**Note:** It is preferable to restart the service by using Cloudera Manager or Hortonworks Ambari.

- 9 Verify that an `init.d` service with a `sas.ep4hadoop` file was created in the following directory.

```
/etc/init.d
```

View the `sas.ep4hadoop` file and verify that the SAS Embedded Process home directory is correct.

The `init.d` service is configured to start at level 3 and level 5.

**Note:** The SAS Embedded Process needs to run on all nodes in the Hadoop cluster.

- 10 Verify that the configuration file, `ep-config.xml`, was written to the HDFS file system.

```
hadoop fs -ls /sas/ep/config
```

**Note:** If you are running on a cluster with Kerberos, you need a Kerberos ticket. If not, you can use the WebHDFS browser.

**Note:** The `/sas/ep/config` directory is created automatically when you run the install script.

## How to Merge Configuration File Properties

### Requirements for Configuration File Properties

The SAS Scoring Accelerator for Hadoop requires that some specific configuration files be merged into a single configuration file. This configuration file is used by the %INDHD\_RUN\_MODEL macro. For more information about the %INDHD\_RUN\_MODEL macro, see the *SAS In-Database Products: User's Guide*.

**Note:** In addition to the SAS Scoring Accelerator for Hadoop, a merged configuration file is also required for using PROC HADOOP and the FILENAME Statement Hadoop Access Method with Base SAS.

**Note:** A merged configuration file is not required for the following products:

- SAS/ACCESS Interface to Hadoop
- Scalable Performance Data Engine (SPD Engine)
- High-Performance Analytics
- SAS Code Accelerator for Hadoop

In the August 2014 release, a new environment variable, SAS\_HADOOP\_CONFIG\_PATH, was created to replace the use of a merged configuration file for the preceding list of products. The configuration files are copied to a location on the client machine, and the SAS\_HADOOP\_CONFIG\_PATH variable is set to that location.

**Note:** The configuration file properties that must be merged depend on which version of Hadoop you are using. The following topics are for Cloudera CDH4.x and 5.x and Hortonworks 1.3.2 and 2.x. If you are using a different version of Cloudera or Hortonworks or another vendor's Hadoop distribution, you must use a comparable version of these configuration files. Otherwise, the installation of the SAS Embedded Process will fail.

### How to Merge Configuration File Properties

To merge configuration file properties, follow these steps:

- 1 Retrieve the Hadoop configuration files from this location on your cluster.

/etc/hadoop/conf

- 2 Using a method of your choice, concatenate the required configuration files into one configuration file.

- a If you are using MapReduce 1, merge the properties from the Hadoop core (core-site.xml), Hadoop HDFS (hdfs-site.xml), and MapReduce (mapred-site.xml) configuration files into one single configuration file.
- b If you are using MapReduce 2 or YARN, merge the properties from the Hadoop core (core-site.xml), Hadoop HDFS (hdfs-site.xml), MapReduce (mapred-site.xml), and YARN (yarn-site.xml) configuration files into one single configuration file.

**Note:** The merged configuration file must have one beginning <configuration> tag and one ending </configuration> tag. Only properties



should exist between the `<configuration>...</configuration>` tags. Here is a Cloudera example:

```
<?xml version="1.0" encoding="UTF-8"?>

<configuration>
  <property>
    <name>hive.metastore.local</name>
    <value>false</value>
  </property>

  <!-- lines omitted for sake of brevity -->

  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value></value>
  </property>
</configuration>
```

**Note:** For information about setting a configuration property to avoid out of memory exceptions, see [“Configuring the Number of Proactive Reads in the SAS Embedded Process”](#) on page 97.

- 3 Save the configuration file to a location of your choosing.

### Configuring the Number of Proactive Reads in the SAS Embedded Process

The SAS Embedded Process for Hadoop implements a mechanism that proactively reads and caches records from the input files while the DS2 program is processing the current block of records. By default, the number of proactive reads is set to 100. On files with very large records, the default value of 100 might cause memory exhaustion in the Java Virtual Machine in which the SAS Embedded Process MapReduce task is running.

**TIP** Java Virtual Machine out of memory exceptions might be avoided if the number of proactive reads is reduced.

To avoid out of memory exceptions, we recommend setting the following property in the merged configuration file:

```
<property>
  <name>sas.ep.superreader.proactive.reader.capacity</name>
  <value>10</value>
</property>
```

## Copying Hadoop JAR Files to the Client Machine

For SAS components that interface with Hadoop, a specific set of common and core Hadoop JAR files must be in one location on the client machine. Examples of those components are the SAS Scoring Accelerator and SAS High-Performance Analytics.

When you run the `sasep-servers.sh -add` script to install the SAS Embedded Process, the script detects the Hadoop distribution and creates a `HADOOP_JARS.zip` file in the `$SASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/` directory. This file

contains the common and core Hadoop JAR files that are required for the SAS Embedded Process. For more information, see [“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92.](#)

To get the Hadoop JAR files on your client machine, follow these steps:

- 1 Move the HADOOP\_JARS.zip file to a directory on your client machine and unzip the file.

```
unzip HADOOP_JARS.zip
```

- 2 Set the SAS\_HADOOP\_JAR\_PATH environment variable to point to the directory that contains the core and common Hadoop JAR files.

**Note:** You can run the `sasep-servers.sh -getjars` script at any time to create a new ZIP file and refresh the JAR file list.

**Note:** The MapReduce 1 and MapReduce 2 JAR files cannot be on the same Java classpath.

**Note:** The JAR files in the SAS\_HADOOP\_JAR\_PATH directory must match the Hadoop server to which SAS is connected. If multiple Hadoop servers are running different Hadoop versions, then create and populate separate directories with version-specific Hadoop JAR files for each Hadoop version. Then dynamically set SAS\_HADOOP\_JAR\_PATH, based on the target Hadoop server to which each SAS job or SAS session is connected. One way to dynamically set SAS\_HADOOP\_JAR\_PATH is to create a wrapper script associated with each Hadoop version. Then invoke SAS via a wrapper script that sets SAS\_HADOOP\_JAR\_PATH appropriately to pick up the JAR files that match the target Hadoop server. Upgrading your Hadoop server version might involve multiple active Hadoop versions. The same multi-version instructions apply.

---

## SASEP-SERVERS.SH Script

### Overview of the SASEP-SERVERS.SH Script

The `sasep-servers.sh` script enables you to perform the following actions.

- Install or uninstall the SAS Embedded Process and SAS Hadoop MapReduce JAR files on a single node or a group of nodes.
- Start or stop the SAS Embedded Process on a single node or on a group of nodes.
- Determine the status of the SAS Embedded Process on a single node or on a group of nodes.
- Write the installation output to a log file.
- Pass options to the SAS Embedded Process.
- Create a HADOOP\_JARS.zip file in the local folder. This ZIP file contains all required client JAR files.

**Note:** The `sasep-servers.sh` script can be run from any folder on any node in the cluster. You can also add its location to the PATH environment variable.

**Note:** You must have sudo access to run the `sasep-servers.sh` script.

## SASEP-SERVERS.SH Syntax

### sasep-servers.sh

```
-add | -remove | -start | -stop | -status | -restart
<-mrhome path-to-mr-home>
<-hdfsuser user-id>
<-epuser>epuser-id
<-epgroup>epgroup-id
<-hostfile host-list-filename | -host <">host-list<">>
<-epscript path-to-ep-install-script>
<-mrscript path-to-mr-jar-file-script>
<-options "option-list">
<-log filename>
<-version apache-version-number>
<-getjars>
```

### Arguments

#### -add

installs the SAS Embedded Process.

**Note** The -add argument also starts the SAS Embedded Process (same function as -start argument). You are prompted and can choose whether to start the SAS Embedded Process.

**Tip** You can specify the hosts on which you want to install the SAS Embedded Process by using the -hostfile or -host option. The -hostfile or -host options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

#### -remove

removes the SAS Embedded Process.

**CAUTION!** If you are using SAS Data Loader's Cleanse Data in Hadoop directives, you should remove the QKB from the Hadoop nodes before removing the SAS Embedded Process. The QKB is removed by running the QKBPUH script. For more information, see the *SAS Data Loader for Hadoop: Administrator's Guide*.

**Tip** You can specify the hosts for which you want to remove the SAS Embedded Process by using the -hostfile or -host option. The -hostfile or -host options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

#### -start

starts the SAS Embedded Process.

**Tip** You can specify the hosts on which you want to start the SAS Embedded Process by using the -hostfile or -host option. The -hostfile or -host options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

**-stop**

stops the SAS Embedded Process.

**Tip** You can specify the hosts on which you want to stop the SAS Embedded Process by using the `-hostfile` or `-host` option. The `-hostfile` or `-host` options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

**-status**

provides the status of the SAS Embedded Process on all hosts or the hosts that you specify with either the `-hostfile` or `-host` option.

**Tips** The status also shows the version and path information for the SAS Embedded Process.

You can specify the hosts for which you want the status of the SAS Embedded Process by using the `-hostfile` or `-host` option. The `-hostfile` or `-host` options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

**-restart**

restarts the SAS Embedded Process.

**Tip** You can specify the hosts on which you want to restart the SAS Embedded Process by using the `-hostfile` or `-host` option. The `-hostfile` or `-host` options are mutually exclusive.

**See** [-hostfile and -host option on page 100](#)

**-mrhome *path-to-mr-home***

specifies the path to the MapReduce home.

**-hdfsuser *user-id***

specifies the user ID that has Write access to HDFS root directory.

**Default** hdfs

**Note** The user ID is used to copy the SAS Embedded Process configuration files to HDFS.

**-epuser *epuser-name***

specifies the name for the SAS Embedded Process user.

**Default** sasep

**-epgroup *epgroup-name***

specifies the name for the SAS Embedded Process group.

**Default** sasep

**-hostfile *host-list-filename***

specifies the full path of a file that contains the list of hosts where the SAS Embedded Process is installed, removed, started, stopped, or status is provided.

**Default** If you do not specify `-hostfile`, the `sasep-servers.sh` script will discover the cluster topology and use the retrieved list of data nodes.

**Tip** You can also assign a host list filename to a UNIX variable, `sas_ephosts_file`.  
`export sasep_hosts=/etc/hadoop/conf/slaves`

**See** [“-hdfsuser user-id” on page 100](#)

**Example** `-hostfile /etc/hadoop/conf/slaves`

### **-host <">host-list<">**

specifies the target host or host list where the SAS Embedded Process is installed, removed, started, stopped, or status is provided.

