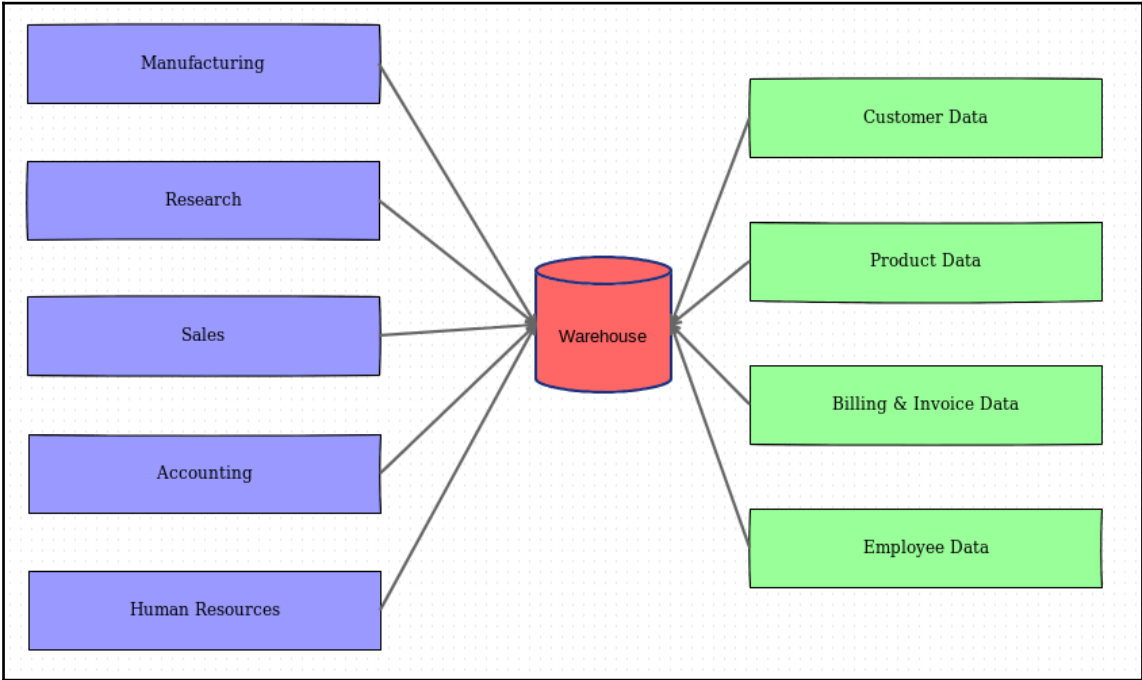
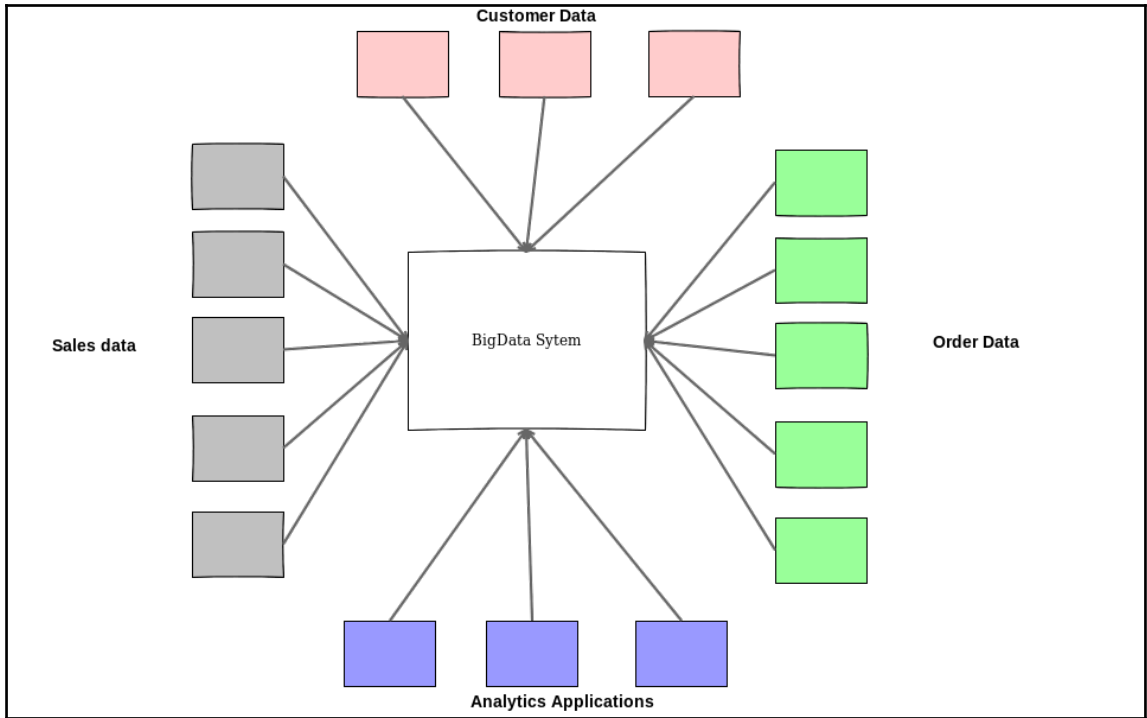
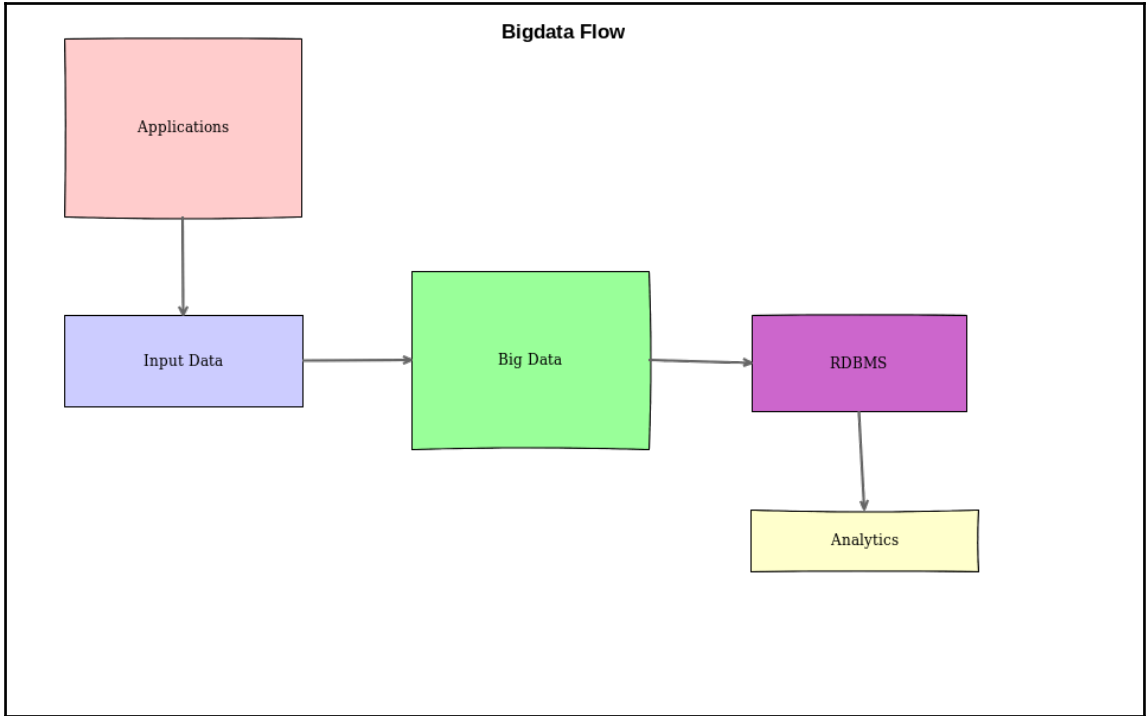
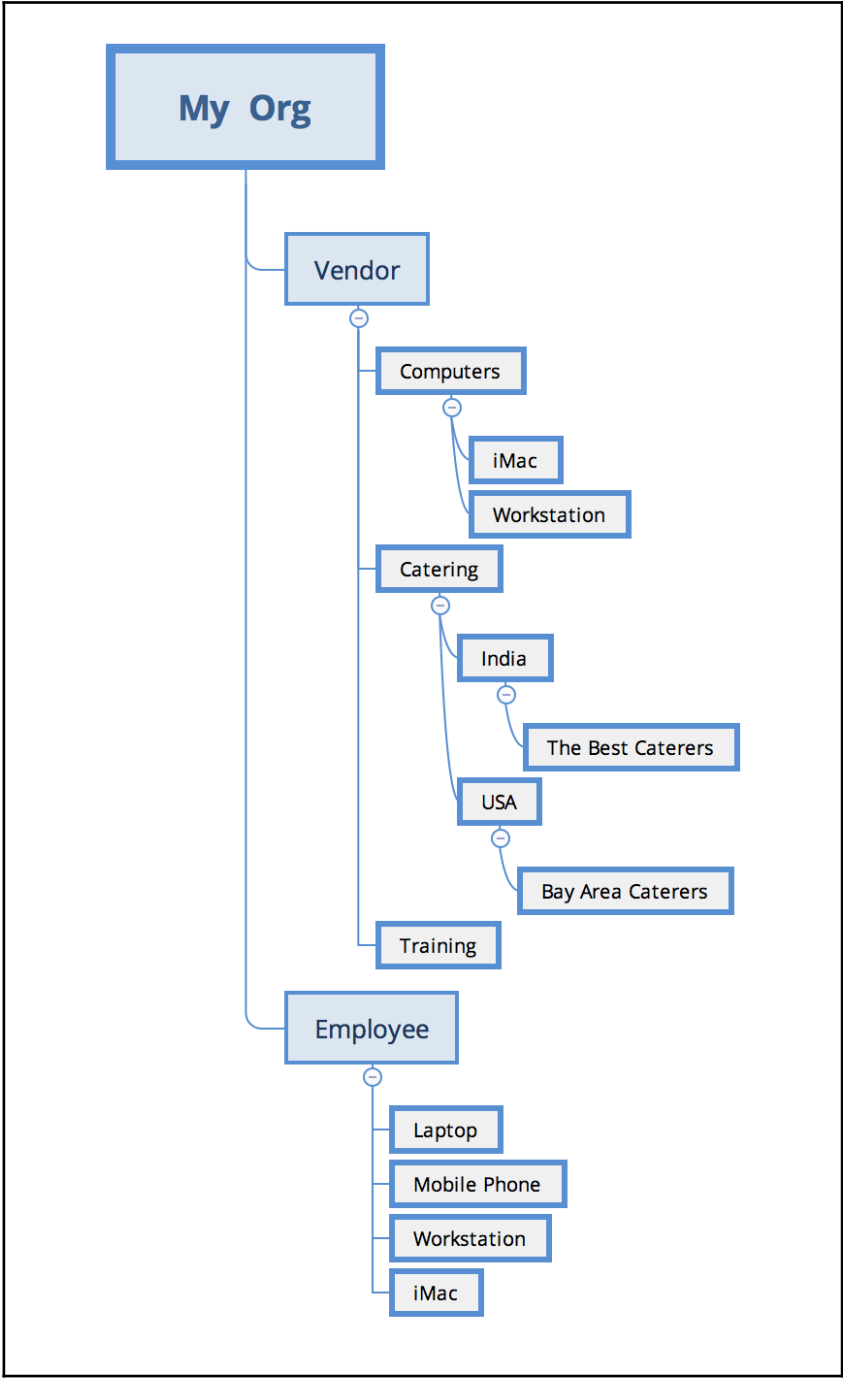


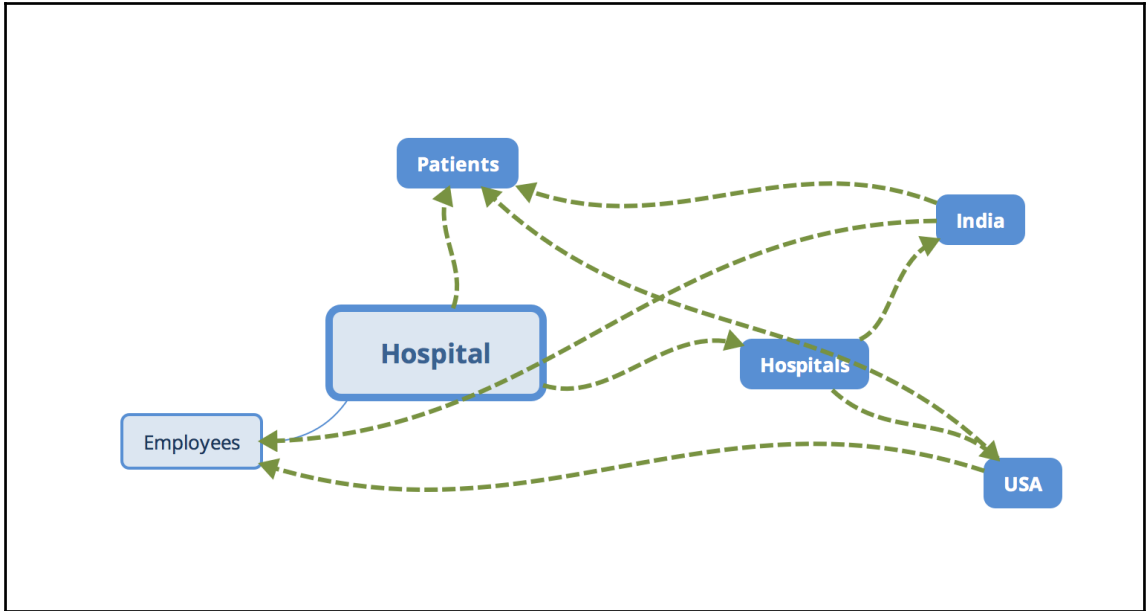
Chapter 1: Enterprise Data Architecture Principles

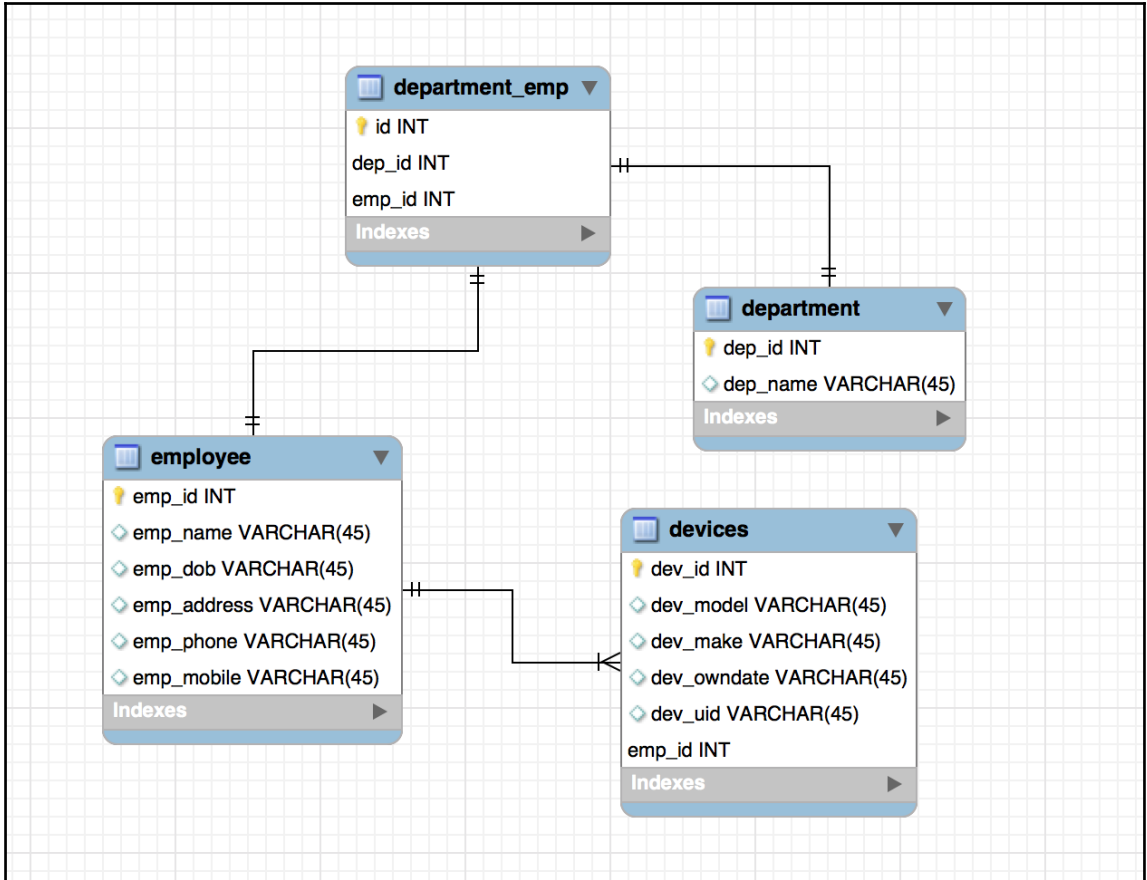


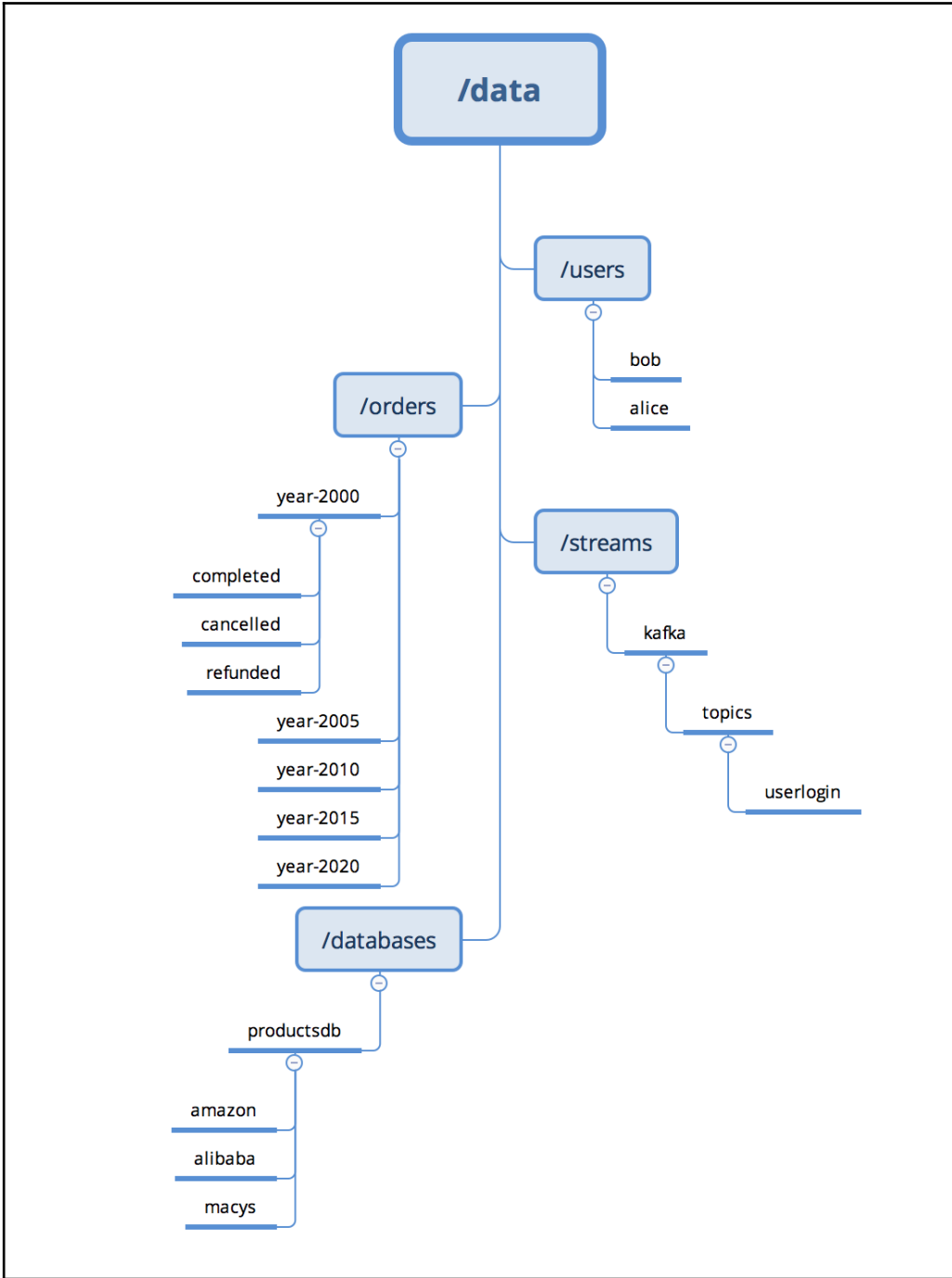






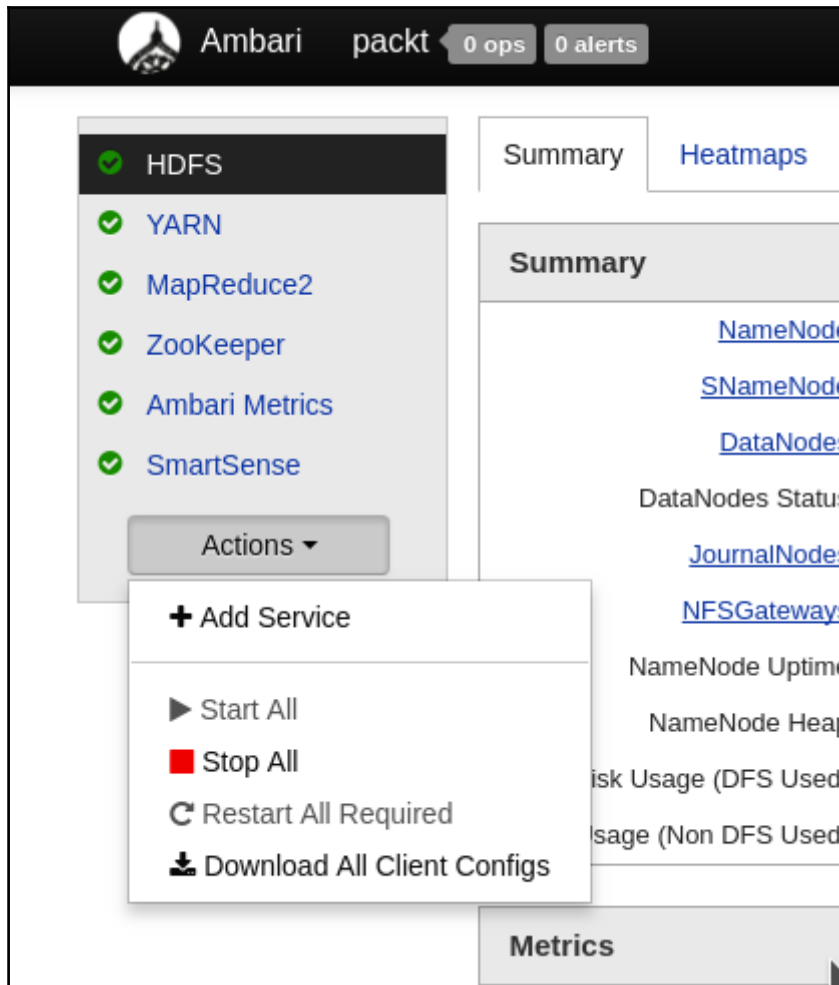






Chapter 2: Hadoop Life Cycle Management

| Input Data | | | | |
|----------------------|------------------|---------------|--------------|---------------|
| First Name | Last Name | Gender | Phone | Salary |
| Myra | Riley | Female | 550-7605-00 | 6502 |
| Preston | Davis | Male | 979-8283-17 | 4376 |
| Kimberly | Perkins | Female | 166-3607-48 | 7195 |
| Jordan | Myers | Male | 513-6151-62 | 5525 |
| Ryan | Harper | Male | 863-7471-20 | 2657 |
| | | | | |
| Shuffled Data | | | | |
| First Name | Last Name | Gender | Phone | Salary |
| Myra | Riley | Female | 550-7605-00 | 7195 |
| Preston | Davis | Male | 979-8283-17 | 4376 |
| Kimberly | Perkins | Female | 166-3607-48 | 5525 |
| Jordan | Myers | Male | 513-6151-62 | 2657 |
| Ryan | Harper | Male | 863-7471-20 | 6502 |



Add Service Wizard

| | | | |
|-------------------------------------|----------------|-------------------|---|
| <input type="checkbox"/> | ZooKeeper | 3.4.0 | Centralized service which provides highly reliable distributed coordination |
| <input type="checkbox"/> | Falcon | 0.10.0 | Data management and processing platform |
| <input type="checkbox"/> | Storm | 1.1.0 | Apache Hadoop Stream processing framework |
| <input type="checkbox"/> | Flume | 1.5.2 | A distributed service for collecting, aggregating, and moving large amounts of streaming data into HDFS |
| <input type="checkbox"/> | Accumulo | 1.7.0 | Robust, scalable, high performance distributed key/value store. |
| <input type="checkbox"/> | Ambari Infra | 0.1.0 | Core shared service used by Ambari managed components. |
| <input checked="" type="checkbox"/> | Ambari Metrics | 0.1.0 | A system for metrics collection that provides storage and retrieval capability for metrics collected from the cluster |
| <input type="checkbox"/> | Atlas | 0.8.0 | Atlas Metadata and Governance platform |
| <input type="checkbox"/> | Kafka | 0.10.1 | A high-throughput distributed messaging system |
| <input type="checkbox"/> | Knox | 0.12.0 | Provides a single point of authentication and access for Apache Hadoop services in a cluster |
| <input type="checkbox"/> | Log Search | 0.5.0 | Log aggregation, analysis, and visualization for Ambari managed services. This service is Technical Preview . |
| <input checked="" type="checkbox"/> | Ranger | 0.7.0 | Comprehensive security for Hadoop |
| <input type="checkbox"/> | Ranger KMS | 0.7.0 | Key Management Server |
| <input checked="" type="checkbox"/> | SmartSense | 1.4.3.2.6.0.0-267 | SmartSense - Hortonworks SmartSense Tool (HST) helps quickly gather configuration, metrics, logs from common HDP services that aids to quickly troubleshoot support cases and receive cluster-specific recommendations. |

Add Service Wizard

Review

Install, Start and Test

Summary

ResourceManager:

App Timeline Server:

History Server:

ZooKeeper Server:

ZooKeeper Server:

ZooKeeper Server:

Metrics Collector:

Grafana:

Ranger Usersync:

Ranger Admin:

HST Server:

Activity Explorer:

Activity Analyzer:

HST Server

Activity Explorer

Activity Analyzer

node-2.c.coastal-airlock-197705.internal (12.6 GB, 2 cores)

SNameNode

ResourceManager

App Timeline Server

History Server

ZooKeeper Server

node-3.c.coastal-airlock-197705.internal (12.6 GB, 2 cores)

ZooKeeper Server

Metrics Collector

Ranger Usersync

Ranger Admin

Add Service Wizard X

ADD SERVICE WIZARD

- [Choose Services](#)
- [Assign Masters](#)
- Assign Slaves and Clients
- [Customize Services](#)
- [Configure Identities](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on.
Hosts that are assigned master components are shown with *.

| Host | all none | all none | all none | all none |
|---------------------------------|--|-------------------------------------|---|--|
| node-1.c.coastal-airlock-1... * | <input type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input type="checkbox"/> NodeManager | <input type="checkbox"/> Ranger Tagsync |
| node-2.c.coastal-airlock-1... * | <input type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input type="checkbox"/> NodeManager | <input type="checkbox"/> Ranger Tagsync |
| node-3.c.coastal-airlock-1... * | <input checked="" type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input checked="" type="checkbox"/> NodeManager | <input checked="" type="checkbox"/> Ranger Tagsync |

Show: 25
1 - 3 of 3
⏪ ⏩ ↺ ↻

← Back
Next →

```

welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MariaDB connection id is 17
Server version: 5.5.56-MariaDB MariaDB Server

Copyright (c) 2000, 2017, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MariaDB [(none)]> create database ranger;
Query OK, 1 row affected (0.00 sec)

MariaDB [(none)]> use ranger;
Database changed
MariaDB [ranger]> grant all privileges on ranger.* to ranger@localhost identified by 'ranger';
Query OK, 0 rows affected (0.01 sec)

MariaDB [ranger]> grant all privileges on ranger.* to ranger@%' identified by 'ranger';
Query OK, 0 rows affected (0.00 sec)

MariaDB [ranger]> flush privileges;
Query OK, 0 rows affected (0.00 sec)

MariaDB [ranger]> quit
Bye
[collabrecruit_in@master ~]$

```

Add Service Wizard

Ranger Admin

DB FLAVOR

Ranger DB name

Ranger DB username

JDBC connect string for a Ranger database

Ranger DB host

Driver class name for a JDBC Ranger database

Ranger DB password

Connection OK

Setup Database and Database User

Dependent Configurations

Recommended Changes

Based on your configuration changes, Ambari is recommending the following dependent configuration changes. Ambari will update all checked configuration changes to the **Recommended Value**. Uncheck any configuration to retain the **Current Value**.

| <input checked="" type="checkbox"/> Property | Service | Config Group | File Name | Current Value | Recommended Value |
|--|---------|--------------|-------------------------------|--------------------|---|
| <input checked="" type="checkbox"/> ranger-hdfs-plugin-enabled | HDFS | Default | ranger-hdfs-plugin-properties | Property undefined | No |
| <input checked="" type="checkbox"/> ranger.plugin.hdfs.policy.rest.url | HDFS | Default | ranger-hdfs-security | Property undefined | http://node-3.c.coastal-airlock-197705.internal:6080 |
| <input checked="" type="checkbox"/> xasecure.audit.destination.hdfs | HDFS | Default | ranger-hdfs-audit | Property undefined | true |
| <input checked="" type="checkbox"/> xasecure.audit.destination.hdfs.dir | HDFS | Default | ranger-hdfs-audit | Property undefined | hdfs://node-1.c.coastal-airlock-197705.internal:8020/ranger/audit |
| <input checked="" type="checkbox"/> xasecure.audit.destination.solr | HDFS | Default | ranger-hdfs-audit | Property undefined | false |
| <input checked="" type="checkbox"/> xasecure.audit.destination.solr.urls | HDFS | Default | ranger-hdfs-audit | Property undefined | |
| <input checked="" type="checkbox"/> xasecure.audit.destination.solr.zookeepers | HDFS | Default | ranger-hdfs-audit | Property undefined | NONE |
| <input checked="" type="checkbox"/> ranger-yarn-plugin-enabled | YARN | Default | ranger-yarn-plugin-properties | Property undefined | No |

Add Service Wizard

- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Configure Identities](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Please review the configuration before installation

Admin Name : admin

Cluster Name : packt

Total Hosts : 3 (0 new)

Repositories:

```

debian7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/debian7/2.x/updates/2.6.3.0
debian7 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/debian7
redhat-ppc7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos7-ppc/2.x/updates/2.6.3.0
redhat-ppc7 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/centos7-ppc
redhat6 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos6/2.x/updates/2.6.3.0
redhat6 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/centos6
redhat7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos7/2.x/updates/2.6.3.0

```

-- Back
Print
Deploy --

Add Service Wizard

- ADD SERVICE WIZARD
- [Choose Services](#)
- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Configure Identities](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Install, Start and Test

Please wait while the selected services are installed and started.

100 % overall

Show: **All (3)** | [In Progress \(0\)](#) | [Warning \(0\)](#) | [Success \(3\)](#) | [Fail \(0\)](#)

| Host | Status | Message |
|--|--|---------|
| node-1.c.coastal-airlock-197705.internal | 100% | Success |
| node-2.c.coastal-airlock-197705.internal | 100% | Success |
| node-3.c.coastal-airlock-197705.internal | 100% | Success |

3 of 3 hosts showing - [Show All](#) Show: 25 | 1 - 3 of 3

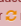
Successfully installed and started the services.

Next --

Add Service Wizard

- ADD SERVICE WIZARD
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary**

Summary


Important: You may also need to restart other services for the newly added services to function properly (for example, HDFS and YARN/MapReduce need to be restarted after adding Oozie). After closing this wizard, please restart all services that have the restart indicator  next to the service name.


Here is the summary of the install process.