**Default** If you do not specify `-host`, the `sasep-servers.sh` script will discover the cluster topology and use the retrieved list of data nodes.

**Requirement** If you specify more than one host, the hosts must be enclosed in double quotation marks and separated by spaces.

**Tip** You can also assign a list of hosts to a UNIX variable, `sas_ephosts`.  
`export sasep_hosts="server1 server2 server3"`

**See** [“-hdfsuser user-id” on page 100](#)

**Example** `-host "server1 server2 server3"`  
`-host bluesvr`

### **-epscript path-to-ep-install-script**

copies and unpacks the SAS Embedded Process install script file to the host.

**Restriction** Use this option only with the `-add` option.

**Requirement** You must specify either the full or relative path of the SAS Embedded Process install script, `tkindbsrv-9.42-n_lax.sh` file.

**Example** `-epscript /home/hadoop/image/current/tkindbsrv-9.42-1_lax.sh`

### **-mrscript path-to-mr-jar-file-script**

copies and unpacks the SAS Hadoop MapReduce JAR files install script on the hosts.

**Restriction** Use this option only with the `-add` option.

**Requirement** You must specify either the full or relative path of the SAS Hadoop MapReduce JAR file install script, `hadoopmrjars-9.42-n_lax.sh` file.

**Example** `-mrscript /home/hadoop/image/current/hadoopmrjars-9.42-1_lax.sh`

### **-options "option-list"**

specifies options that are passed directly to the SAS Embedded Process. The following options can be used.

**-trace *trace-level***

specifies what type of trace information is created.

- 0 no trace log
- 1 fatal error
- 2 error with information or data value
- 3 warning
- 4 note
- 5 information as an SQL statement
- 6 critical and command trace
- 7 detail trace, lock
- 8 enter and exit of procedures
- 9 tedious trace for data types and values
- 10 trace all information

**Default** 02

**Note** The trace log messages are stored in the MapReduce job log.

**-port *port-number***

specifies the TCP port number where the SAS Embedded Process accepts connections.

**Default** 9261

**Requirement** The options in the list must be separated by spaces, and the list must be enclosed in double quotation marks.

**-log *filename***

writes the installation output to the specified filename.

**-version *apache-version-number***

specifies the Hadoop version of the JAR file that you want to install on the cluster. The *apache-version-number* can be one of the following values.

**0.23**

installs the SAS Hadoop MapReduce JAR files that are built from Apache Hadoop 0.23 (sas.hadoop.ep.apache023.jar and sas.hadoop.ep.apache023.nls.jar).

**1.2**

installs the SAS Hadoop MapReduce JAR files that are built from Apache Hadoop 1.2.1 (sas.hadoop.ep.apache121.jar and sas.hadoop.ep.apache121.nls.jar).

**2.0**

installs the SAS Hadoop MapReduce JAR files that are built from Apache Hadoop 0.2.3 (sas.hadoop.ep.apache023.jar and sas.hadoop.ep.apache023.nls.jar).

**2.1**

installs the SAS Hadoop MapReduce JAR files that are built from Apache Hadoop 2.0.5 (sas.hadoop.ep.apache205.jar and sas.hadoop.ep.apache205.nls.jar).

Default	If you do not specify the <code>-version</code> option, the <code>sasep.servers.sh</code> script will detect the version of Hadoop that is in use and install the JAR files associated with that version. For more information, see <a href="#">“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92</a> .
Interaction	The <code>-version</code> option overrides the version that is automatically detected by the <code>sasep.servers.sh</code> script.

**-getjars**

creates a `HADOOP_JARS.zip` file in the local folder. This ZIP file contains all required client JAR files.

You need to move this ZIP file to your client machine and unpack it. If you want to replace the existing JAR files, move it to the same directory where you previously unpacked the existing JAR files.

See For more information, see [“Copying Hadoop JAR Files to the Client Machine” on page 97](#).

## Starting the SAS Embedded Process

There are three ways to manually start the SAS Embedded Process.

**Note:** Root authority is required to run the `sasep-servers.sh` script.

- Run the `sasep-servers.sh` script with the `-start` option on the master node.

This starts the SAS Embedded Process on all nodes. For more information about running the `sasep-servers.sh` script, see [“SASEP-SERVERS.SH Syntax” on page 99](#).

- Run `sasep-server-start.sh` on a node.

This starts the SAS Embedded Process on the local node only. The `sasep-server-start.sh` script is located in the `$ASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/` directory. For more information, see [“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92](#).

- Run the UNIX `service` command on a node.

This starts the SAS Embedded Process on the local node only. The `service` command calls the init script that is located in the `/etc/init.d` directory. A symbolic link to the init script is created in the `/etc/rc3.d` and `/etc/rc5.d` directories, where 3 and 5 are the run level at which you want the script to be executed.

Because the SAS Embedded Process init script is registered as a service, the SAS Embedded Process is started automatically when the node is rebooted.

## Stopping the SAS Embedded Process

The SAS Embedded Process continues to run until it is manually stopped. The ability to control the SAS Embedded Process on individual nodes could be useful when performing maintenance on an individual node.

There are three ways to stop the SAS Embedded Process.

**Note:** Root authority is required to run the `sasep-servers.sh` script.

- Run the `sasep-servers.sh` script with the `-stop` option from the master node.

This stops the SAS Embedded Process on all nodes. For more information about running the `sasep-servers.sh` script, see [“SASEP-SERVERS.SH Syntax” on page 99](#).

- Run `sasep-server-stop.sh` on a node.

This stops the SAS Embedded Process on the local node only. The `sasep-server-stop.sh` script is located in the `$ASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/` directory. For more information, see [“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92](#).

- Run the UNIX `service` command on a node.

This stops the SAS Embedded Process on the local node only.

## Determining the Status of the SAS Embedded Process

You can display the status of the SAS Embedded Process on one node or all nodes. There are three ways to display the status of the SAS Embedded Process.

**Note:** Root authority is required to run the `sasep-servers.sh` script.

- Run the `sasep-servers.sh` script with the `-status` option from the master node.

This displays the status of the SAS Embedded Process on all nodes. For more information about running the `sasep-servers.sh` script, see [“SASEP-SERVERS.SH Syntax” on page 99](#).

- Run `sasep-server-status.sh` from a node.

This displays the status of the SAS Embedded Process on the local node only. The `sasep-server-status.sh` script is located in the `$ASEPHome/SAS/SASTKInDatabaseServerForHadoop/9.42/bin/` directory. For more information, see [“Installing the SAS Embedded Process and SAS Hadoop MapReduce JAR Files” on page 92](#).

- Run the UNIX `service` command on a node.

This displays the status of the SAS Embedded Process on the local node only.

---

## Hadoop Permissions

The person who installs the SAS Embedded Process must have `sudo` access.



---

## Documentation for Using In-Database Processing in Hadoop

For information about using in-database processing in Hadoop, see the following publications:

- *SAS In-Database Products: User's Guide*
- High-performance procedures in various SAS publications
- *SAS Data Integration Studio: User's Guide*
- SAS/ACCESS Interface to Hadoop and PROC HDMD in *SAS/ACCESS for Relational Databases: Reference*
- *SAS High-Performance Analytics Infrastructure: Installation and Configuration Guide*
- *SAS Intelligence Platform: Data Administration Guide*
- PROC HADOOP in *Base SAS Procedures Guide*
- FILENAME Statement, Hadoop Access Method in *SAS Statements: Reference*
- *SAS Data Loader for Hadoop: User's Guide*



## 7

## Configuring the Analytics Environment for a Remote Parallel Connection

<b>Infrastructure Deployment Process Overview</b>	<b>107</b>
<b>Overview of Configuring the Analytics Environment for a Remote Parallel Connection</b>	<b>108</b>
<b>Preparing for a Remote Parallel Connection</b>	<b>109</b>
Overview of Preparing for a Remote Parallel Connection	109
Prepare for Hadoop	110
Prepare for a Greenplum Data Computing Appliance	111
Prepare for a HANA Cluster	111
Prepare for an Oracle Exadata Appliance	112
Prepare for a Teradata Managed Server Cabinet	112
<b>Configuring for Access to a Data Store with a SAS Embedded Process</b>	<b>113</b>
Overview of Configuring for Access to a Data Store with a SAS Embedded Process	113
How the Configuration Script Works	113
Configure for Access to a Data Store with a SAS Embedded Process	114

---

## Infrastructure Deployment Process Overview

Configuring your data storage is the last of eight steps for deploying the SAS High-Performance Analytics infrastructure.

1. Create a SAS Software Depot.
2. Check for documentation updates.
3. Prepare your analytics cluster.
4. (Optional) Deploy SAS High-Performance Computing Management Console.
5. (Optional) Deploy co-located Hadoop.
6. Deploy the SAS High-Performance Analytics environment.
7. (Optional) Deploy the SAS Embedded Process for Hadoop.