The cluster consists of 3 hosts
Installed and started services successfully on 3 new hosts
Install and start completed in 6 minutes and 42 seconds

[Complete](#)

Ranger

 Username:

 Password:

[Sign In](#)

Ranger Access Manager Audit Settings admin

Service Manager

Service Manager

| | | |
|---|---|--|
| HDFS + [lock] [refresh] [delete] packt_hadoop [lock] [delete] | HBASE + [lock] [refresh] [delete] | HIVE + [lock] [refresh] [delete] |
| YARN + [lock] [refresh] [delete] | KNOX + [lock] [refresh] [delete] | STORM + [lock] [refresh] [delete] |
| SOLR + [lock] [refresh] [delete] | KAFKA + [lock] [refresh] [delete] packt_kafka [lock] [delete] | NIFI + [lock] [refresh] [delete] |
| ATLAS + [lock] [refresh] [delete] | | |

Ranger Access Manager Audit Settings

Service Manager > Edit Service

Edit Service

Service Details :

Service Name *

Description

Active Status Enabled Disabled

Select Tag Service

Config Properties :

Username *

Password *

Namenode URL *

Authorization Enabled

Authentication Type *

hadoop.security.auth_to_local

dfs.datanode.kerberos.principal

dfs.namenode.kerberos.principal

dfs.secondary.namenode.kerberos.principal

RPC Protection Type

Common Name for Certificate

Add New Configurations

| Name | Value | |
|---|--|----------------------------------|
| <input type="text" value="tag.download.auth.users"/> | <input type="text" value="hdfs"/> | <input type="button" value="x"/> |
| <input type="text" value="policy.download.auth.users"/> | <input type="text" value="hdfs"/> | <input type="button" value="x"/> |
| <input type="text" value="ambari.service.check.user"/> | <input type="text" value="ambari-qa"/> | <input type="button" value="x"/> |

Ranger Access Manager Audit Settings admin

Create Policy

Policy Details:

Policy Type: **Access**

Policy Name: PROJECT-A enabled

Resource Path: /project-a recursive

Audit Logging: YES

Description: project a

Allow Conditions:

| Select Group | Select User | Permissions | Delegate Admin | |
|--------------|----------------|--|--------------------------|-------------------------------------|
| Select Group | [X] hdfs-alice | Execute Read Write <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Select Group | [X] hdfs-bob | Read <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Select Group | [X] hdfs-tom | Read <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

```
[hdfs@node-2 ~]$ hdfs dfs -mkdir /projects
[hdfs@node-2 ~]$ hdfs dfs -ls /projects
[hdfs@node-2 ~]$
```

hdfs-alice@node-3:~ - Chromium

Secure | <https://ssh.cloud.google.com/projects/coastal-airlock-197705/zones/asia-south1-a/instances/n...>

```
[hdfs-alice@node-3 ~]$ hdfs dfs -mkdir /projects/1
[hdfs-alice@node-3 ~]$
```

hdfs-tom@node-1:~ - Chromium

Secure | <https://ssh.cloud.google.com/projects/coastal-airlock-197705/zones/asia-south1-a/instances/n...>

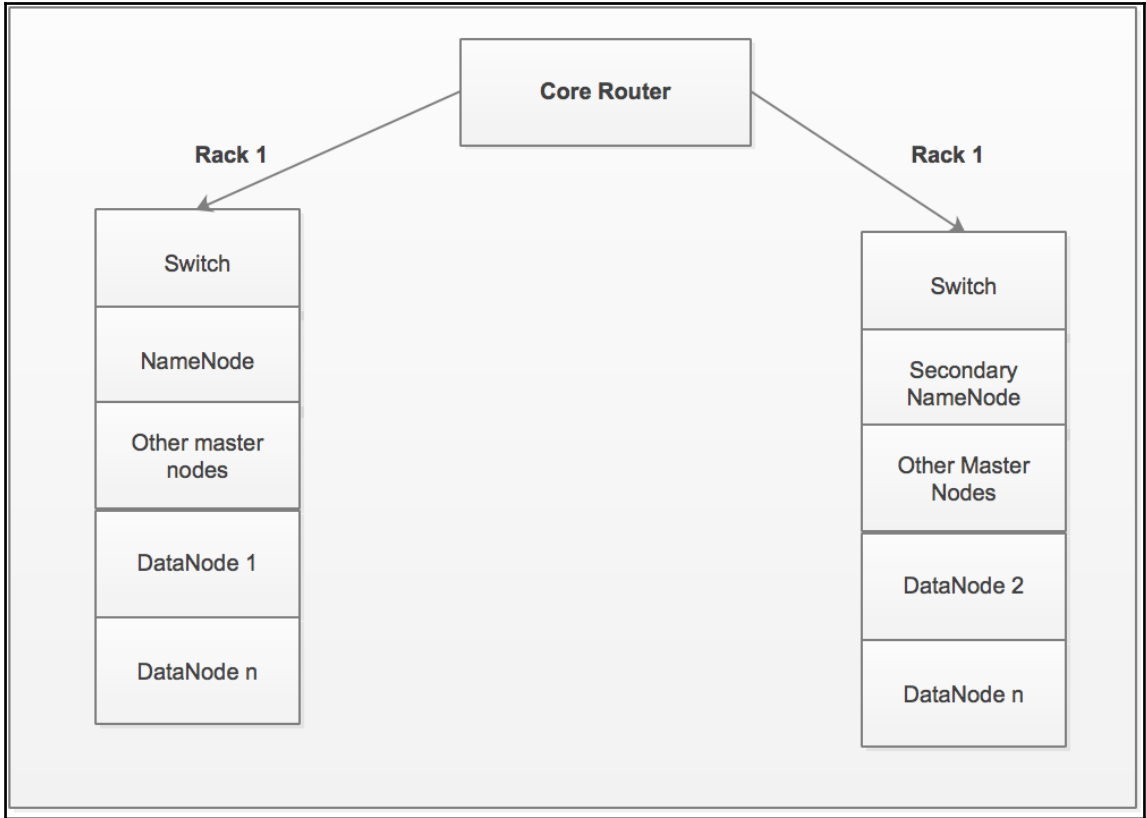
```
[hdfs-tom@node-1 ~]$ hdfs dfs -mkdir /projects/1
mkdir: Permission denied: user=hdfs-tom, access=WRITE, inode="/projects/1":hdfs:hdfs:drwxr-xr-x
[hdfs-tom@node-1 ~]$
```

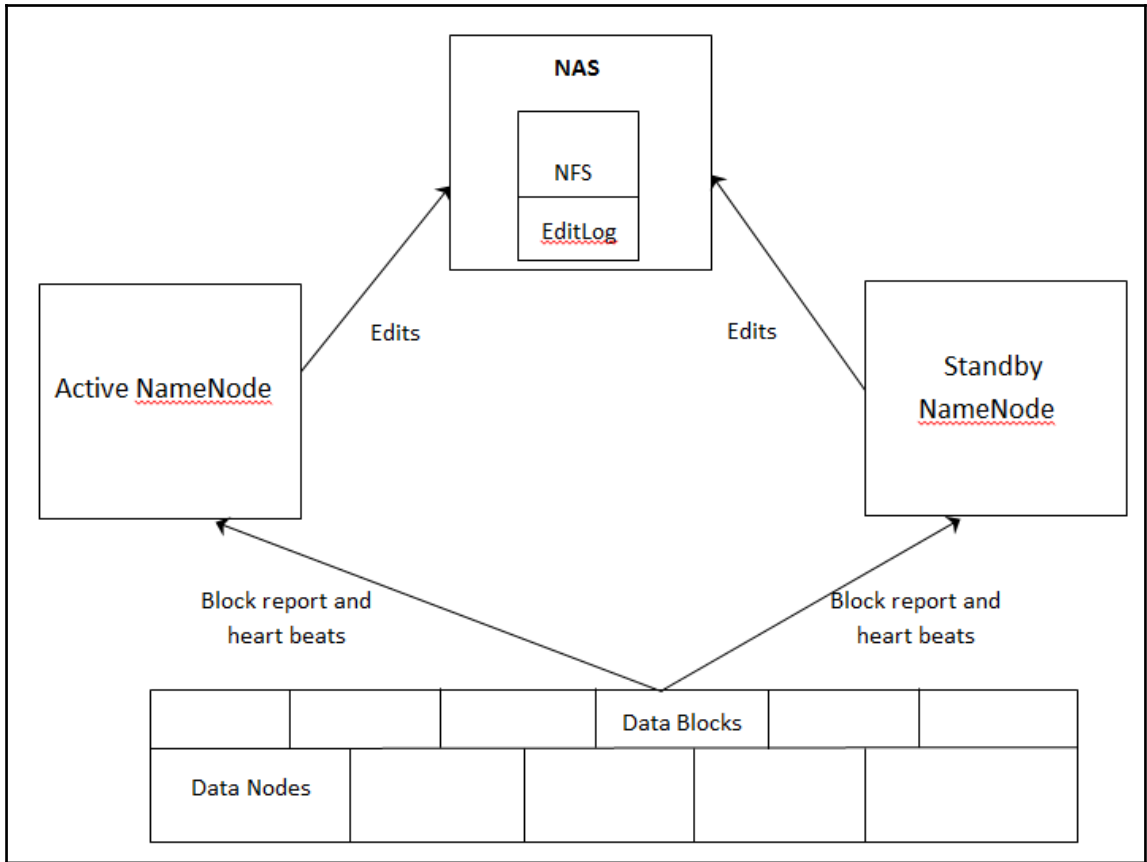
Ranger Access Manager Audit Settings admin

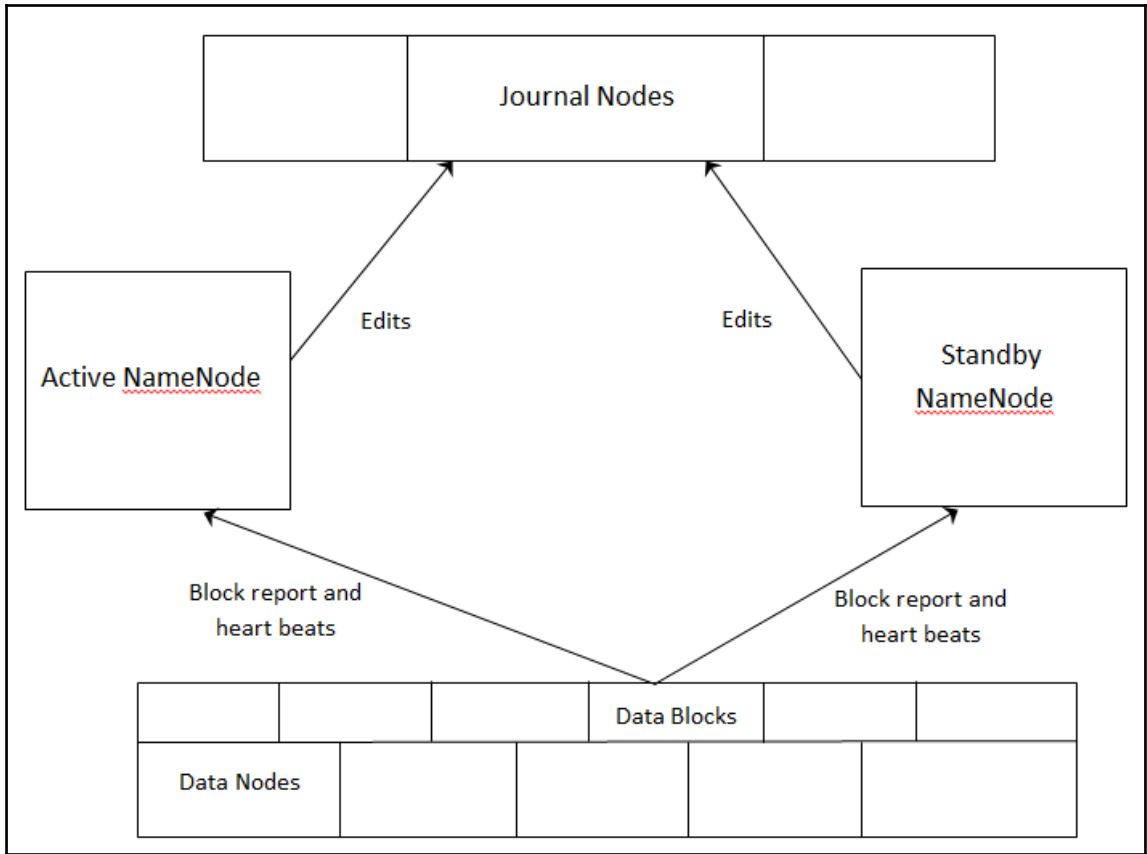
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | | |
|-----|------------------------|------------|--------------|------------------|-------------|--------------|---------|------------|------------|------------|-------|-----|-----|
| ... | 03/13/2018 02:39:03 PM | mapred | packt_hadoop | /tmp-history/tmp | path | READ_EXECUTE | Allowed | hadoop-act | 10.160.0.4 | packt | 1 | ... | |
| 9 | 03/13/2018 02:38:41 PM | hdfs-alice | packt_hadoop | hdfs | /projects/1 | path | WRITE | Allowed | ranger-act | 10.160.0.5 | packt | 1 | ... |
| ... | 03/13/2018 02:38:17 PM | hdfs-tom | packt_hadoop | hdfs | /projects/1 | path | WRITE | Denied | hadoop-act | 10.160.0.3 | packt | 1 | ... |

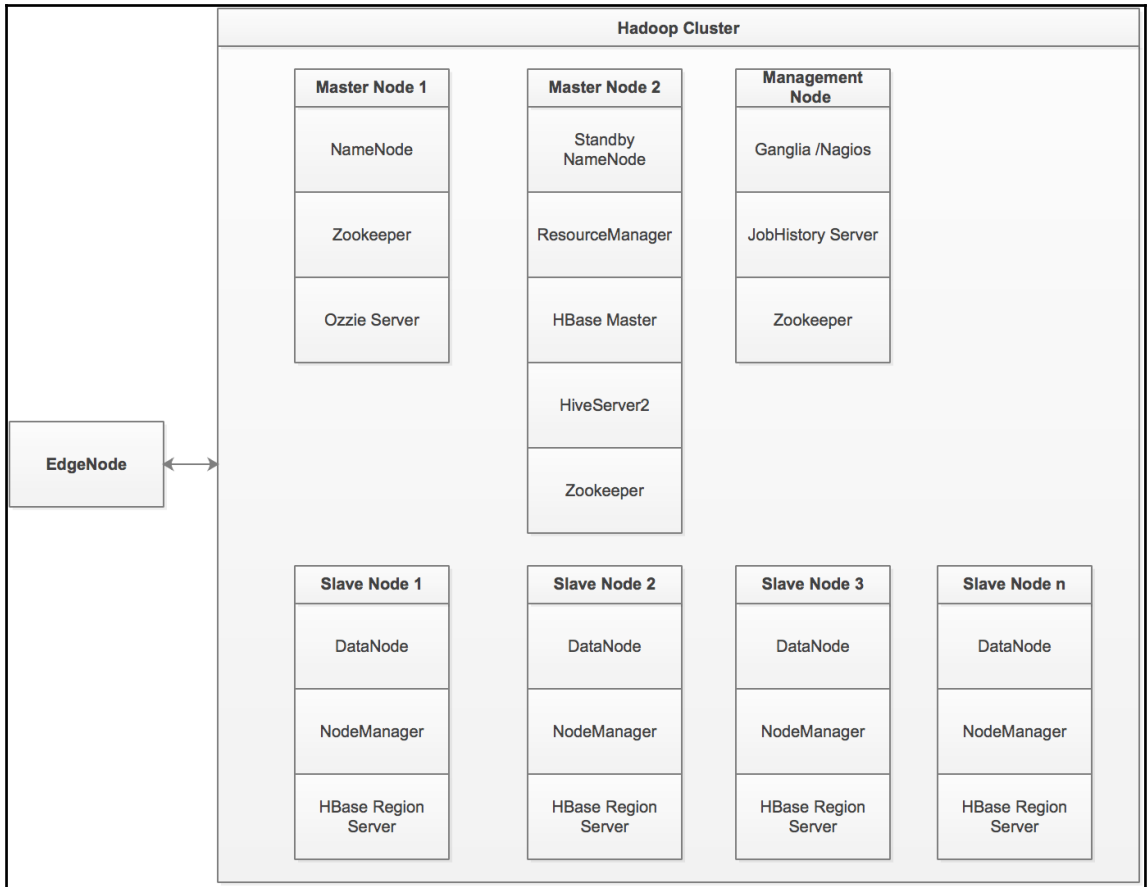
Chapter 3: Hadoop Design Consideration

| | |
|----------|---|
| Volume | Data size varies from terabytes to petabytes |
| Velocity | Data is generated every time tick, no limit |
| Variety | Data format – Structured, Semi-structured, Unstructured |
| Veracity | Data dependability, cleanliness |
| Value | Data conversion into value |