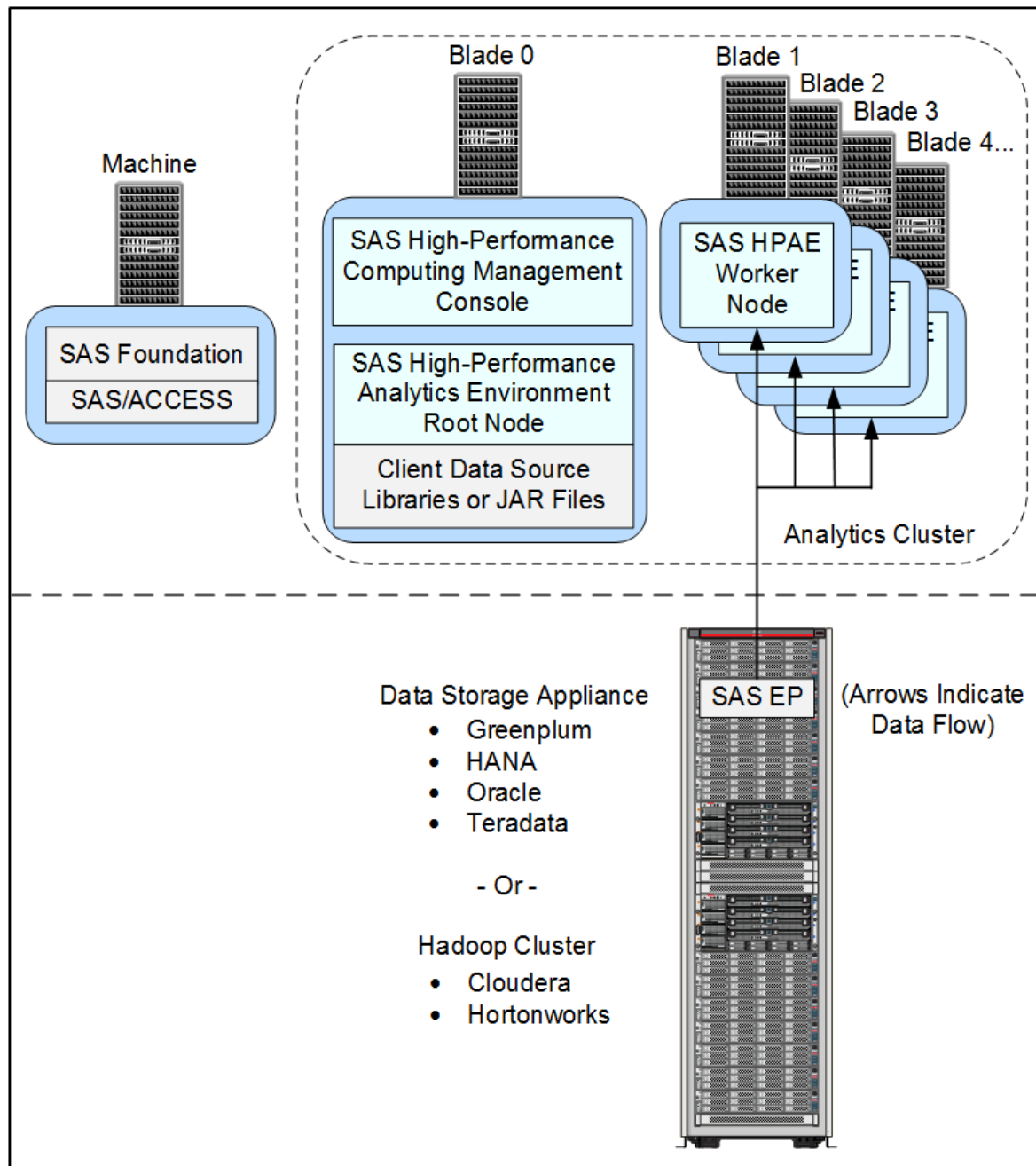
- **8. (Optional) Configure the analytics environment for a remote parallel connection.**

---

## **Overview of Configuring the Analytics Environment for a Remote Parallel Connection**

The SAS/ACCESS Interface and SAS Embedded Process provide a high-speed parallel connection that delivers data from the co-located SAS data source to the SAS-High Performance Analytics environment on the analytic cluster. After you have installed SAS/ACCESS and its embedded process, you configure the analytics environment for the particular access interface that you will use with a shell script, TKGrid\_REP.

Figure 7.1 Analytics Cluster Remote from Your Data Store (Parallel Connection)



## Preparing for a Remote Parallel Connection

### Overview of Preparing for a Remote Parallel Connection

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your data store,

you must locate particular JAR files and gather particular information about your data provider.

From the following list, choose the topic for your respective remote data provider:

- 1 “Prepare for Hadoop” on page 110 .
- 2 “Prepare for a Greenplum Data Computing Appliance” on page 111.
- 3 “Prepare for a HANA Cluster” on page 111.
- 4 “Prepare for an Oracle Exadata Appliance” on page 112.
- 5 “Prepare for a Teradata Managed Server Cabinet” on page 112.

## Prepare for Hadoop

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your Hadoop data store, there are certain requirements that must be met.

**Note:** In the second maintenance release of SAS 9.4, the SAS Embedded Process supports the Cloudera and Hortonworks distributions of Hadoop. For more detailed information, see the SAS Foundation system requirements documentation for your operating environment, available at <http://support.sas.com/resources/sysreq/index.html>.

- 1 Record the path to the Hadoop JAR files required by SAS in the table that follows:

**Table 7.1** Record the Location of the Hadoop JAR Files Required by SAS

Example	Actual Path of the Required Hadoop JAR Files on Your System
/opt/hadoop_jars (common and core JAR files)	
/opt/hadoop_jars/MR1 (Map Reduce JAR files)	
/opt/hadoop_jars/MR2 (Map Reduce JAR files)	

**Note:** The location of the common and core JAR files listed in [Table 7.1](#) should also be the same location that you unzip HADOOP\_JARS.zip in “Copying Hadoop JAR Files to the Client Machine” on page 97.

- 2 Record the location (JAVA\_HOME) of the 64-bit Java Runtime Engine (JRE) on your Hadoop cluster in the table that follows:

**Table 7.2** Record the Location of the JRE

Example	Actual Path of the JRE on Your System
<code>/opt/java/jre1.7.0_07</code>	

## Prepare for a Greenplum Data Computing Appliance

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your Greenplum Data Computing Appliance, there are certain requirements that must be met.

- 1 Install the Greenplum client on the Greenplum Master Server (blade 0) in your analytics cluster.

For more information, refer to your Greenplum documentation.

- 2 Record the path to the Greenplum client in the table that follows:

**Table 7.3** Record the Location of the Greenplum Client

Example	Actual Path of the Greenplum Client on Your System
<code>/usr/local/greenplum-db</code>	

## Prepare for a HANA Cluster

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your HANA cluster, there are certain requirements that must be met.

- 1 Install the HANA client on blade 0 in your analytics cluster.

For more information, refer to your HANA documentation.

- 2 Record the path to the HANA client in the table that follows:

**Table 7.4** Record the Location of the HANA Client

Example	Actual Path of the HANA Client on Your System
<code>/usr/local/lib/hdbclient</code>	

## Prepare for an Oracle Exadata Appliance

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your Oracle Exadata appliance, there are certain requirements that must be met.

- 1 Install the Oracle client on blade 0 in your analytics cluster.  
For more information, refer to your Oracle documentation.
- 2 Record the path to the Oracle client in the table that follows. (This should be the absolute path to libclntsh.so):

**Table 7.5** Record the Location of the Oracle Client

Example	Actual Path of the Oracle Client on Your System
	/usr/local/ora11gr2/product/11.2.0/client_1/lib

- 3 Record the value of the Oracle TNS\_ADMIN environment variable in the table that follows. (Typically, this is the directory that contains the tnsnames.ora file):

**Table 7.6** Record the Value of the Oracle TNS\_ADMIN Environment Variable

Example	Oracle TNS_ADMIN Environment Variable Value on Your System
	/my_server/oracle

## Prepare for a Teradata Managed Server Cabinet

Before you can configure the SAS High-Performance Analytics environment to use the SAS Embedded Process for a parallel connection with your Teradata Managed Server Cabinet, there are certain requirements that must be met.

- 1 Install the Teradata client on blade 0 in your analytics cluster.  
For more information, refer to your Teradata documentation.
- 2 Record the path to the Teradata client in the table that follows. (This should be the absolute path to the directory that contains the odbc\_64 subdirectory):

**Table 7.7** Record the Location of the Teradata Client

Example	Actual Location of the Teradata Client on Your System
	/opt/teradata/client/13.10



---

## Configuring for Access to a Data Store with a SAS Embedded Process

### Overview of Configuring for Access to a Data Store with a SAS Embedded Process

The process involved for configuring the SAS High-Performance Analytics environment with a SAS Embedded Process consists of the following steps:

- 1 Prepare for the data provider that the analytics environment will query.

For more information, see [“Preparing for a Remote Parallel Connection” on page 109](#).

**Note:** Other third-party data providers besides Hadoop are supported. For more information, see the *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.

- 2 Review the considerations for configuring the analytics environment for use with a remote data store.

For more information, see [“How the Configuration Script Works” on page 113](#).

- 3 Configure the analytics environment for a remote data store.

For more information, see [“Configure for Access to a Data Store with a SAS Embedded Process” on page 114](#).

### How the Configuration Script Works

You configure the SAS High-Performance Analytics environment with a SAS Embedded Process using a shell script. The script enables you to configure the environment for the various third-party data stores supported by the SAS Embedded Process.

The Analytics environment is designed on the principle, install once, configure many. For example, suppose that your site has three remote data stores from three different third-party vendors whose data you want to analyze. You run the analytics environment configuration script one time and provide the information for each data store vendor as you are prompted for it. (When prompted for a data store vendor that you do not have, simply ignore that set of prompts.)

When you have different versions of the same vendor's data store, specifying the vendor's *latest* client data libraries usually works. However, this choice can be problematic for different versions of Hadoop, where a later set of JAR files is not typically backwardly compatible with earlier versions, or for sites that use Hadoop implementations from more than one vendor. (The configuration script does not delineate between different Hadoop vendors.) In these situations, you must run the analytics environment configuration script once for each different Hadoop version or vendor. As the configuration script creates a `TKGrid_REP`

directory underneath the current directory, it is important to run the script a second time from a different directory.

To illustrate how you might manage configuring the analytics environment for two different Hadoop vendors, consider this example: suppose your site uses Cloudera Hadoop 4 and Hortonworks Data Platform 2. When running the analytics environment script to configure for Cloudera 4, you would create a directory similar to:

```
cdh4
```

When configuring the analytics environment for Cloudera, you would run the script from the `cdh4` directory. When complete, the script creates a `TKGrid_REP` child directory:

```
cdh4/TKGrid_REP
```

For Hortonworks, you would create a directory similar to:

```
hdp2
```

When configuring the analytics environment for Hortonworks, you would run the script from the `hdp2` directory. When complete, the script creates a `TKGrid_REP` child directory:

```
hdp2/TKGrid_REP
```

## Configure for Access to a Data Store with a SAS Embedded Process

To configure the High-Performance Analytics environment for a remote data store, follow these steps:

- 1 Make sure that you have reviewed all of the information contained in the section [“Preparing for a Remote Parallel Connection”](#) on page 109.
- 2 Make sure that you understand how the analytics environment configuration script works, as described in [“How the Configuration Script Works”](#) on page 113.
- 3 The software that is needed for the analytics environment is available from within the SAS Software Depot that was created by the site depot administrator: `depot-installation-location/standalone_installs/SAS_High-Performance_Node_Installation/2_94/Linux_for_x64`.
- 4 Copy the `TKGrid_REP` file that is appropriate for your operating system to the `/tmp` directory of the root node of the analytic cluster.
- 5 Log on to the machine that will serve as the root node of the cluster with a user account that has the necessary permissions.

For more information, see [“User Accounts for the SAS High-Performance Analytics Environment”](#) on page 28.

- 6 Change directories to the desired installation location, such as `/opt`.
- 7 Run the shell script in this directory.

The shell script creates the `TKGrid_REP` subdirectory and places all files under that directory.

- 8 Respond to the prompts from the configuration program:

**Table 7.8** Configuration Parameters for the TGrid\_REP Shell Script

Parameter	Description
Do you want to configure remote access to Teradata? (yes/NO)	If you are using a Teradata Managed Cabinet for your data provider, specify <b>yes</b> and press Enter. Otherwise, specify <b>no</b> and press Enter.
Enter path of Teradata client install. i.e.: /opt/teradata/client/13.10	If you specified <b>no</b> in the previous step, specify the path where the Teradata client was installed and press Enter. (This path was recorded earlier in <a href="#">Table 7.7 on page 112.</a> )
Do you want to configure remote access to Greenplum? (yes/NO)	If you are using a Greenplum Data Computing Appliance for your data provider, specify <b>yes</b> and press Enter. Otherwise, specify <b>no</b> and press Enter.
Enter path of Greenplum client install. i.e.: /usr/local/greenplum-db	If you specified <b>no</b> in the previous step, specify the path where the Greenplum client was installed and press Enter. (This path was recorded earlier in <a href="#">Table 7.3 on page 111.</a> )
Do you want to configure remote access to Hadoop? (yes/NO)	If you are using a Hadoop machine cluster for your data provider, specify <b>yes</b> and press Enter. Otherwise, specify <b>no</b> and press Enter.
Enter path of 64 bit JRE i.e.: /usr/java/jdk1.7.0_09/jre	If you chose <b>no</b> in the previous step, specify the path where the JRE is installed and press Enter. (This path was recorded earlier in <a href="#">Table 7.2 on page 111.</a> )
Enter path of the directory containing the Hadoop and client jars.	Specify the path where the Cloudera Hadoop JAR files required by SAS reside and press Enter. (This path was recorded earlier in <a href="#">Table 7.1 on page 110.</a> )
Do you want to configure remote access to Oracle? (yes/NO)	If you are using an ORACLE Exadata appliance for your data provider, specify <b>yes</b> and press Enter. Otherwise, specify <b>no</b> and press Enter.
Enter path of Oracle client libraries. i.e.: /usr/local/ora11gr2/product/11.2.0/client_1/lib	Enter the path where the Oracle client libraries reside and press Enter. (This path was recorded earlier in <a href="#">Table 7.5 on page 112.</a> )
Enter path of TNS_ADMIN, or just enter if not needed.	Enter the value of the Oracle TNS_ADMIN environment variable and press Enter. (This value was recorded earlier in <a href="#">Table 7.6 on page 112.</a> )
Do you want to configure remote access to SAP HANA? (yes/NO)	If you are using a HANA cluster for your data provider, specify <b>yes</b> and press Enter. Otherwise, specify <b>no</b> and press Enter.
Enter path of HANA client install. i.e.: /usr/local/lib/hdbclient	Enter the path where the HANA client libraries reside and press Enter. (This path was recorded earlier in <a href="#">Table 7.4 on page 111.</a> )
Shared install or replicate to each node? (Y=SHARED/n=replicated)	If you are installing to a local drive on each node, then select <b>no</b> and press Enter to indicate that this is a replicated installation. If you are installing to a drive that is shared across all the nodes (for example, NFS), then specify <b>yes</b> and press Enter.

Parameter	Description
Enter path to TKGrid install	Specify the absolute path to where the SAS High-Performance Analytics environment is installed and press Enter. This should be the directory in which the analytics environment install program was run with <b>TKGrid</b> appended to it (for example, <code>/opt/TKGrid</code> ).  For more information, see <a href="#">Step 6 on page 76</a> .
Enter additional paths to include in LD_LIBRARY_PATH, separated by colons (:)	If you have any external library paths that you want to be accessible to the SAS High-Performance Analytics environment, specify the paths here and press Enter. Separate paths with a colon (:). If you have no paths to specify, press Enter.

- 9** If you selected a replicated installation at the first prompt, you are now prompted to choose the technique for distributing the contents to the appliance nodes:

```
The install can now copy this directory to all the machines
listed in 'pathname' using scp, skipping the first entry. Perform copy? (YES/no)
```

Press Enter if you want the installation program to perform the replication. Enter **no** if you are distributing the contents of the installation directory by some other technique.

- 10** You have finished deploying the analytics environment for a remote data source. If you have not done so already, install the appropriate SAS Embedded Process on the remote data appliance or machine cluster for your respective data provider.

For more information, see *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.

- 11** To validate your analytics environment, proceed to “Validating the Analytics Environment Deployment” on page 80.

# Appendix 1

## Updating the SAS High-Performance Analytics Infrastructure

<i>Overview of Updating the Analytics Infrastructure</i> .....	<b>117</b>
<i>Updating the SAS High-Performance Computing Management Console</i> ....	<b>117</b>
Overview of Updating the Management Console .....	117
Update the Management Console Using RPM .....	118
<i>Updating SAS High-Performance Deployment of Hadoop</i> .....	<b>118</b>
Overview of Updating SAS High-Performance Deployment of Hadoop .....	118
Preparing to Update Hadoop .....	119
Update SAS High-Performance Deployment of Hadoop (SAS LASR Adapter Components Only) .....	120
Update SAS High-Performance Deployment of Hadoop .....	121
<i>Update the Analytics Environment</i> .....	<b>127</b>

### Overview of Updating the Analytics Infrastructure

Here are some considerations for updating the SAS High-Performance Analytics infrastructure:

- Because of dependencies, if you update the analytics environment, you must also update SAS High-Performance Deployment of Hadoop.
- Update Hadoop first, followed by the analytics environment.

### Updating the SAS High-Performance Computing Management Console

#### Overview of Updating the Management Console

Starting in version 2.6 of SAS High-Performance Computing Management Console, there is no longer support for memory management through CGroups.

Before upgrading the management console to version 2.6, make sure that you manually record any memory settings and then clear them on the **CGroup**

**Resource Management** page. You can manually transfer these memory settings to the SAS High-Performance Analytics environment resource settings file. Or, if you are implementing YARN, transfer these settings to YARN. For more information, see *SAS LASR Analytic Server: Reference Guide*, available at <http://support.sas.com/documentation/solutions/va/index.html>.

## Update the Management Console Using RPM

To update your deployment of SAS High-Performance Computing Management Console, follow these steps:

- 1 Make sure that you have manually recorded and then cleared any memory settings in the management console. For more information, see “[Overview of Updating the Management Console](#)” on page 117.

- 2 Stop the server by entering the following command as the `root` user:

```
service sashpcmc stop
```

- 3 Update the management console using the following RPM command:

```
rpm -U --prefix=install-directory  
/SAS-Software-Depot-root-directory/standalone_installs/  
SAS_High-Performance_Computing_Management_Console/2_6/Linux_for_x64/  
sashpcmc-2.6.x86_64.rpm
```

In this command, *install-directory* is the location where the management console is installed and *SAS-Software-Depot-root-directory* is the location where your SAS Software Depot resides.

- 4 Log on to the console to validate your update.

---

## Updating SAS High-Performance Deployment of Hadoop

### Overview of Updating SAS High-Performance Deployment of Hadoop

The SAS High-Performance Deployment of Hadoop package consists of the following major components:

- Apache Hadoop
- LASR Analytic Server Hadoop adapter components (JAR files and shared libraries)

SAS gives you two options for updating SAS High-Performance Deployment of Hadoop:

- Update LASR Analytic Server Hadoop adapter components only:

You update the LASR Analytic Server Hadoop adapter components (JAR files and shared libraries) only. Apache Hadoop and the HDFS file system are not modified.

This approach is simpler than a full Hadoop upgrade, and has a lesser impact from a change management perspective.

For more information, see “Update SAS High-Performance Deployment of Hadoop (SAS LASR Adapter Components Only)” on page 120.

- Update SAS High-Performance Deployment of Hadoop:

You update the LASR Analytic Server Hadoop adapter components (JAR files and shared libraries), Apache Hadoop, and the HDFS file system. Your data that resides in your current version of Hadoop will be upgraded in place. The new version of Hadoop will access that data.

This approach is more complicated than updating the LASR Hadoop adapter components only, and has a greater impact from a change management perspective.

For more information, see “Update SAS High-Performance Deployment of Hadoop” on page 121.

## Preparing to Update Hadoop

Prior to starting the SAS High-Performance Deployment of Hadoop update, perform the following steps:

**Note:** The following steps also apply when you are upgrading SAS LASR adapter components only.

- 1 If one does not already exist, create a SAS Software Depot that contains the installation software that you will use to update Hadoop.

For more information, see “Creating a SAS Software Depot” in the *SAS Intelligence Platform: Installation and Configuration Guide*, available at <http://support.sas.com/documentation/cdl/en/biig/63852/HTML/default/p03intellplatform00installgd.htm>.

- 2 Log on to the Hadoop NameNode as the `hdfs` user.
- 3 Run the following command to make sure that the Hadoop file system is healthy: `hadoop fsck /`  
Correct any issues before proceeding.
- 4 Stop any other processes, such as YARN, running on the Hadoop cluster.  
Confirm that all processes have stopped across all the cluster machines. (You might have to become another user to have the necessarily privileges to stop all processes.)
- 5 As the `hdfs` user, run the command `$HADOOP_HOME/sbin/stop-dfs.sh` to stop HDFS daemons, and confirm that all processes have ceased on all the machines in the cluster.

**TIP** Check that there are no Java processes owned by `hadoop` running on any machine: `ps -ef | grep hadoop`. If you find any Java processes owned by the `hdfs` user account, terminate them. You can issue a single `simsh` command to simultaneously check all the machines in the cluster: `/HPA-environment-installation-directory/bin/simsh ps -ef | grep hdfs`.

- 6 Back up the Hadoop name directory (`hadoop-name` by default).

Perform a file system backup using tar (or whatever tool or process that your site uses for backups).

## Update SAS High-Performance Deployment of Hadoop (SAS LASR Adapter Components Only)

The Hadoop install script gives you the option of upgrading of the LASR Analytic Server Hadoop adapter components (JAR files and shared libraries) only. When you choose this option, the install script does *not* modify Apache Hadoop and the HDFS file system.

To update LASR Analytic Server Hadoop adapter components (JAR files and shared libraries) only, follow these steps:

- 1 Make sure that you have performed the steps listed in the section [“Preparing to Update Hadoop” on page 119](#).