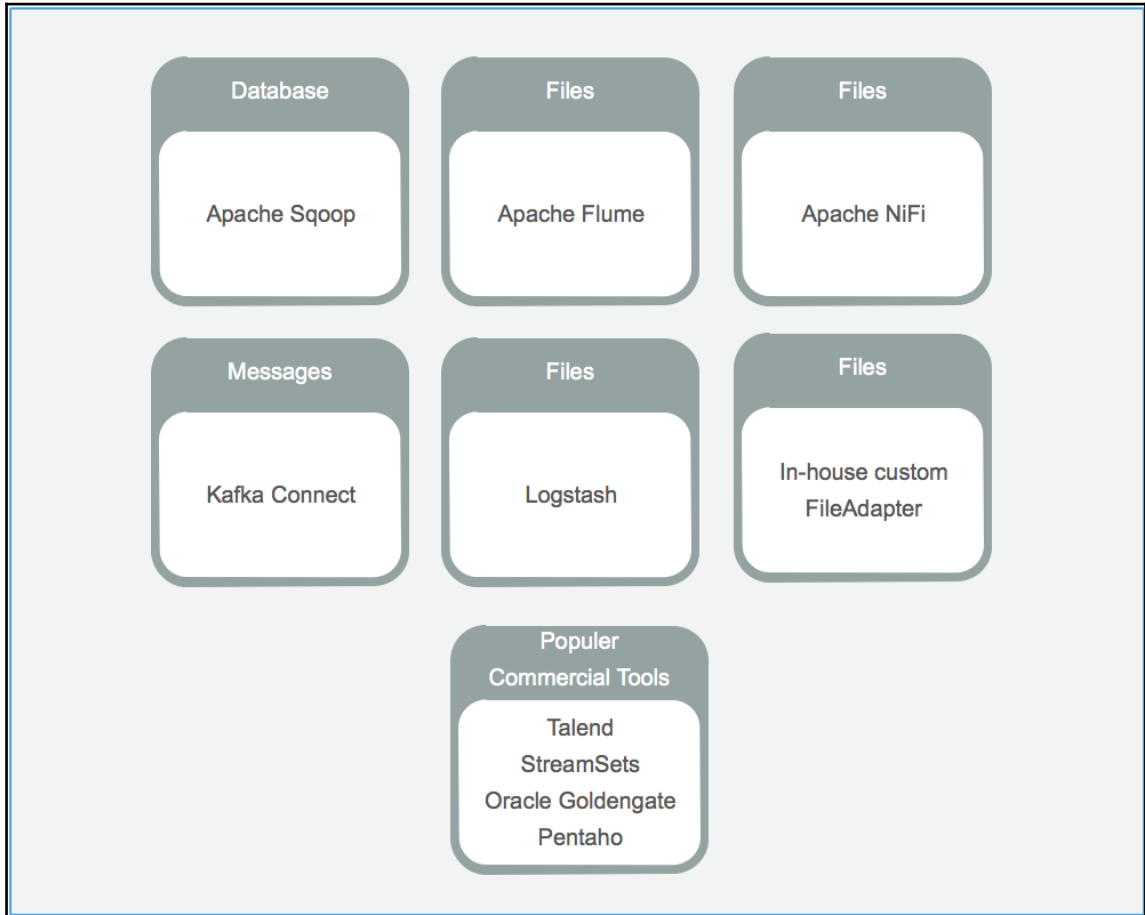


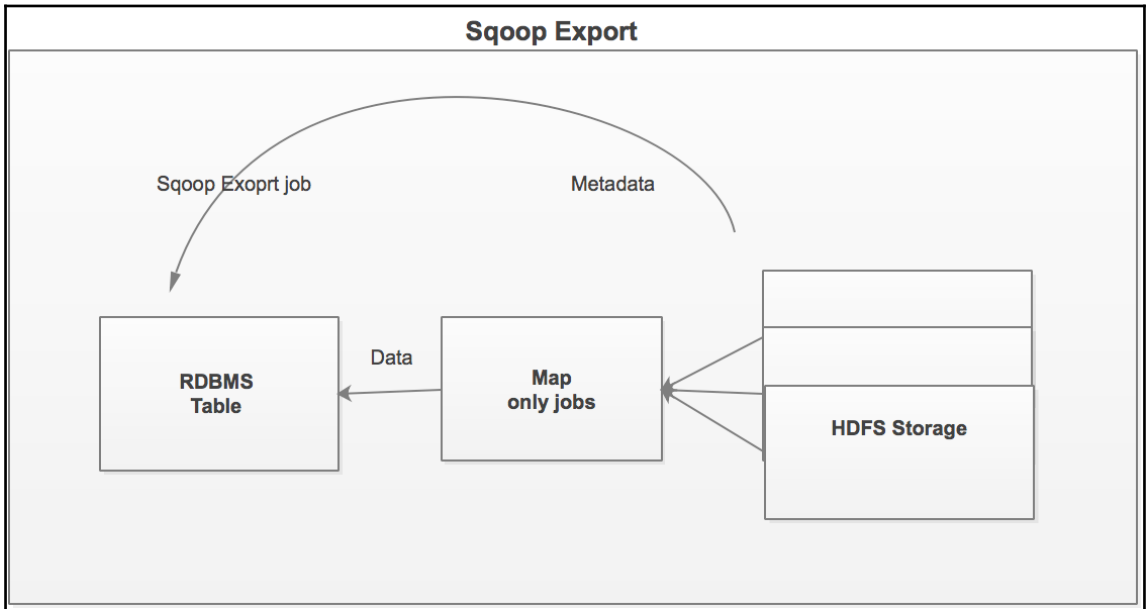
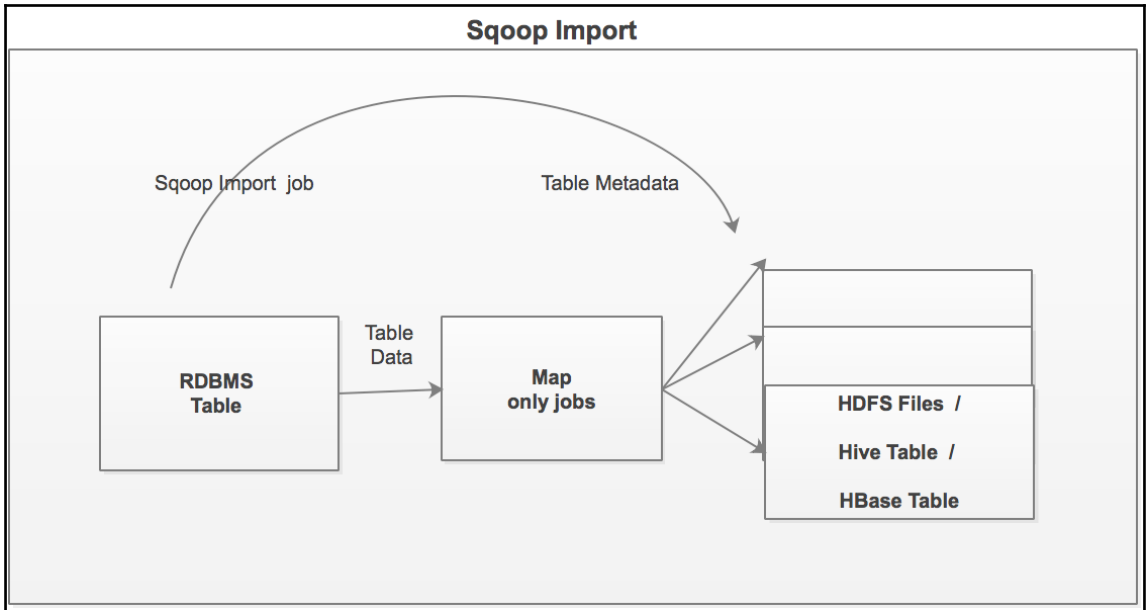


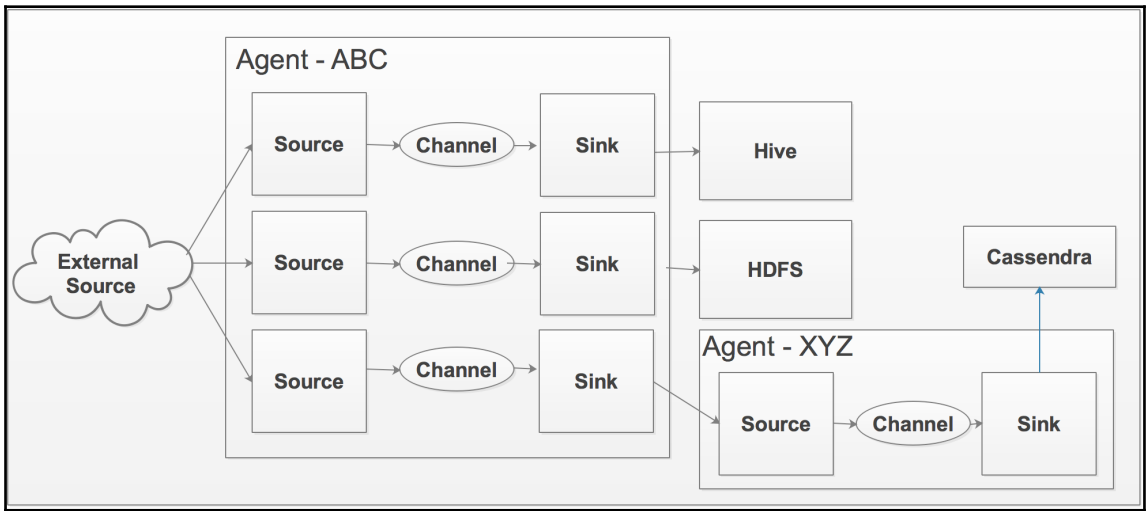
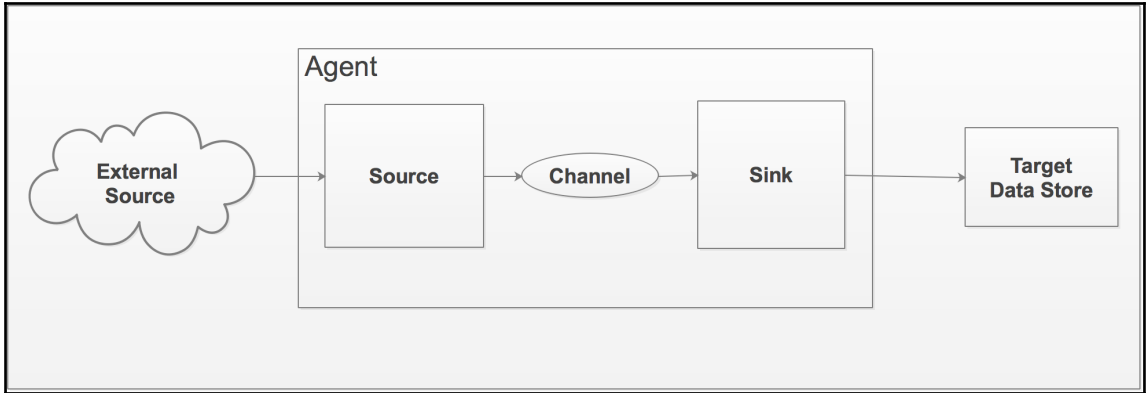


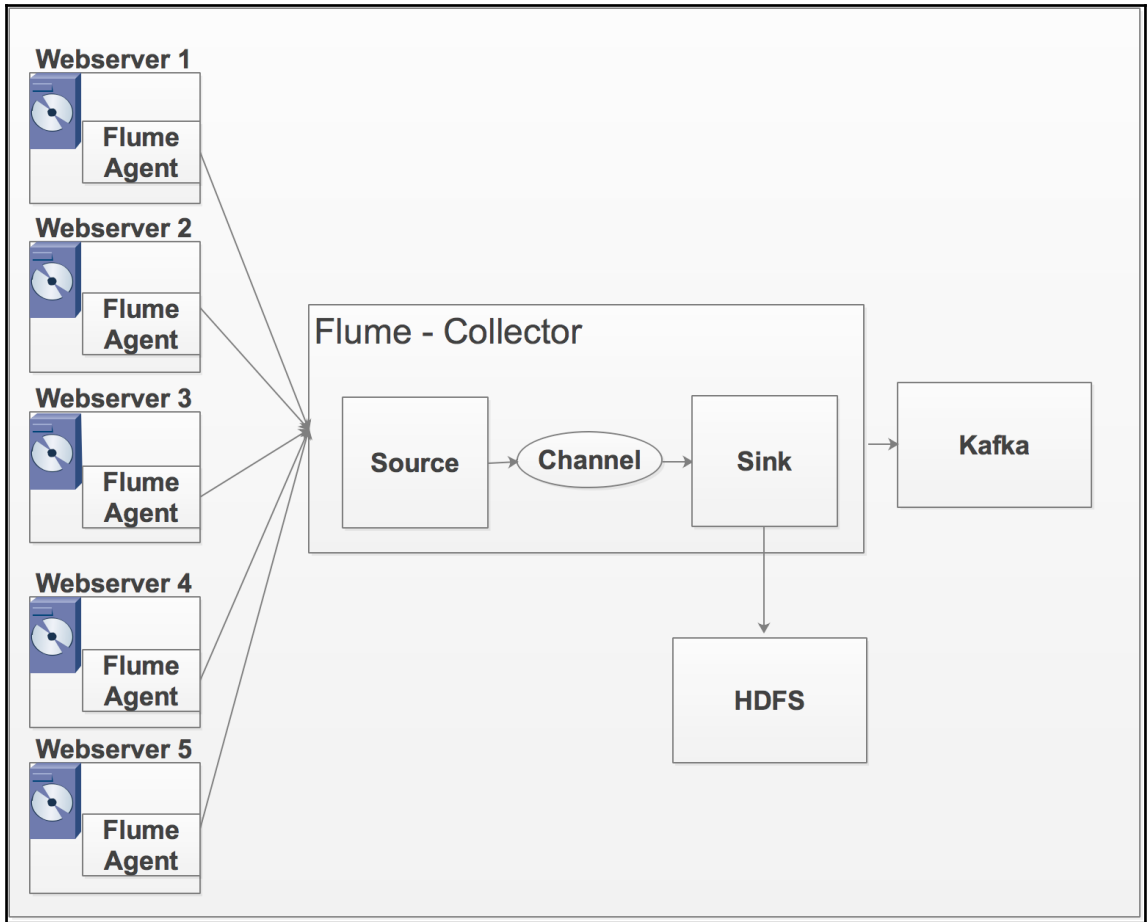


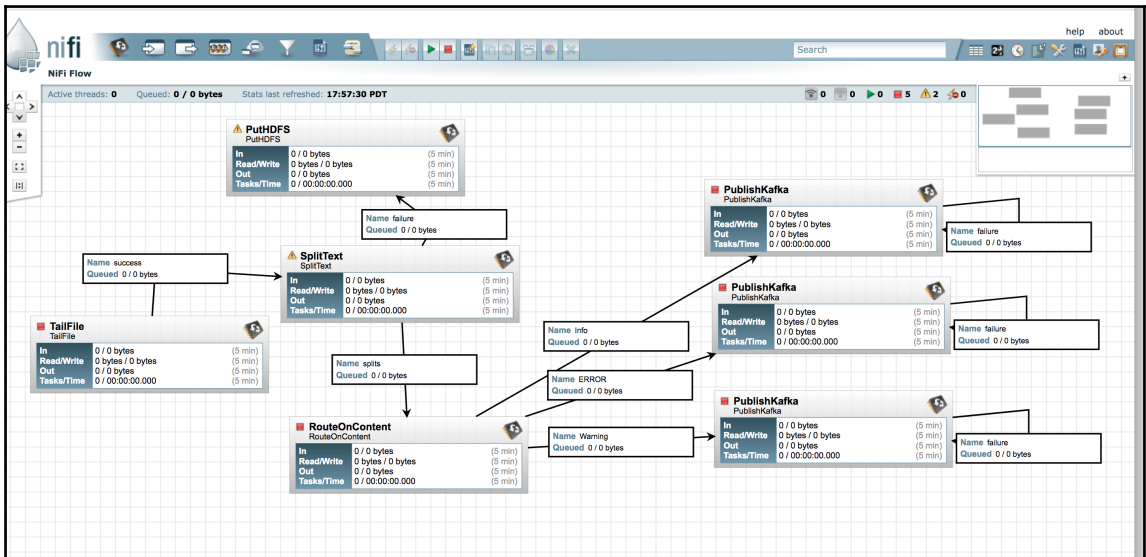
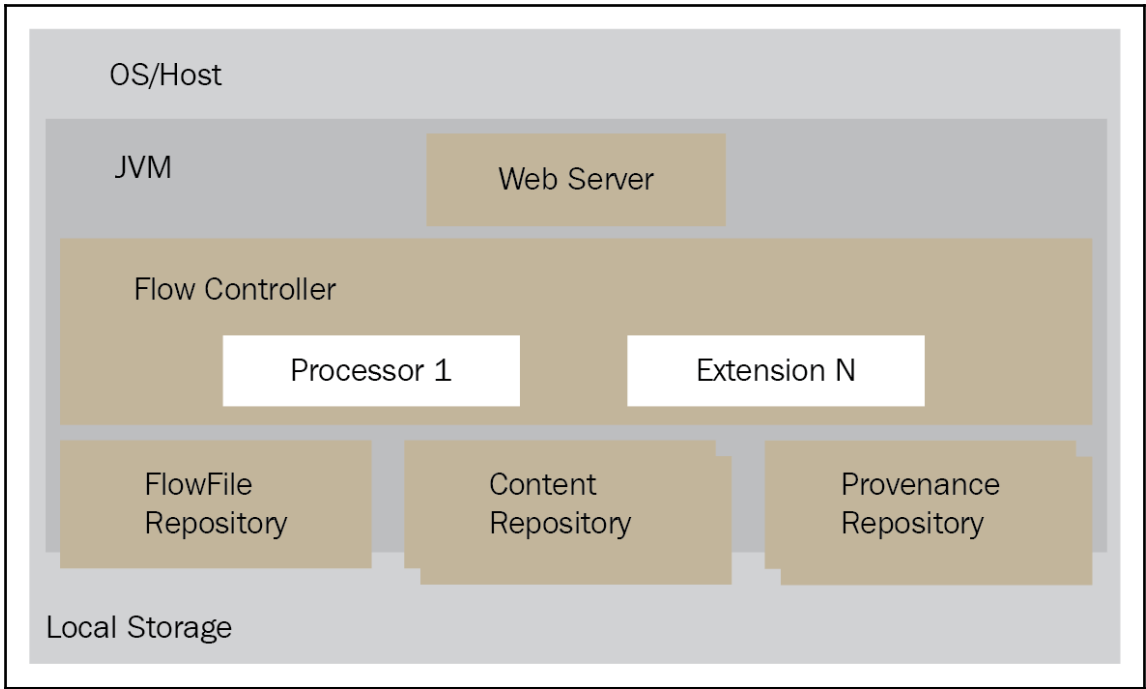
Chapter 4: Data Movement Techniques

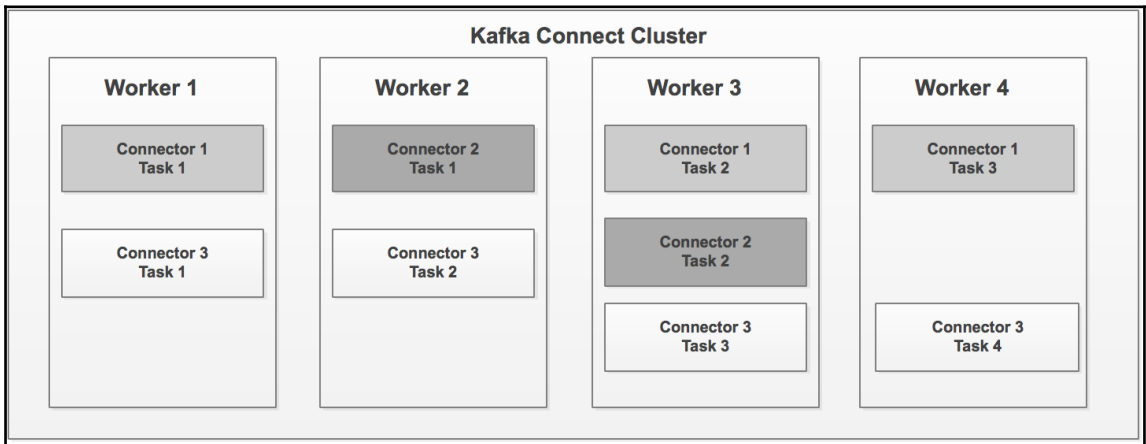
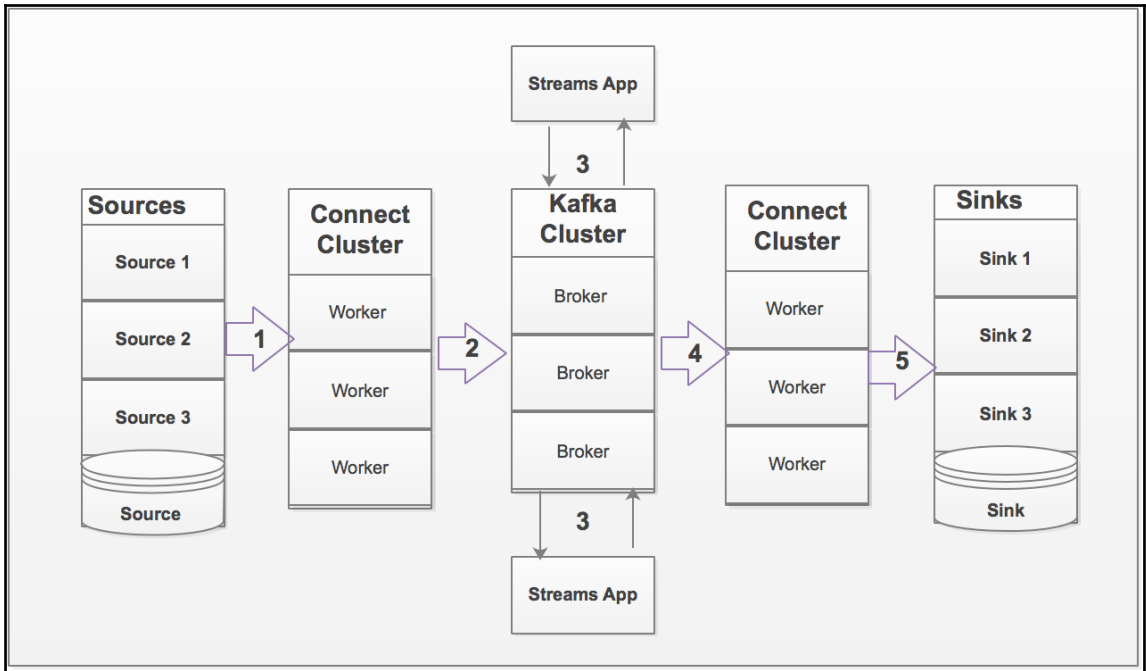


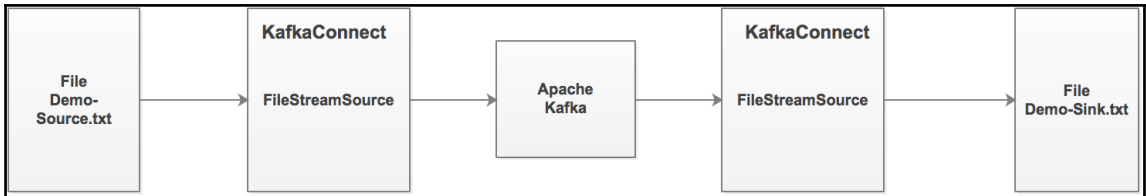
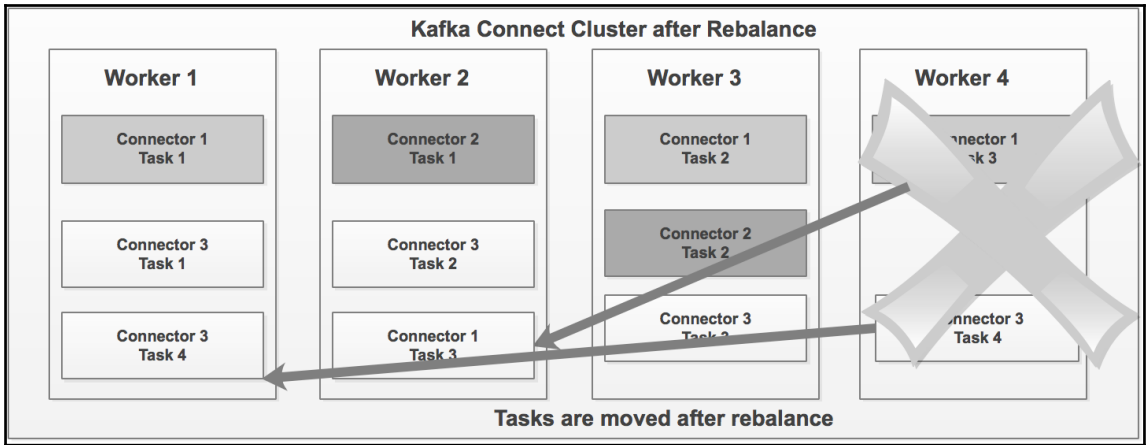




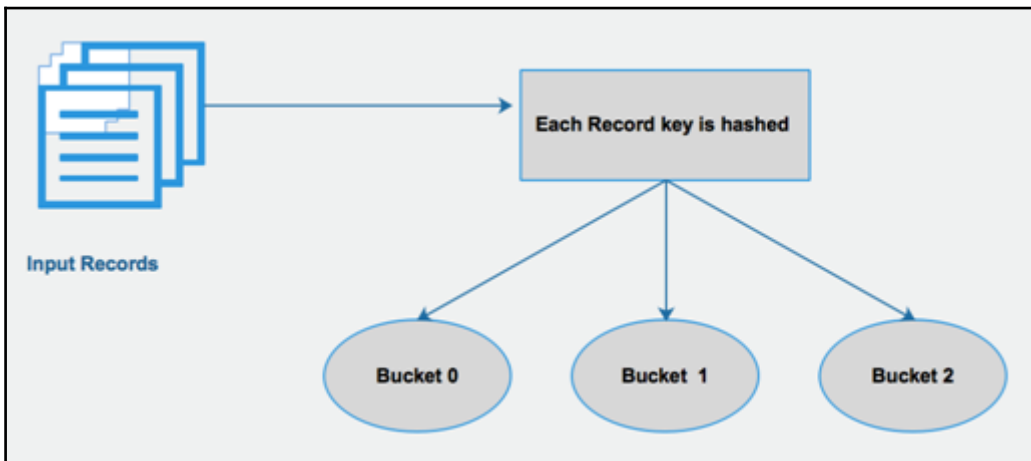
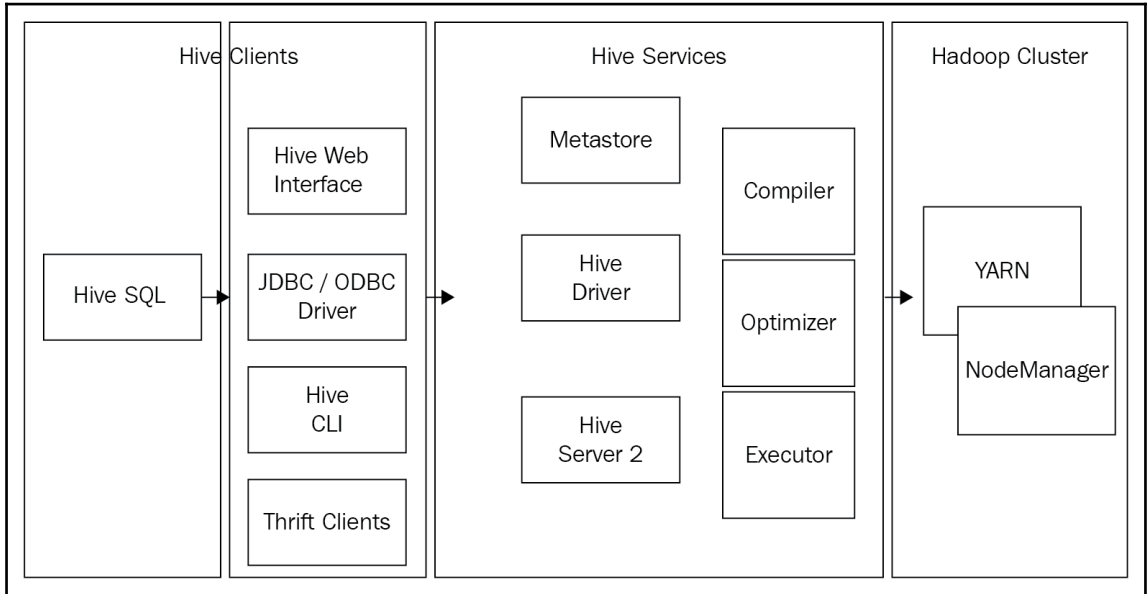


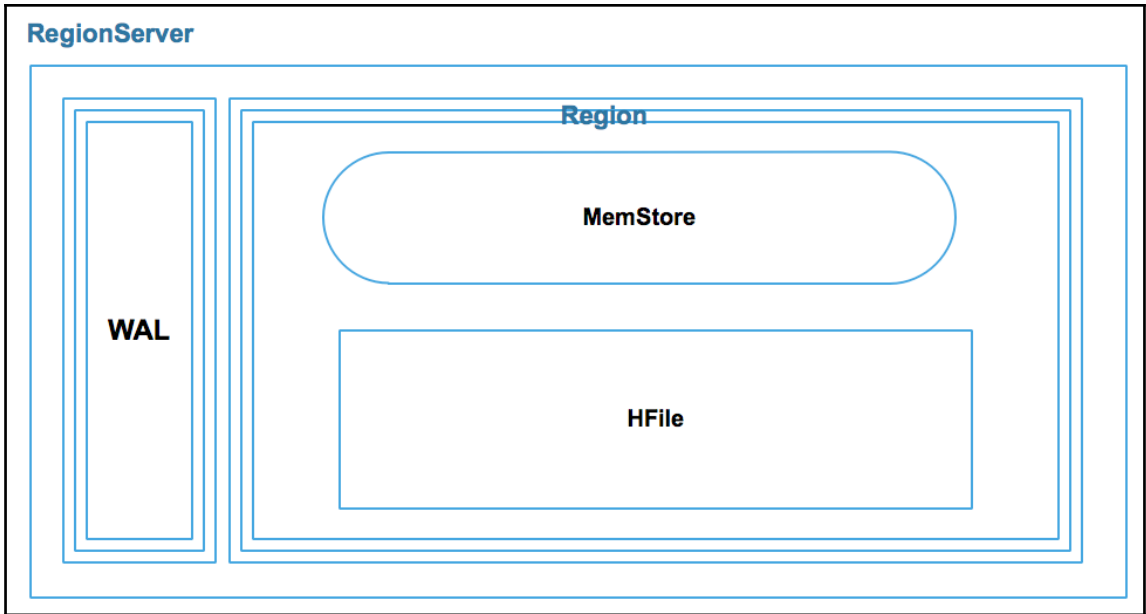
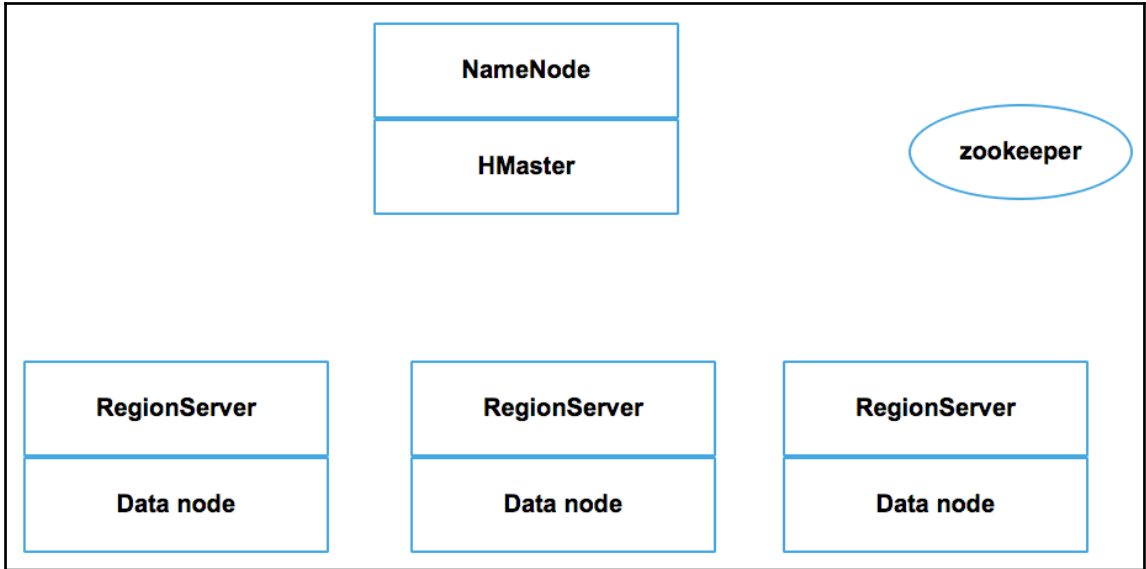