- 2 Log on to the Hadoop NameNode as the user ID that owns your current Hadoop installation directories.

- 3 Copy the `sashadoop.tar.gz` file to a temporary location and extract it:

```
cp sashadoop.tar.gz /tmp
cd /tmp
tar xzf sashadoop.tar.gz
```

A directory that is named `sashadoop` is created.

- 4 Change directory to the `sashadoop` directory and run the `hadoopInstall` command:

```
cd sashadoop
./hadoopInstall
```

- 5 Respond to the prompts from the configuration program:

**Table A1.1** SAS High-Performance Deployment of Hadoop Configuration Parameters

Parameter	Description
Choose the type of installation to perform: 1) New installation of SAS Apache Hadoop 2.4.0 with new HDFS. 2) Add the latest LASR support to an existing SAS Apache Hadoop. Leave existing HDFS unmodified. 3) New installation of SAS Apache Hadoop 2.4.0 with upgrade of your existing HDFS directory structure. 4) Quit. [This utility is not used with 3rd-party Hadoop distributions.] Enter choice (1-4). Default is 4: (1/2/3/4)?	Specify 2 and press Enter to perform a new installation.  If you want to upgrade Hadoop (option 3), see <a href="#">“Update SAS High-Performance Deployment of Hadoop” on page 121</a> .
Enter path to existing Hadoop installation.	Specify the value of <code>HADOOP_HOME</code> (for example, <code>/opt/hadoop/hadoop-0.23.1</code> ) and press Enter.



Parameter	Description
Supported version of Hadoop found at: '/opt/hadoop/hadoop-0.23.1' Updating Hadoop install at: '/opt/hadoop/hadoop-0.23.1' Stop Hadoop server at: '/opt/hadoop/hadoop-0.23.1', and Hit Return.	Be sure that the Hadoop server is stopped (\$HADOOP_HOME/sbin/stop-dfs.sh) and press Enter.

The install script outputs messages similar to the following:

```
Verify that the following lines are in '/opt/hadoop/hadoop-0.23.1/etc/hadoop/hdfs-site.xml'.
```

```
<property>
  <name>dfs.permissions.enabled</name>
  <value>true</value>
</property>
<property>
  <name>dfs.namenode.plugins</name>
  <value>com.sas.lasr.hadoop.NameNodeService</value>
</property>
<property>
  <name>dfs.datanode.plugins</name>
  <value>com.sas.lasr.hadoop.DataNodeService</value>
</property>
<property>
  <name>com.sas.lasr.hadoop.fileinfo</name>
  <value>ls -l {0}</value>
  <description>The command used to get the user, group, and permission
  information for a file.
  </description>
</property>
<property>
  <name>com.sas.lasr.service.allow.put</name>
  <value>true</value>
  <description>Flag indicating whether the PUT command is enabled when
  running as a service. The default is false.
  </description>
</property>
```

```
Installation complete. Please restart your Hadoop server.
```

- 6 Verify that each on node, the hdfs-site.xml file contains the earlier listed properties.
- 7 Restart Hadoop by entering the following command:

```
$HADOOP_HOME/sbin/start-dfs.sh
```

## Update SAS High-Performance Deployment of Hadoop

Version 2.6 of SAS High-Performance Deployment of Hadoop represents a version upgrade of Apache Hadoop (version 0.23.1 to version 2.4). This newer version includes new features such as YARN. During an upgrade, the install

script installs a new version of Hadoop. Your data that resides in your current version of Hadoop will be upgraded in place. The new version of Hadoop will access that data.

Before you update Hadoop, you must gather the following information listed in [Table A1.2](#). You can find most of this information in your current Hadoop configuration file, `$HADOOP_HOME/etc/hadoop/hdfs-site.xml`:

**Table A1.2** Hadoop Installation Checklist

Installation Prompt	Requirement / How to Locate	Actual Value
Hadoop install directory	One level above the current HADOOP_HOME value. For example: <code>/hadoop</code>	
Replication factor	Refer to <code>hdfs-site.xml</code> .	
Port for <code>fs.defaultFS</code>	Refer to <code>core-site.xml</code> .	
Port for <code>mapred.job.tracker</code>	Refer to <code>mapred_site.xml</code>	
Port for <code>dfs.datanode.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.namenode.backup.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.namenode.https-address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.datanode.https.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.datanode.ipc.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.namenode.http-address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.datanode.http.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.secondary.http.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.namenode.backup.address</code>	Refer to <code>hdfs-site.xml</code> .	
Port for <code>dfs.namenode.backup.http-address</code>	Refer to <code>hdfs-site.xml</code> .	

Installation Prompt	Requirement / How to Locate	Actual Value
Port for com.sas.lasr.hadoop.service.namenode.port	Refer to hdfs-site.xml.	
Port for com.sas.lasr.hadoop.service.datanode.port	Refer to hdfs-site.xml.	
HDFS server process user	Must be the same user as current Hadoop user.	
Path for JAVA_HOME directory	Location of your JRE installation (default: /usr/lib/jvm/jre).	
Path for Hadoop data directory	Same as the current Hadoop data directory. Refer to hdfs-site.xml.	
Path for Hadoop name directory	Same as the current Hadoop name directory. Refer to hdfs-site.xml.	
Path to machine list	See <a href="#">“List the Machines in the Cluster or Appliance”</a> .	

To update SAS High-Performance Deployment of Hadoop, follow these steps:

- 1 Make sure that you have performed the steps listed in the section [“Preparing to Update Hadoop” on page 119](#).
- 2 Log on to the Hadoop NameNode as the root user.
- 3 Copy the sashadoop.tar.gz file to a temporary location and extract it:

```
cp sashadoop.tar.gz /tmp
cd /tmp
tar xzf sashadoop.tar.gz
```

A directory that is named **sashadoop** is created.

- 4 Change directory to the **sashadoop** directory and run the **hadoopInstall** command:

```
cd sashadoop
./hadoopInstall
```

- 5 Using the information that you gathered earlier in [Table A1.2](#), respond to the prompts from the configuration program:

**Table A1.3** SAS High-Performance Deployment of Hadoop Configuration Parameters

Parameter	Description
<p>Choose the type of installation to perform:</p> <p>1) New installation of SAS Apache Hadoop 2.4.0 with new HDFS.</p> <p>2) Add the latest LASR support to an existing SAS Apache Hadoop. Leave existing HDFS unmodified.</p> <p>3) New installation of SAS Apache Hadoop 2.4.0 with upgrade of your existing HDFS directory structure.</p> <p>4) Quit.</p> <p>[This utility is not used with 3rd-party Hadoop distributions.]</p> <p>Enter choice (1-4). Default is 4: (1/2/3/4)?</p>	<p>Specify 3 and press Enter to perform a new installation.</p> <p>If you want to upgrade SAS LASR adapter components only (option 2), see <a href="#">“Update SAS High-Performance Deployment of Hadoop (SAS LASR Adapter Components Only)”</a> on page 120.</p>
<p>Enter path to install Hadoop. The directory 'hadoop-2.4.0' will be created in the path specified.</p>	<p>Specify the directory one level above the current HADOOP_HOME and press Enter. Refer to <a href="#">Table A1.2</a>.</p>
<p>Do you wish to use Yarn and MR Jobhistory Server? (y/N)</p>	<p>If you plan to use YARN and MapReduce, specify <b>y</b> and press Enter. If you are using YARN, be sure to review <a href="#">“Preparing for YARN (Experimental)”</a> on page 27 before proceeding.</p> <p>Otherwise, specify <b>n</b> and press Enter.</p>
<p>Enter replication factor. Default 2</p>	<p>Specify the replication factor used for your current Hadoop deployment and press Enter. Refer to <a href="#">Table A1.2</a>.</p>
<p>Enter port number for fs.defaultFS. Default 54310</p> <p>Enter port number for dfs.namenode.https-address. Default 50470</p> <p>Enter port number for dfs.datanode.https.address. Default 50475</p> <p>Enter port number for dfs.datanode.address. Default 50010</p> <p>Enter port number for dfs.datanode.ipc.address. Default 50020</p> <p>Enter port number for dfs.namenode.http-address. Default 50070</p> <p>Enter port number for dfs.datanode.http.address. Default 50075</p> <p>Enter port number for dfs.secondary.http.address. Default 50090</p> <p>Enter port number for dfs.namenode.backup.address. Default 50100</p> <p>Enter port number for dfs.namenode.backup.http-address. Default 50105</p> <p>Enter port number for com.sas.lasr.hadoop.service.namenode.port. Default 15452</p> <p>Enter port number for com.sas.lasr.hadoop.service.datanode.port. Default 15453</p>	<p>Specify each port and press Enter. Refer to <a href="#">Table A1.2</a>.</p>

Parameter	Description
<p>[The following port prompts are displayed when you choose to deploy YARN:]</p> <p>Enter port number for mapreduce.jobhistory.admin.address. Default 10033</p> <p>Enter port number for mapreduce.jobhistory.webapp.address. Default 19888</p> <p>Enter port number for mapreduce.jobhistory.address. Default 10021</p> <p>Enter port number for yarn.resourcemanager.scheduler.address. Default 8030</p> <p>Enter port number for yarn.resourcemanager.resource- tracker.address. Default 8031</p> <p>Enter port number for yarn.resourcemanager.address. Default 8032</p> <p>Enter port number for yarn.resourcemanager.admin.address. Default 8033</p> <p>Enter port number for yarn.resourcemanager.webapp.address. Default 8088</p> <p>Enter port number for yarn.nodemanager.localizer.address. Default 8040</p> <p>Enter port number for yarn.nodemanager.webapp.address. Default 8042</p> <p>Enter port number for yarn.web-proxy.address. Default 10022</p>	<p>Specify each port and press Enter. Refer to <a href="#">Table A1.2</a>.</p>
Enter maximum memory allocation per Yarn container. Default 5905	This is the maximum amount of memory (in MB) that YARN can allocate on a particular machine in the cluster. Press Enter to accept the default or specify a different value and press Enter.
Enter user that will be running the HDFS server process.	Specify the user name (for example, hdfs) and press Enter. Refer to <a href="#">Table A1.2</a> .
Enter user that will be running Yarn services.	Specify the user name (for example, yarn) and press Enter. For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter user that will be running the Map Reduce Job History Server.	Specify the user name (for example, mapred) and press Enter. For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter common primary group for users running Hadoop services.	Apache recommends that the hdfs, mapred, and YARN users share the same primary Linux group. Enter a group name and press Enter. For more information, see <a href="#">“Preparing for YARN (Experimental)” on page 27</a> .
Enter path for JAVA_HOME directory. (Default: /usr/lib/jvm/jre)	<p>Specify the path to the JRE or JDK and press Enter. Refer to <a href="#">Table A1.2</a>.</p> <p><b>Note:</b> The configuration program does not verify that a JRE is installed at <code>/usr/lib/jvm/jre</code>, which is the default path for some Linux vendors.</p>

Parameter	Description
Enter path for Hadoop data directory. This should be on a large drive. Default is '/hadoop/hadoop-data'.	Specify the paths to your current Hadoop data and name directories and press Enter. Refer to <a href="#">Table A1.2</a> .
Enter path for Hadoop name directory. Default is '/hadoop/hadoop-name'.	<p><b>Note:</b> The data directory cannot be the root directory of a partition or mount.</p> <p><b>Note:</b> If you have more than one data device, enter one of the data directories now, and refer to “(Optional) Deploy with Multiple Data Devices” on page 53 after the installation.</p>
Enter full path to machine list. The NameNode 'host' should be listed first.	Specify the path to your current machine list and press Enter. Refer to <a href="#">Table A1.2</a> .