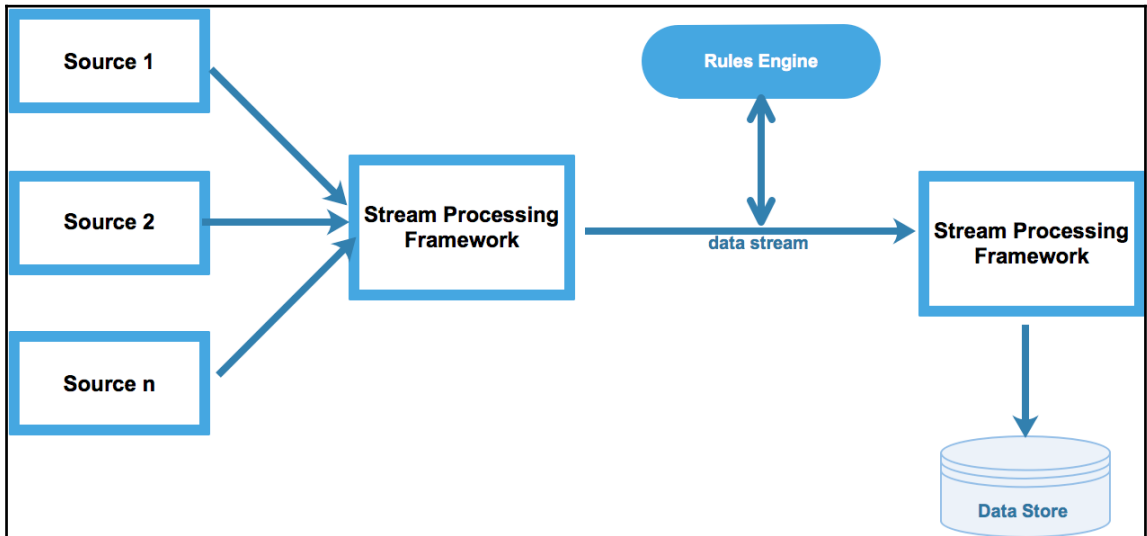
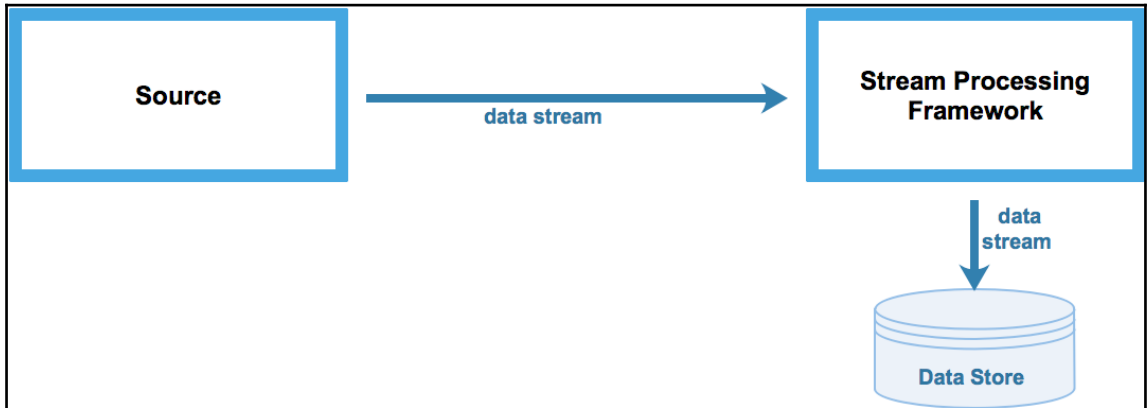


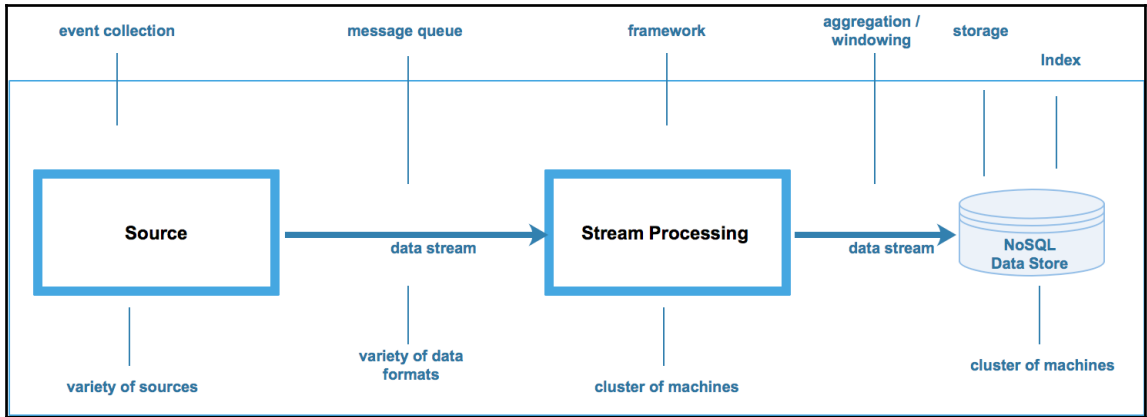
Chapter 5: Data Modeling in Hadoop

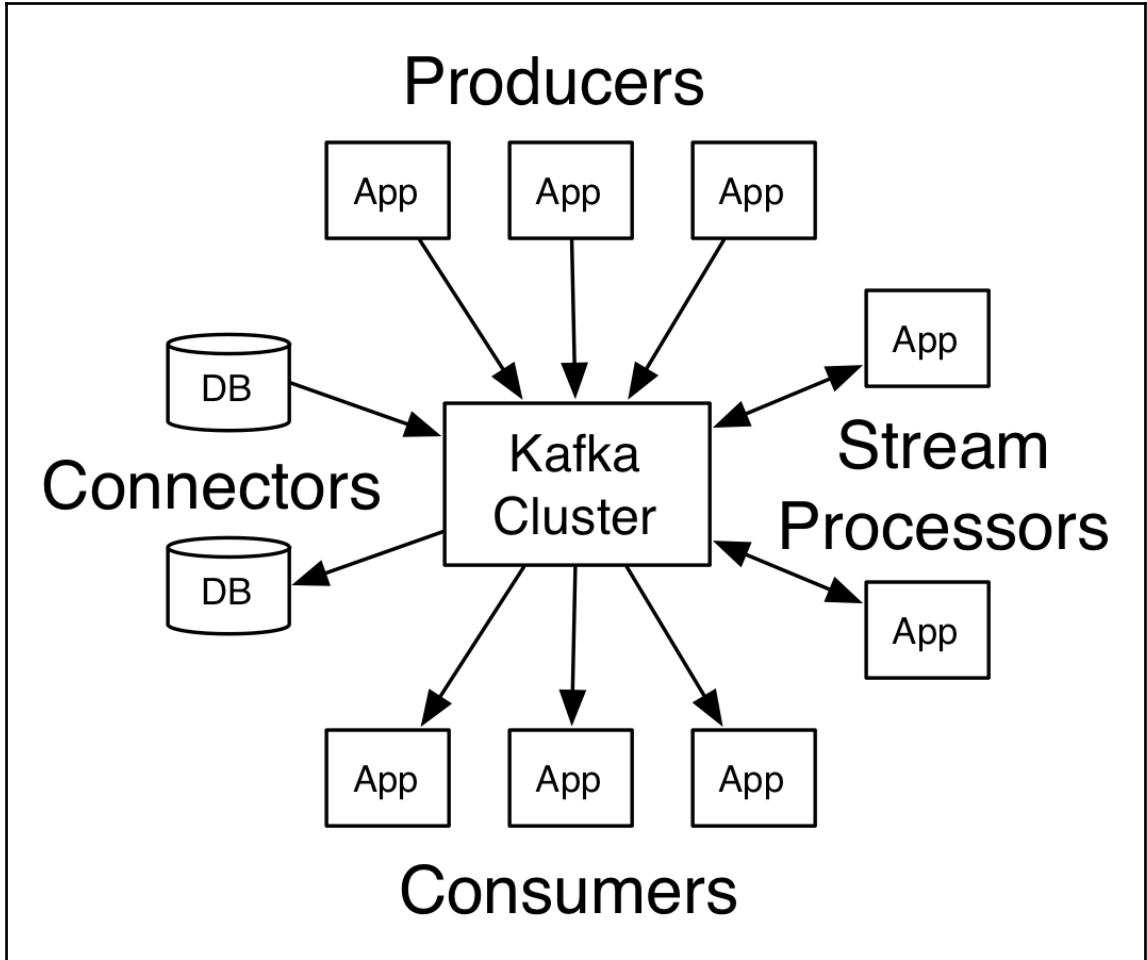


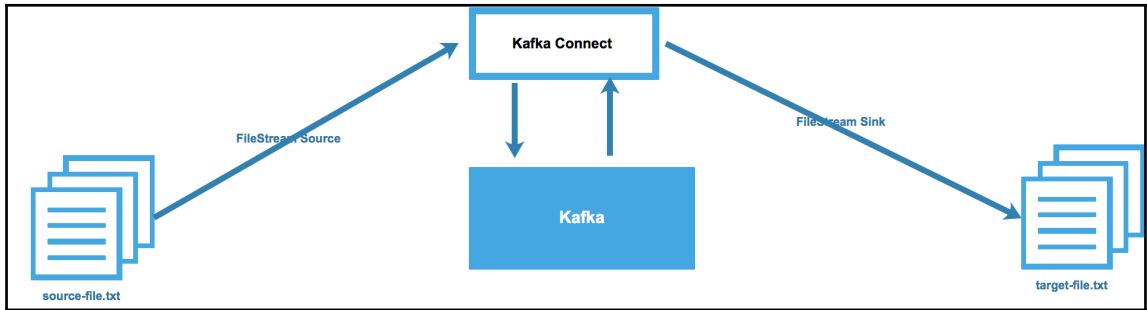
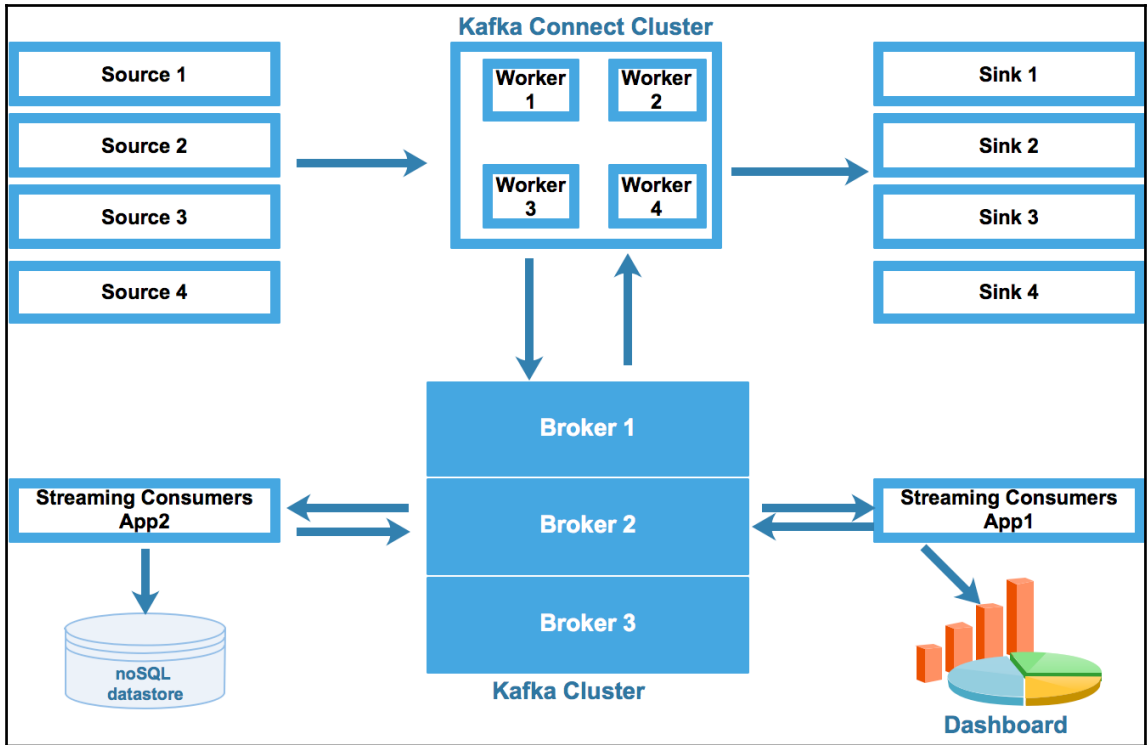


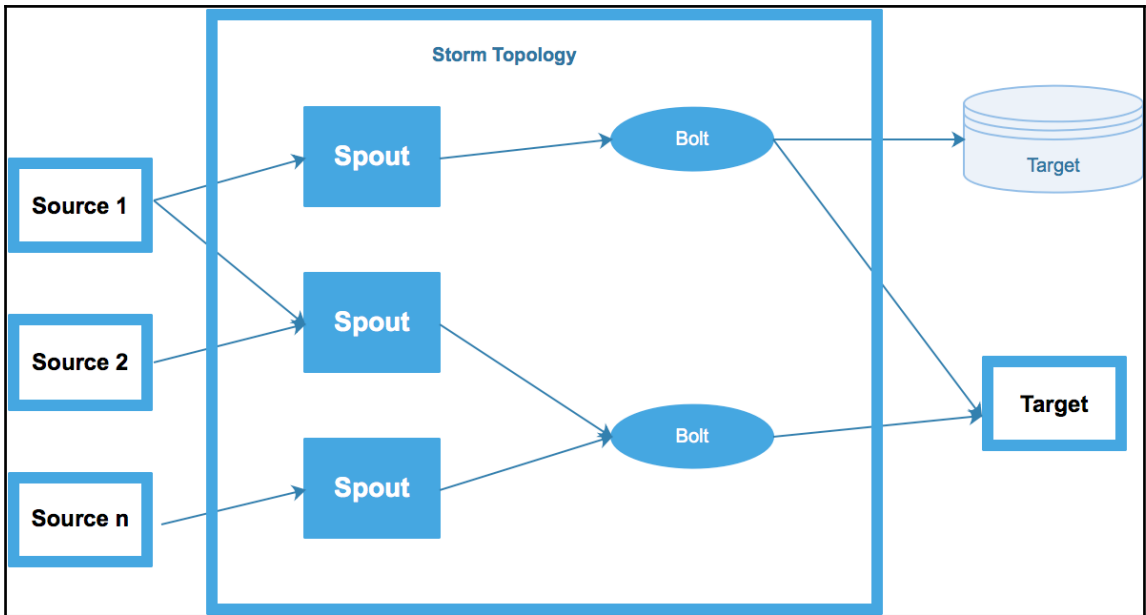
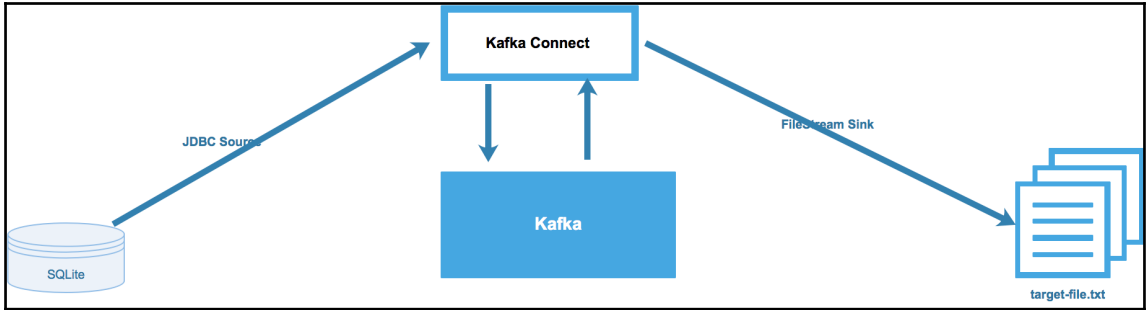
Chapter 6: Designing Real-Time Streaming Data Pipelines

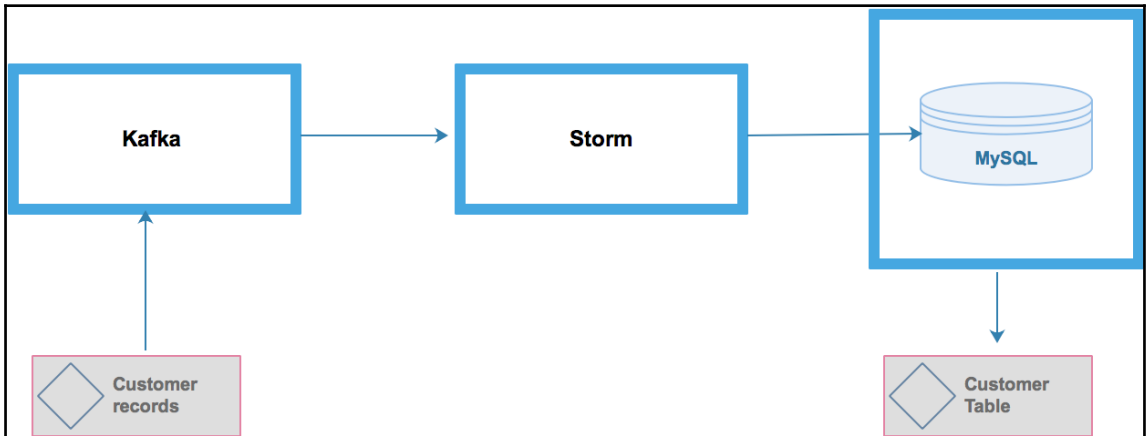
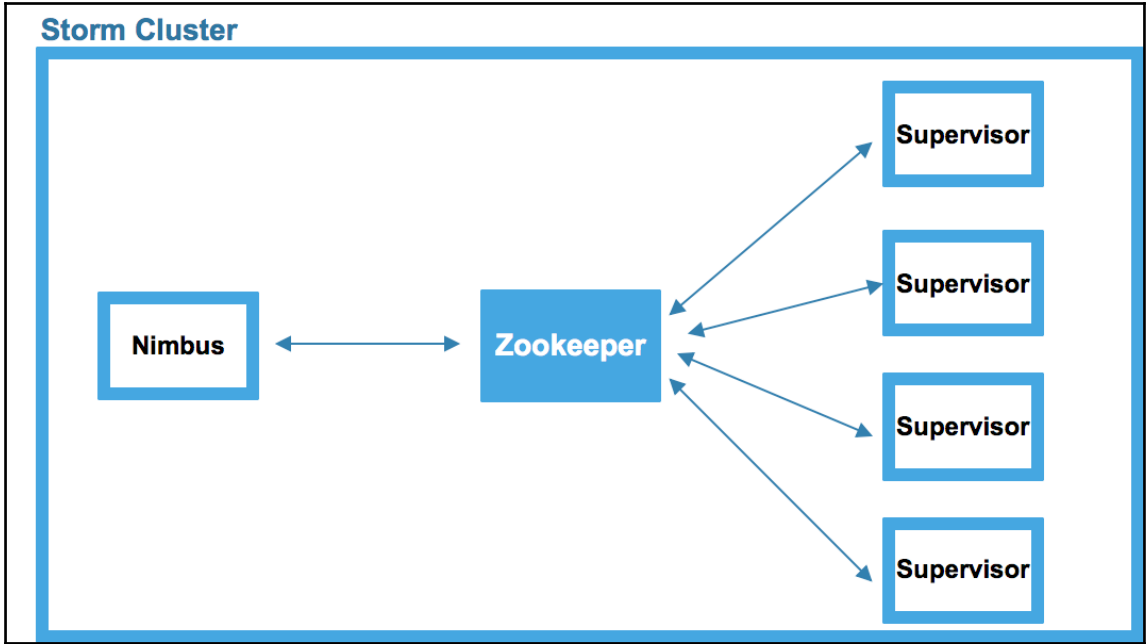


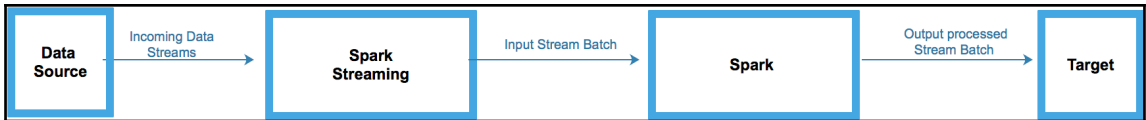
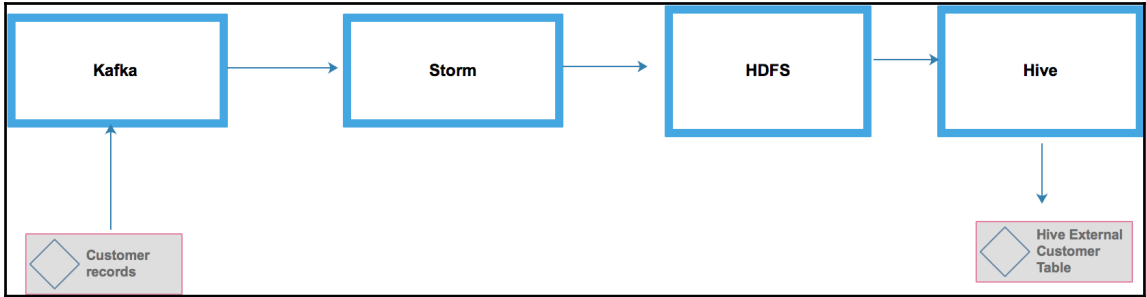




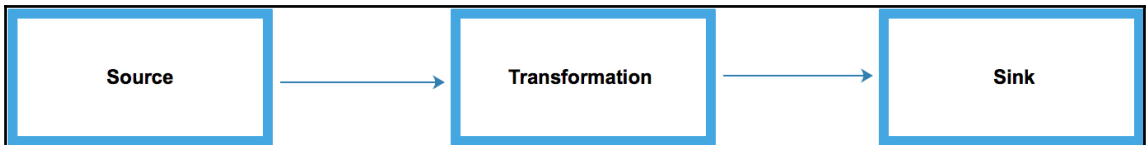




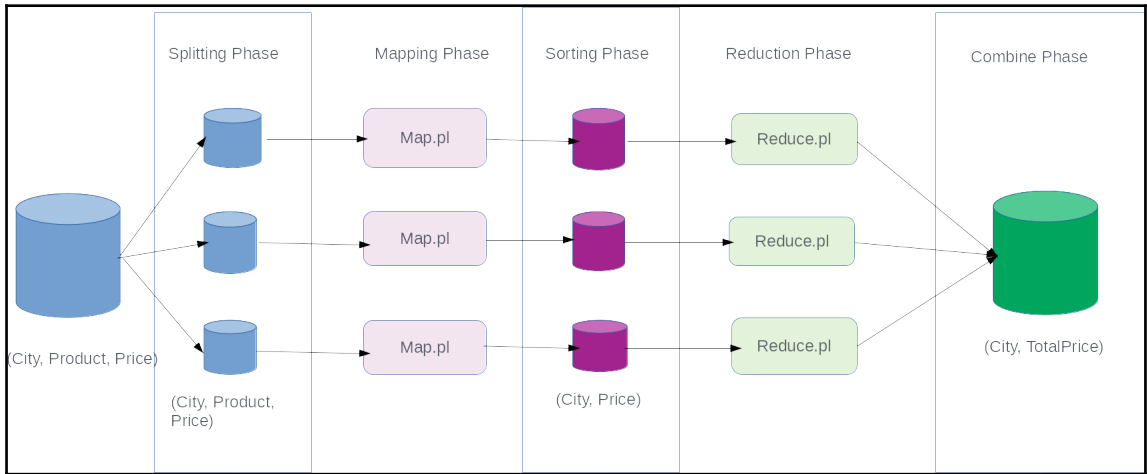
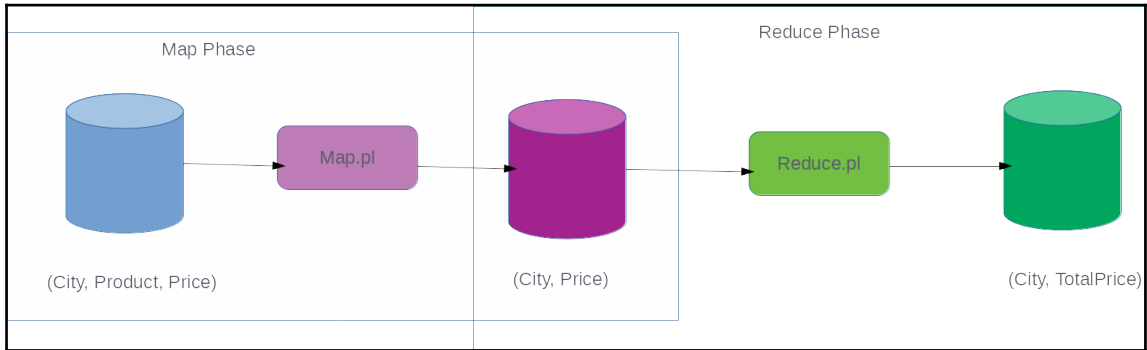




| | | | | | | | |
|-----------------------------|--|------------------------------------|--|--|----------------------------------|----------------------------|--|
| APIs & Libraries | CEP Event Processing | Table Relational | | FlinkML Machine Learning | Gelly Graph Processing | Table Relational | |
| | DataStream API Stream Processing | | | DataSet API Batch Processing | | | |
| Core | Runtime Distributed Streaming Dataflow | | | | | | |
| Deploy | Local Single JVM | Cluster Standalone, YARN | | | Cloud GCE, EC2 | | |



Chapter 7: Large-Scale Data Processing Frameworks



The screenshot displays the Ambari service management interface. On the left, a list of services is shown with green checkmarks indicating they are running. The services listed are HDFS, ZooKeeper, Ambari Metrics, SmartSense, Druid, and Superset. The SmartSense service is highlighted with a black background and a red notification badge containing the number '1'. Below the list is an 'Actions' dropdown menu. The dropdown menu is open, showing the following options: '+ Add Service', '▶ Start All', '■ Stop All', '⌂ Restart All Required', and '⬇ Download All Client Configs'. To the right of the service list, there are two 'Summary' panels, one above the other.

Add Service Wizard

| Service | Version | Description |
|--|-----------------|---|
| <input type="checkbox"/> Ranger | 0.11.0 | Comprehensive security for Hadoop |
| <input type="checkbox"/> Ranger KMS | 0.7.0 | Key Management Server |
| <input checked="" type="checkbox"/> SmartSense | 1.4.4.2.6.1.5-3 | SmartSense - Hortonworks SmartSense Tool (HST) helps quickly gather configuration, metrics, logs from common HDP services that aids to quickly troubleshoot support cases and receive cluster-specific recommendations. |
| <input type="checkbox"/> Spark | 1.6.3 | Apache Spark is a fast and general engine for large-scale data processing. |
| <input checked="" type="checkbox"/> Spark2 | 2.2.0 | Apache Spark is a fast and general engine for large-scale data processing |
| <input type="checkbox"/> Zeppelin Notebook | 0.7.3 | A web-based notebook that enables interactive data analytics. It enables you to make beautiful data-driven, interactive and collaborative documents with SQL, Scala and more. |
| <input checked="" type="checkbox"/> Druid | 0.10.1 | A fast column-oriented distributed data store. |
| <input type="checkbox"/> Mahout | 0.9.0 | Project of the Apache Software Foundation to produce free implementations of distributed or otherwise scalable machine learning algorithms focused primarily in the areas of collaborative filtering, clustering and classification |
| <input type="checkbox"/> Slider | 0.92.0 | A framework for deploying, managing and monitoring existing distributed applications on YARN. |
| <input checked="" type="checkbox"/> Superset | 0.15.0 | Superset is a data exploration platform designed to be visual, intuitive and interactive. |

[Next >>](#)

Add Service Wizard

- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Assign master components to hosts you want to run them on.
 * HiveServer2 and WebHCat Server will be hosted on the same host.

SNameNode:

NameNode:

App Timeline Server:

ResourceManager:

History Server:

Hive Metastore:

WebHCat Server: node-2

HiveServer2:

ZooKeeper Server:

ZooKeeper Server:

ZooKeeper Server:

Grafana:

node-1 (12.6 GB, 2 cores)

NameNode, ZooKeeper Server, Grafana, Activity Analyzer, HST Server, Activity Explorer, Spark2 History Server

node-2 (12.6 GB, 2 cores)

SNameNode, App Timeline Server, ResourceManager, History Server, Hive Metastore, WebHCat Server, HiveServer2, ZooKeeper Server

node-3 (12.6 GB, 2 cores)

ZooKeeper Server, Metrics Collector, Superset, Druid Router, Druid Broker, Druid Coordinator, Druid Overlord

Add Service Wizard

ADD SERVICE WIZARD

[Choose Services](#)

[Assign Masters](#)

Assign Slaves and Clients

[Customize Services](#)

[Configure Identities](#)

[Review](#)

[Install, Start and Test](#)

[Summary](#)

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on.
 Hosts that are assigned master components are shown with *.
 Client will install YARN Client, MapReduce2 Client, Tez Client, HCat Client, Hive Client, Pig Client, Spark2 Client and Slider Client.

| all | none | all | none | all | none | all | none | all | none | all | none |
|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|--------------------------|-------------------------------------|--------------------------|-------------------------------------|-------------------------------------|-------------------------------------|--------------------------|
| <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

← Show: 25 1 - 3 of 3 →

[← Back](#) [Next →](#)

Add Service Wizard

ADD SERVICE WIZARD

[Choose Services](#)

[Assign Masters](#)

Assign Slaves and Clients

[Customize Services](#)

[Configure Identities](#)

[Review](#)

[Install, Start and Test](#)

[Summary](#)

Settings **Advanced**

Hive Metastore

Hive Metastore host: node-2

Hive Database:

- New MySQL Database
- Existing MySQL / MariaDB Database
- Existing PostgreSQL Database
- Existing Oracle Database

To use MySQL with Hive, you must download the [MySQL Connector/J JDBC Driver from MySQL](#). Once downloaded to the Ambari Server host, run:
ambari-server setup --jdbc-db=mysql --jdbc-driver=/path/to/mysql/mysql-connector-java.jar

Database Name: 🔒 🔄

Database Username: 🔒 🔄

Database Password: 🔒

JDBC Driver Class: 🔒 🔄

Database URL: 🔒 🔄

[Test Connection](#) Connection OK ✔

Add Service Wizard

- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Please review the configuration before installation

WARNING (HDP-UTILS-1.1.0.22):
<http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.22/repos/ubuntu16>

Services:

YARN + MapReduce2
 App Timeline Server : node-2.c.coastal-airlock-197705.internal
 NodeManager : 1 host
 ResourceManager : node-2.c.coastal-airlock-197705.internal

Tez
 Clients : 3 hosts

Hive
 Metastore : node-2.c.coastal-airlock-197705.internal
 HiveServer2 : node-2.c.coastal-airlock-197705.internal
 WebHCat Server : node-2.c.coastal-airlock-197705.internal
 Database : Existing MySQL / MariaDB Database

Pig
 Clients : 3 hosts

Spark2
 Livy for Spark2 Server : 1 host
 History Server : node-1.c.coastal-airlock-197705.internal
 Thrift Server : 1 host