- 6 You will see failure to create directory errors for a directory other than the `hadoop-2.4.0` directory. These are normal, since the directories being created already exist. These errors occur on all nodes after you confirm that you want the installation program to copy the installation to all nodes.

**CAUTION!** After the installation is complete, do *not* reformat the NameNode. Reformatting the Hadoop NameNode deletes your data in the HDFS cluster.

- 7 Log out as the root user. Log in as the `hdfs` user.
- 8 Run this command to define `HADOOP_HOME` in the Hadoop user's (`hdfs`) environment:
- ```
export HADOOP_HOME=/installation-directory/hadoop/hadoop-2.4.0
```
- where *installation-directory* is the location where you installed Hadoop (for example, `/opt/hadoop/hadoop-2.4.0`).
- 9 Run the following command to start Hadoop:
- ```
$HADOOP_HOME/sbin/start-dfs.sh -upgrade.
```
- 10 Run the following command: `$HADOOP_HOME/bin/hadoop fsck /`.
- You should see a healthy file system and the correct number of DataNodes.
- 11 The `initial-sas-hdfs-setup.sh` script makes modifications required for Hadoop, such as creating some new directories that support YARN and applying permissions that improve security. Review the `hdfs fs` commands that are listed in `$HADOOP_HOME/sbin/initial-sas-hdfs-setup.sh` and then run this script once.
- Alternatively, you can run individual commands from the script if you understand how the commands modify HDFS.
- 12 Confirm that Hadoop is running successfully by opening a browser to `http://namenode:50070/dfshealth.html`. Review the information in the cluster summary section of the page. Confirm that the number of live nodes equals the number of DataNodes and that the number of dead nodes is zero.
- 13 If you do *not* plan to update the SAS High-Performance Analytics environment, then you must manually update the analytics environment to reflect the new `HADOOP_HOME` value. Do this by editing

`$GRIDINSTALLLOC/tkmpirsh.sh`. Then copy this file to the same location across all the machines in the cluster. For example:

```
/opt/TKGrid/bin/simcp $GRIDINSTALLLOC/tkmpirsh.sh $GRIDINSTALLLOC/tkmpirsh.sh
```

## Update the Analytics Environment

You have the following options for managing updates to the SAS High-Performance Analytics environment:

- Delete the SAS High-Performance Analytics environment and install the newer version.

See the procedure later in this topic.

- Rename the root installation directory for the current SAS High-Performance Analytics environment, and install the newer version under the previous root installation directory.

See “Install the Analytics Environment” on page 76.

- Do nothing to the current SAS High-Performance Analytics environment, and install the new version under a new installation directory.

See “Install the Analytics Environment” on page 76.

When you change the path of the SAS High-Performance Analytics environment, you have to also have to reconfigure the SAS LASR Analytic Server to point to the new path. See “Add a SAS LASR Analytic Server” in Chapter 5 of *SAS Visual Analytics: Administration Guide* available at <http://support.sas.com/documentation/solutions/va/index.html>.

Updating your deployment of the SAS High-Performance Analytics environment consists of deleting the deployment and reinstalling the newer version. To update the SAS High-Performance Analytics environment, follow these steps:

- 1 Check that there are no analytics environment processes running on any machine:

```
ps -ef | grep TKGrid
```

If you find any TKGrid processes, terminate them.

**TIP** You can issue a single `simsh` command to simultaneously check all the machines in the cluster: `/HPA-environment-installation-directory/bin/simsh ps -ef | grep TKGrid`.

- 2 Delete the analytics environment installation directory on every machine in the cluster:

```
rm -r -f /HPA-environment-install-dir
```

**TIP** You can issue a single `simsh` command to simultaneously remove the environment install directories on all the machines in the cluster: `/HPA-environment-installation-directory/bin/simsh rm -r -f /HPA-environment-installation-directory`.

- 3** Re-install the analytics environment using the shell script as described in [“Install the Analytics Environment” on page 76](#).



# Appendix 2

## SAS High-Performance Analytics Infrastructure Command Reference

The `simsh` and `simcp` commands are installed with SAS High-Performance Computing Management Console and the SAS High-Performance Analytics environment. The default path to the commands is `/HPCMC-installation-directory/webmin/utilbin` and `/HPA-environment-installation-directory/bin`, respectively. Any user account that can access the commands and has passwordless secure shell configured can use them.

**TIP** Add one of the earlier referenced installation paths to your system PATH variable to make invoking `simsh` and `simcp` easier.

The `simsh` command uses secure shell to invoke the specified command on every machine that is listed in the `/etc/gridhosts` file. The following command demonstrates invoking the `hostname` command on each machine in the cluster:

```
/HPCMC-install-dir/webmin/utilbin/simsh hostname
```

**TIP** You can use SAS High-Performance Computing Management Console to create and manage your grid hosts file. For more information, see *SAS High-Performance Computing Management Console: User's Guide*, available at <http://support.sas.com/documentation/onlinedoc/va/index.html>.

The `simcp` command is used to copy a file from one machine to the other machines in the cluster. Passwordless secure shell and an `/etc/gridhosts` file are required. The following command is an example of copying the `/etc/hosts` file to each machine in the cluster:

```
/HPA-environment-installation-directory/bin/simcp /etc/hosts /etc
```



# Appendix 3

## SAS High-Performance Analytics Environment Client-Side Environment Variables

The following environment variables can be used on the client side to control the connection to the SAS High-Performance Analytics environment. You can set these environment variables in the following ways:

- invoke them in your SAS program using `options set=`
- add them to your shell before running the SAS program
- add them to your `sasenv_local` configuration file, if you want them used in all SAS programs

### GRIDHOST=

identifies the root node on the SAS High-Performance Analytics environment to which the client connects.

The values for GRIDHOST and GRIDINSTALLLOC can both be specified in the GRIDHOST variable, separated by a colon (similar to the format used by `scp`). For example:

```
GRIDHOST=my_machine_cluster_001:/opt/TKGrid
```

### GRIDINSTALLLOC=

identifies the location on the machine cluster where the SAS High-Performance Analytics environment is installed. For example:

```
GRIDINSTALLLOC=/opt/TKGrid
```

### GRIDMODE=SYM | ASYM

toggles the SAS High-Performance Analytics environment between symmetric (default) and asymmetric mode.

### GRIDRSHCOMMAND= " " | " *ssh-path* "

(optional) specifies `rsh` or `ssh` used to launch the SAS High-Performance Analytics environment.

If unspecified or a null value is supplied, a SAS implementation of the SSH protocol is used.

*ssh-path* specifies the path to the SSH executable that you want to use. This can be useful in deployments where export controls restrict SAS from delivering software that uses cryptography. For example:

```
option set=GRIDRSHCOMMAND="/usr/bin/ssh";
```

**GRIDPORTRANGE=**

identifies the port range for the client to open. The root node connects back to the client using ports in the specified range. For example:

```
option set=GRIDPORTRANGE=7000-8000;
```

**GRIDREPLYHOST=**

specifies the name of the client machine to which the SAS High-Performance Analytics environment connects. `GRIDREPLYHOST` is used when the client has more than one network card or when you need to specify a full network name.

`GRIDREPLYHOST` can be useful when you need to specify a fully qualified domain name, when the client has more than one network interface card, or when you need to specify an IP address for a client with a dynamically assigned IP address that domain name resolution has not registered yet. For example:

```
GRIDREPLYHOST=myclient.example.com
```

# Appendix 4

## Deploying on SELinux and IPTables

<i>Overview of Deploying on SELinux and IPTables</i> .....	<b>133</b>
<i>Prepare the Management Console</i> .....	<b>134</b>
SELinux Modifications for the Management Console .....	134
IPTables Modifications for the Management Console .....	134
<i>Prepare Hadoop</i> .....	<b>134</b>
SELinux Modifications for Hadoop .....	134
IPTables Modifications for Hadoop .....	134
<i>Prepare the Analytics Environment</i> .....	<b>135</b>
SELinux Modifications for the Analytics Environment .....	135
IPTables Modifications for the Analytics Environment .....	135
<i>Analytics Environment Post-Installation Modifications</i> .....	<b>135</b>
<i>iptables File</i> .....	<b>136</b>

---

## Overview of Deploying on SELinux and IPTables

This document describes how to prepare Security Enhanced Linux (SELinux) and IPTables for a SAS High-Performance Analytics infrastructure deployment.

Security Enhanced Linux (SELinux) is a feature in some versions of Linux that provides a mechanism for supporting access control security policies. IPTables is a firewall—a combination of a packet-filtering framework and generic table structure for defining rulesets. SELinux and IPTables is available in most new distributions of Linux, both community-based and enterprise-ready. For sites that require added security, the use of SELinux and IPTables is an accepted approach for many IT departments.

Because of the limitless configuration possibilities, this document is based on the default configuration for SELinux and IPTables running on RedHat Enterprise Linux (RHEL) 6.3. You might need to adjust the directions accordingly, especially for complex SELinux and IPTables configurations.

---

## Prepare the Management Console

### SELinux Modifications for the Management Console

After generating and propagating root's SSH keys throughout the cluster or data appliance, you must run the following command on every machine or blade to restore the security context on the files in `/root/.ssh`:

```
restorecon -R -v /root/.ssh
```

### IPTables Modifications for the Management Console

Add the following line to `/etc/sysconfig/iptables` to allow connections to the port on which the management console is listening (10020 by default). Open the port only on the machine on which the management console is running:

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 10020 -j ACCEPT
```

---

## Prepare Hadoop

### SELinux Modifications for Hadoop

After generating and propagating root's SSH keys throughout the cluster or data appliance, you must run the following command on every machine or blade to restore the security context on the files in `/root/.ssh`:

```
restorecon -R -v /root/.ssh
```

### IPTables Modifications for Hadoop

The SAS High-Performance Deployment of Hadoop has a number of ports on which it communicates. To open these ports, place the following lines in `/etc/sysconfig/iptables`:

**Note:** The following example uses default ports. Modify as necessary for your site.

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 54310 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 54311 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50470 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50475 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50010 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50020 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50070 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50075 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50090 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50100 -j ACCEPT
```

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50105 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50030 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50060 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 15452 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 15453 -j ACCEPT
```

Edit `/etc/sysconfig/iptables` and then copy this file across the machine cluster or data appliance. Lastly, restart the IPTables service.

---

## Prepare the Analytics Environment

### SELinux Modifications for the Analytics Environment

After generating and propagating root's SSH keys throughout the cluster or data appliance, you must run the following command on every machine or blade to restore the security context on the files in `/root/.ssh`:

```
restorecon -R -v /root/.ssh
```

### IPTables Modifications for the Analytics Environment

If you are deploying the SAS LASR Analytic Server, then you must define one port per server in `/etc/sysconfig/iptables`. (The port number is defined in the SAS code that starts the SAS LASR Analytic server.)

If you have more than one server running simultaneously, you need all these ports defined in the form of a range.

The following is an example of an iptables entry for a single server (one port):

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 10010 -j ACCEPT
```

The following is an example of an iptables entry for five servers (port range):

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 10010:10014 -j ACCEPT
```

MPICH\_PORT\_RANGE must also be opened in IPTables by editing the `/etc/sysconfig/iptables` file and adding the port range.

The following is an example for five servers:

```
-A INPUT -m state --state NEW -m tcp -p tcp --dport 10010:10029 -j ACCEPT
```

Edit `/etc/sysconfig/iptables` and then copy this file across the machine cluster or data appliance. Lastly, restart the IPTables service.

---

## Analytics Environment Post-Installation Modifications

The SAS High-Performance Analytics environment uses Message Passing Interface (MPI) communications, which requires you to define one port range per active job across the machine cluster or data appliance.

(A port range consists of a minimum of four ports per active job. Every running monitoring server counts as a job on the cluster or appliance.)

For example, if you have five jobs running simultaneously across the machine cluster or data appliance, you need a minimum of 20 ports in the range.

The following example is an entry in `tkmpirsh.sh` for five jobs:

```
export MPICH_PORT_RANGE=18401:18420
```

Edit `tkmpirsh.sh` using the number of jobs appropriate for your site. (`tkmpirsh.sh` is located in `/installation-directory/TKGrid/.`) Then, copy `tkmpirsh.sh` across the machine cluster or data appliance.

---

## iptables File

This topic lists the complete `/etc/sysconfig/iptables` file. The additions to iptables described in this document are highlighted.

```
*filter
:INPUT ACCEPT [0:0]
:FORWARD ACCEPT [0:0]
:OUTPUT ACCEPT [0:0]
-A INPUT -m state --state ESTABLISHED,RELATED -j ACCEPT
-A INPUT -p icmp -j ACCEPT
-A INPUT -i lo -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 22 -j ACCEPT
# Needed by SAS HPC MC
-A INPUT -m state --state NEW -m tcp -p tcp --dport 10020 -j ACCEPT
# Needed for HDFS (Hadoop)
A INPUT -m state --state NEW -m tcp -p tcp --dport 54310 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 54311 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50470 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50475 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50010 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50020 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50070 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50075 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50090 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50100 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50105 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50030 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 50060 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 15452 -j ACCEPT
-A INPUT -m state --state NEW -m tcp -p tcp --dport 15453 -j ACCEPT
# End of HDFS Additions
# Needed for LASR Server Ports.
-A INPUT -m state --state NEW -m tcp -p tcp --dport 17401:17405 -j ACCEPT
# End of LASR Additions
# Needed for MPICH.
-A INPUT -m state --state NEW -m tcp -p tcp --dport 18401:18420 -j ACCEPT
# End of MPICH additions.
-A INPUT -j REJECT --reject-with icmp-host-prohibited
-A FORWARD -j REJECT --reject-with icmp-host-prohibited
```



## Recommended Reading

Here is the recommended reading list for this title:

- *Configuration Guide for SAS Foundation for Microsoft Windows for x64*, available at <http://support.sas.com/documentation/installcenter/en/ikfdtnwx6cg/66385/PDF/default/config.pdf>.
- *Configuration Guide for SAS Foundation for UNIX Environments*, available at <http://support.sas.com/documentation/installcenter/en/ikfdtnunxcg/66380/PDF/default/config.pdf>.
- *SAS/ACCESS for Relational Databases: Reference*, <http://support.sas.com/documentation/cdl/en/acrelldb/67473/PDF/default/acrelldb.pdf>.
- *SAS Deployment Wizard and SAS Deployment Manager: User's Guide*, available at <http://support.sas.com/documentation/installcenter/en/ikdeploywizug/66034/PDF/default/user.pdf>.
- *SAS Guide to Software Updates*, available at <http://support.sas.com/documentation/cdl/en/whatsdiff/66129/PDF/default/whatsdiff.pdf>.
- *SAS High-Performance Computing Management Console: User's Guide*, available at <http://support.sas.com/documentation/solutions/hpainfrastructure/>.
- *SAS In-Database Products: Administrator's Guide*, available at <http://support.sas.com/documentation/cdl/en/indbag/67365/PDF/default/indbag.pdf>.
- *SAS Intelligence Platform: Installation and Configuration Guide*, available at <http://support.sas.com/documentation/cdl/en/biig/63852/PDF/default/biig.pdf>.
- *SAS Intelligence Platform: Security Administration Guide*, available at <http://support.sas.com/documentation/cdl/en/bisecag/67045/PDF/default/bisecag.pdf>.

For a complete list of SAS publications, go to [sas.com/store/books](http://sas.com/store/books). If you have questions about which titles you need, please contact a SAS Representative:

SAS Books  
 SAS Campus Drive  
 Cary, NC 27513-2414  
 Phone: 1-800-727-0025  
 Fax: 1-919-677-4444  
 E-mail: [sasbook@sas.com](mailto:sasbook@sas.com)  
 Web address: [sas.com/store/books](http://sas.com/store/books)



# Glossary

**data set**

See SAS data set

**encryption**

the act or process of converting data to a form that is unintelligible except to the intended recipients.

**foundation services**

See SAS Foundation Services

**grid host**

the machine to which the SAS client makes an initial connection in a SAS High-Performance Analytics application.

**Hadoop Distributed File System**

a framework for managing files as blocks of equal size, which are replicated across the machines in a Hadoop cluster to provide fault tolerance.

**HDFS**

See Hadoop Distributed File System

**identity**

See metadata identity

**Integrated Windows authentication**

a Microsoft technology that facilitates use of authentication protocols such as Kerberos. In the SAS implementation, all participating components must be in the same Windows domain or in domains that trust each other.

**Internet Protocol Version 6**

See IPv6

**IPv6**

a protocol that specifies the format for network addresses for all computers that are connected to the Internet. This protocol, which is the successor of Internet Protocol Version 4, uses hexadecimal notation to represent 128-bit address spaces. The format can consist of up to eight groups of four hexadecimal characters, delimited by colons, as in FE80:0000:0000:0000:0202:B3FF:FE1E:8329. As an alternative, a group of consecutive zeros could be replaced with two colons, as in FE80::0202:B3FF:FE1E:8329. Short form: IPv6

**IWA**

See Integrated Windows authentication

**JAR file**

a Java Archive file. The JAR file format is used for aggregating many files into one file. JAR files have the file extension .jar.

**Java**

a set of technologies for creating software programs in both stand-alone environments and networked environments, and for running those programs safely. Java is an Oracle Corporation trademark.

**Java Database Connectivity**

See JDBC

**Java Development Kit**

See JDK

**JDBC**

a standard interface for accessing SQL databases. JDBC provides uniform access to a wide range of relational databases. It also provides a common base on which higher-level tools and interfaces can be built. Short form: JDBC.

**JDK**

a software development environment that is available from Oracle Corporation. The JDK includes a Java Runtime Environment (JRE), a compiler, a debugger, and other tools for developing Java applets and applications. Short form: JDK.

**localhost**

the keyword that is used to specify the machine on which a program is executing. If a client specifies localhost as the server address, the client connects to a server that runs on the same machine.

**login**

a SAS copy of information about an external account. Each login includes a user ID and belongs to one SAS user or group. Most logins do not include a password.

**Message Passing Interface**

is a message-passing library interface specification. SAS High-Performance Analytics applications implement MPI for use in high-performance computing environments.

**metadata identity**

a metadata object that represents an individual user or a group of users in a SAS metadata environment. Each individual and group that accesses secured resources on a SAS Metadata Server should have a unique metadata identity within that server.

**metadata object**

a set of attributes that describe a table, a server, a user, or another resource on a network. The specific attributes that a metadata object includes vary depending on which metadata model is being used.

**middle tier**

in a SAS business intelligence system, the architectural layer in which Web applications and related services execute. The middle tier receives user

requests, applies business logic and business rules, interacts with processing servers and data servers, and returns information to users.

**MPI**

See Message Passing Interface

**object spawner**

a program that instantiates object servers that are using an IOM bridge connection. The object spawner listens for incoming client requests for IOM services. When the spawner receives a request from a new client, it launches an instance of an IOM server to fulfill the request. Depending on which incoming TCP/IP port the request was made on, the spawner either invokes the administrator interface or processes a request for a UUID (Universal Unique Identifier).

**planned deployment**

a method of installing and configuring a SAS business intelligence system. This method requires a deployment plan that contains information about the different hosts that are included in the system and the software and SAS servers that are to be deployed on each host. The deployment plan then serves as input to the SAS Deployment Wizard.

**root node**

in a SAS High-Performance Analytics application, the role of the software that distributes and coordinates the workload of the worker nodes. In most deployments the root node runs on the machine that is identified as the grid host. SAS High-Performance Analytics applications assign the highest MPI rank to the root node.

**SAS Application Server**

a logical entity that represents the SAS server tier, which in turn comprises servers that execute code for particular tasks and metadata objects.

**SAS authentication**

a form of authentication in which the target SAS server is responsible for requesting or performing the authentication check. SAS servers usually meet this responsibility by asking another component (such as the server's host operating system, an LDAP provider, or the SAS Metadata Server) to perform the check. In a few cases (such as SAS internal authentication to the metadata server), the SAS server performs the check for itself. A configuration in which a SAS server trusts that another component has pre-authenticated users (for example, Web authentication) is not part of SAS authentication.

**SAS configuration directory**

the location where configuration information for a SAS deployment is stored. The configuration directory contains configuration files, logs, scripts, repository files, and other items for the SAS software that is installed on the machine.

**SAS data set**

a file whose contents are in one of the native SAS file formats. There are two types of SAS data sets: SAS data files and SAS data views.

**SAS Deployment Manager**

a cross-platform utility that manages SAS deployments. The SAS Deployment Manager supports functions such as updating passwords for your SAS deployment, rebuilding SAS Web applications, and removing configurations.

**SAS Deployment Wizard**

a cross-platform utility that installs and initially configures many SAS products. Using a SAS installation data file and, when appropriate, a deployment plan for its initial input, the wizard prompts the customer for other necessary input at the start of the session, so that there is no need to monitor the entire deployment.

**SAS Foundation Services**

a set of core infrastructure services that programmers can use in developing distributed applications that are integrated with the SAS platform. These services provide basic underlying functions that are common to many applications. These functions include making client connections to SAS application servers, dynamic service discovery, user authentication, profile management, session context management, metadata and content repository access, activity logging, event management, information publishing, and stored process execution.

**SAS installation data file**

See SID file

**SAS installation directory**

the location where your SAS software is installed. This location is the parent directory to the installation directories of all SAS products. The SAS installation directory is also referred to as SAS Home in the SAS Deployment Wizard.

**SAS IOM workspace**

in the IOM object hierarchy for a SAS Workspace Server, an object that represents a single session in SAS.

**SAS Metadata Server**

a multi-user server that enables users to read metadata from or write metadata to one or more SAS Metadata Repositories.

**SAS Pooled Workspace Server**

a SAS Workspace Server that is configured to use server-side pooling. In this configuration, the SAS object spawner maintains a collection of workspace server processes that are available for clients.

**SAS Software Depot**

a file system that consists of a collection of SAS installation files that represents one or more orders. The depot is organized in a specific format that is meaningful to the SAS Deployment Wizard, which is the tool that is used to install and initially configure SAS. The depot contains the SAS Deployment Wizard executable, one or more deployment plans, a SAS installation data file, order data, and product data.

**SAS Stored Process Server**

a SAS IOM server that is launched in order to fulfill client requests for SAS Stored Processes.

**SAS Workspace Server**

a SAS IOM server that is launched in order to fulfill client requests for IOM workspaces.

**SASHDAT file**

the data format used for tables that are added to HDFS by SAS. SASHDAT files are read in parallel by the server.

**SASHOME directory**

the file location where an instance of SAS software is installed on a computer. The location of the SASHOME directory is established at the initial installation of SAS software by the SAS Deployment Wizard. That location becomes the default installation location for any other SAS software you install on the same machine.

**server context**

a SAS IOM server concept that describes how SAS Application Servers manage client requests. A SAS Application Server has an awareness (or context) of how it is being used and makes decisions based on that awareness. For example, when a SAS Data Integration Studio client submits code to its SAS Application Server, the server determines what type of code is submitted and directs it to the correct physical server for processing (in this case, a SAS Workspace Server).

**server description file**

a file that is created by a SAS client when the LASR procedure executes to create a server. The file contains information about the machines that are used by the server. It also contains the name of the server signature file that controls access to the server.

**SID file**

a control file containing license information that is required in order to install SAS.

**spawner**

See object spawner

**worker node**

in a SAS High-Performance Analytics application, the role of the software that receives the workload from the root node.

**workspace**

See SAS IOM workspace





# Index

## A

accounts  
     See [user accounts](#)  
 Authen::PAM PERL [25](#)  
 authorized\_keys file [34](#)

## C

checklists  
     pre-installation for port numbers [30](#)  
 configuration  
     Hadoop [89](#)  
 configuration file  
     FILENAME Statement Hadoop  
         Access Method [96](#)  
     High Performance Analytics for  
         Hadoop [96](#)  
     PROC Hadoop [96](#)  
     SAS Code Accelerator for Hadoop  
         [96](#)  
     SAS Scoring Accelerator [96](#)  
     SAS/ACCESS Interface to Hadoop  
         [96](#)  
     SPD Engine [96](#)

## D

deployment  
     overview [9](#)  
 depot  
     See [SAS Software Depot](#)

## G

gridhosts file [25](#)  
 groups  
     setting up [14](#), [33](#), [71](#), [85](#)

## H

Hadoop  
     client-side JAR files [97](#)  
     configuration file [96](#)  
     in-database deployment package  
         [86](#)  
     installation and configuration [89](#)  
     permissions [104](#)  
     SAS/ACCESS Interface [86](#)  
     starting the SAS Embedded Process  
         [103](#)  
     status of the SAS Embedded  
         Process [104](#)  
     stopping the SAS Embedded  
         Process [103](#)  
     unpacking self-extracting archive  
         files [92](#)

## I

in-database deployment package for  
     Hadoop  
         overview [88](#)  
         prerequisites [86](#)  
 installation [1](#)  
     Hadoop [89](#)  
     SAS Embedded Process (Hadoop)  
         [88](#), [92](#)  
     SAS Hadoop MapReduce JAR files  
         [92](#)

## K

keys  
     See [SSH public key](#)

## M

middle tier shared key  
     propagate [41](#)

**O**

operating system accounts  
     See [user accounts](#)

**P**

perl-Net-SSLeay [25](#)  
 permissions  
     for Hadoop [104](#)  
 ports  
     designating [30](#)  
     reserving for SAS [30](#)  
 pre-installation checklists  
     for port numbers [30](#)  
 publishing  
     Hadoop permissions [104](#)

**R**

reinstalling a previous version  
     Hadoop [89](#)  
 required user accounts [14](#), [33](#), [71](#), [85](#)  
 requirements, system [10](#)  
 reserving ports  
     SAS [30](#)

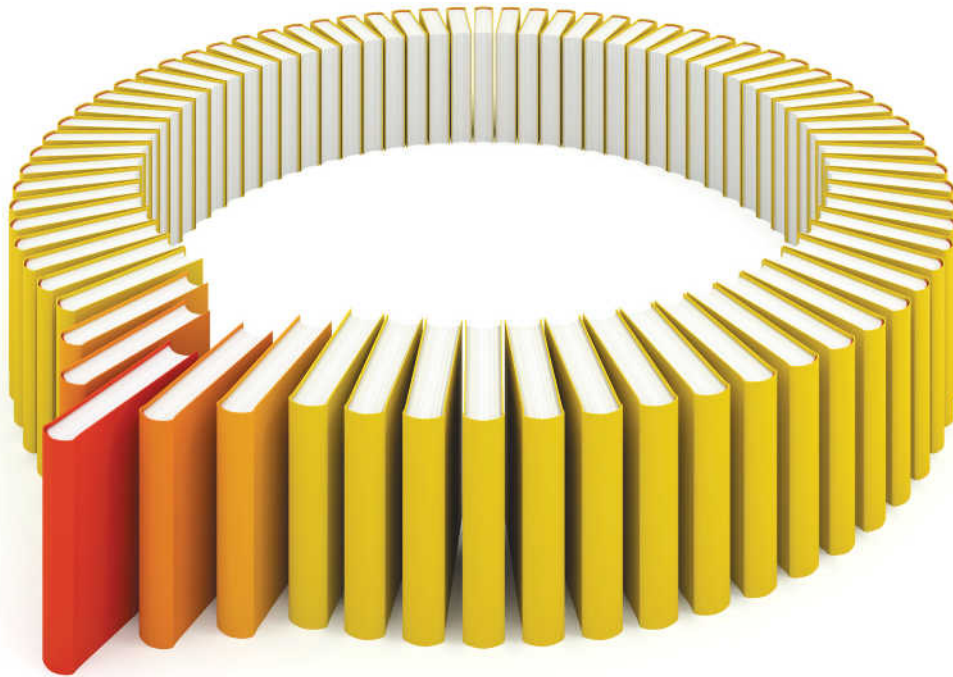
**S**

SAS Embedded Process  
     controlling (Hadoop) [98](#)  
     Hadoop [86](#)  
 SAS Foundation [86](#)  
 SAS Hadoop MapReduce JAR files  
     [92](#)  
 SAS High-Performance Computing  
     Management Console  
     create user accounts [41](#)  
     deployment [34](#)  
     logging on [38](#)

    middle tier shared key [41](#)  
 SAS High-Performance Computing  
     Management Console server  
     starting [35](#)  
 SAS Software Depot [25](#)  
 SAS system accounts [14](#), [33](#), [71](#), [85](#)  
 SAS Visual Analytics  
     deploying [9](#)  
 SAS/ACCESS Interface to Hadoop [86](#)  
 sasep-servers.sh script  
     overview [98](#)  
     syntax [99](#)  
 secure shell [25](#)  
     JBoss Application Server public key  
         [38](#)  
     propagate keys [41](#)  
 self-extracting archive files  
     unpacking for Hadoop [92](#)  
 server  
     SAS High-Performance Computing  
         Management Console [35](#)  
 SSH  
     See [secure shell](#)  
 SSH public key  
     JBoss Application Server [38](#)  
 SSH public keys  
     propagate [41](#)  
 SSL [36](#)  
 system requirements [10](#)

**U**

unpacking self-extracting archive files  
     for Hadoop [92](#)  
 upgrading from a previous version  
     Hadoop [89](#)  
 user accounts [14](#), [33](#), [71](#), [85](#)  
     JBoss Application Server [38](#)  
     SAS system accounts [14](#), [33](#), [71](#),  
         [85](#)  
     setting up required accounts [14](#), [33](#),  
         [71](#), [85](#)



# Gain Greater Insight into Your SAS<sup>®</sup> Software with SAS Books.

Discover all that you need on your journey to knowledge and empowerment.



SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are trademarks of their respective companies. © 2013 SAS Institute Inc. All rights reserved. S107969US.0613