Slider
 Clients : 3 hosts

← Back
Print
Deploy →

Add Service Wizard

- ADD SERVICE WIZARD
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

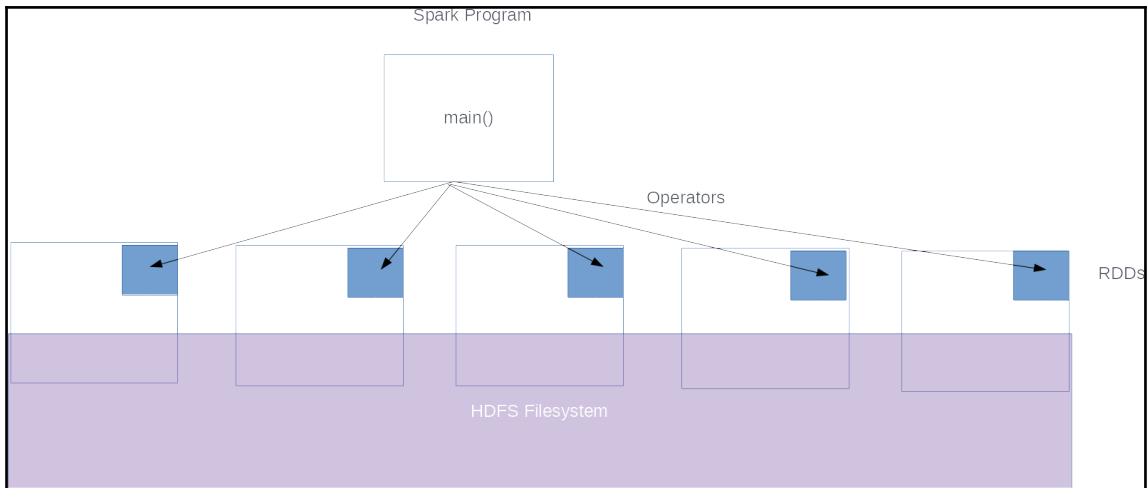
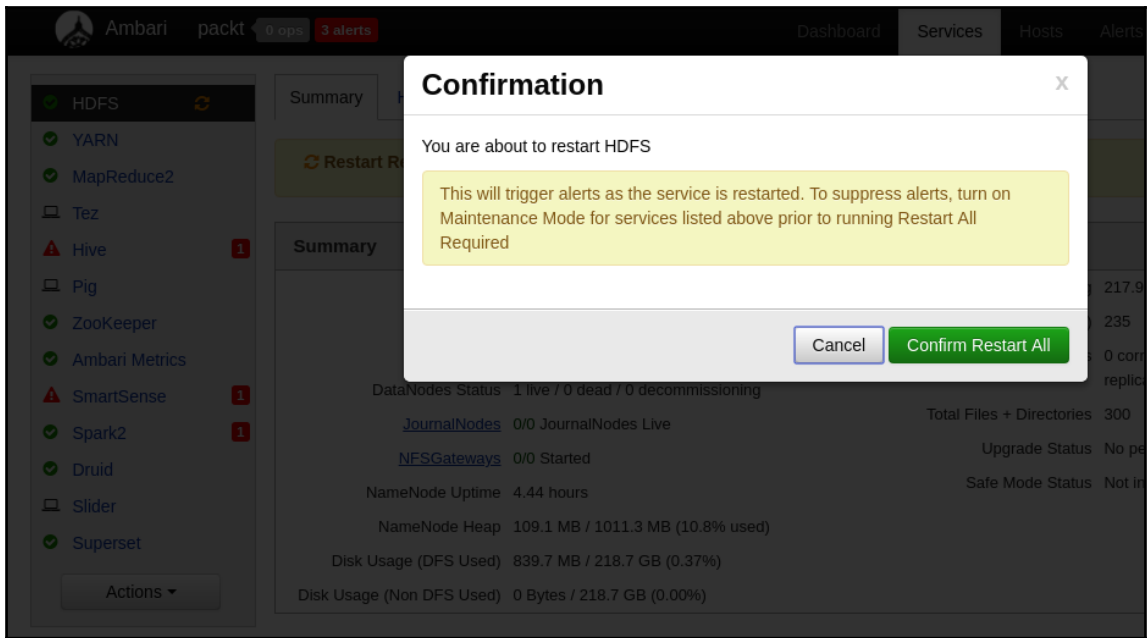
Summary

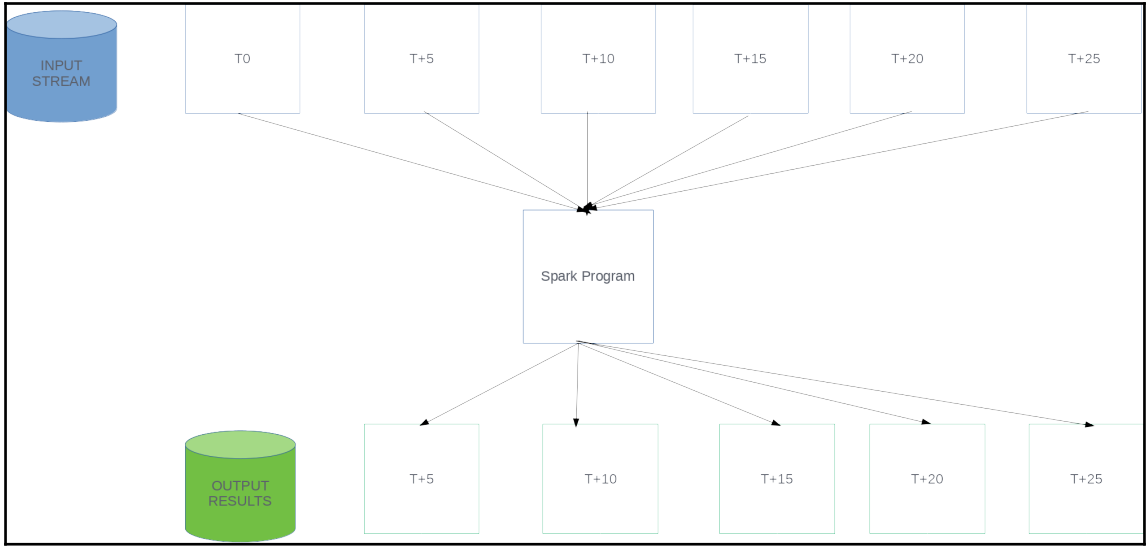
Important: You may also need to restart other services for the newly added services to function properly (for example, HDFS and YARN/MapReduce need to be restarted after adding Oozie). After closing this wizard, please restart all services that have the restart indicator 🔄 next to the service name.

Here is the summary of the install process.

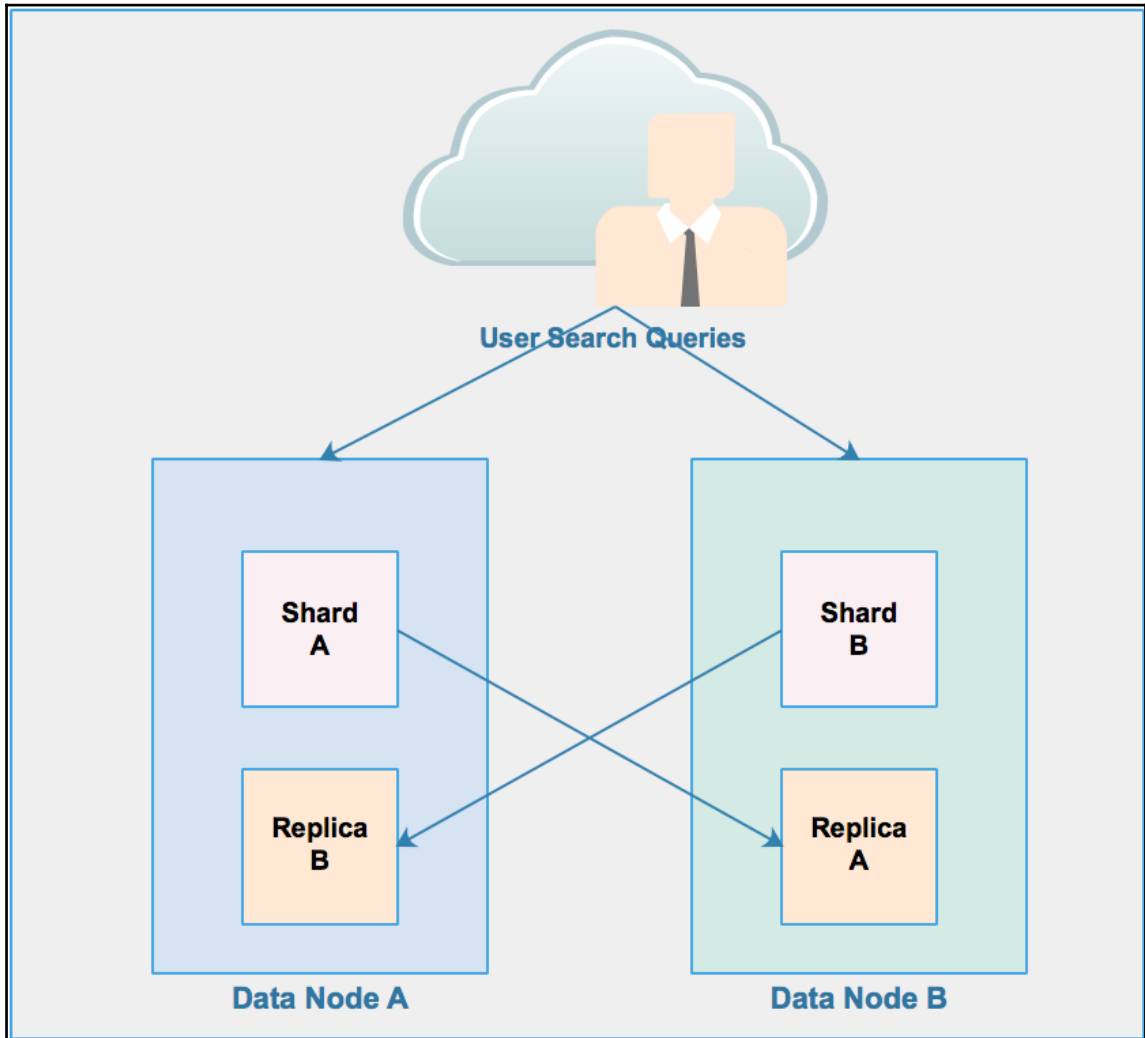
The cluster consists of 3 hosts
 3 warnings

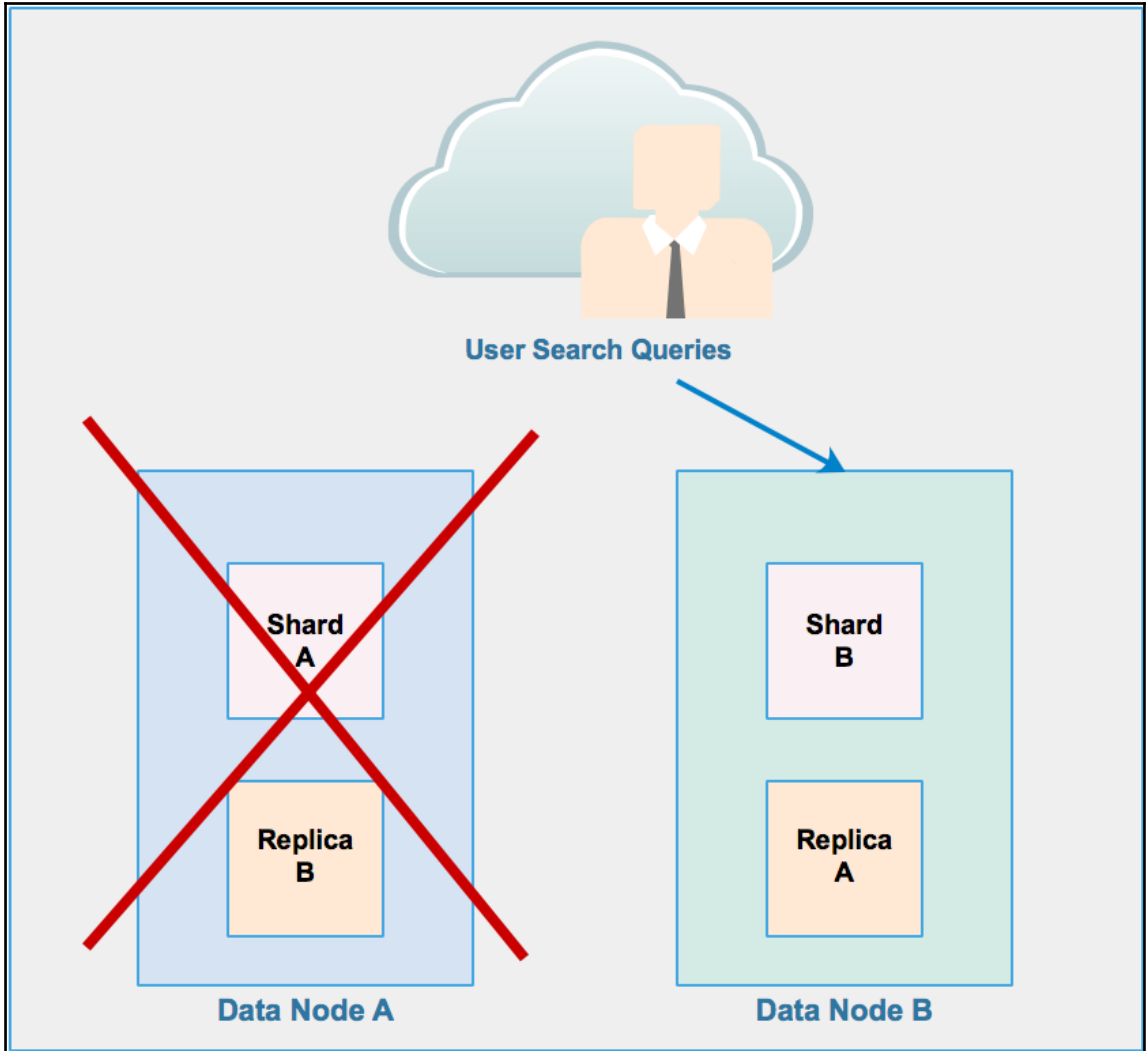
Complete →

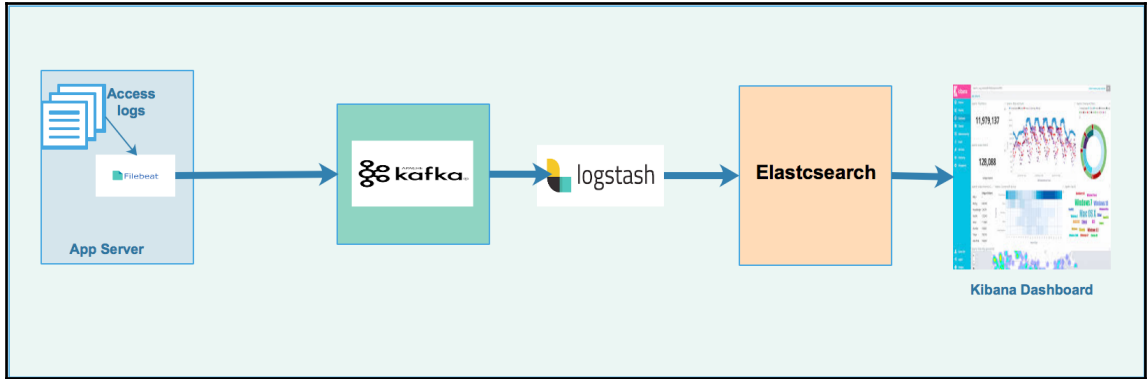




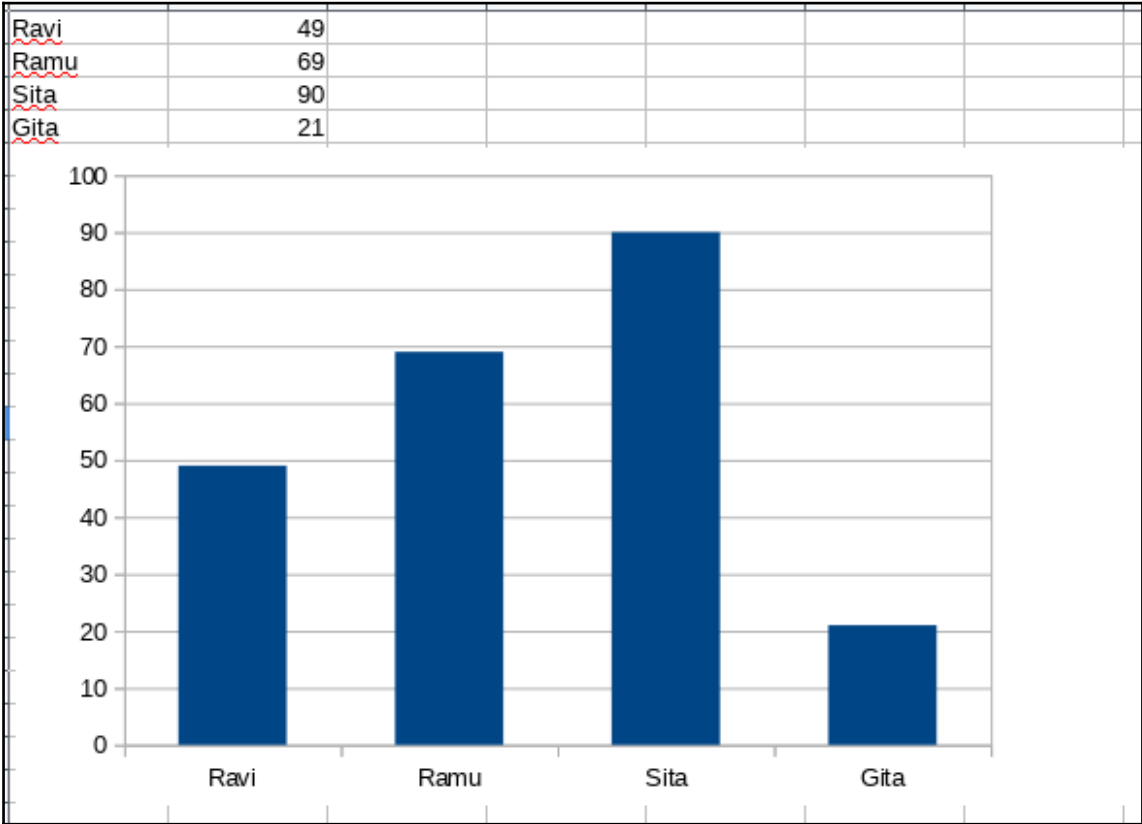
Chapter 8: Building Enterprise Search Platform

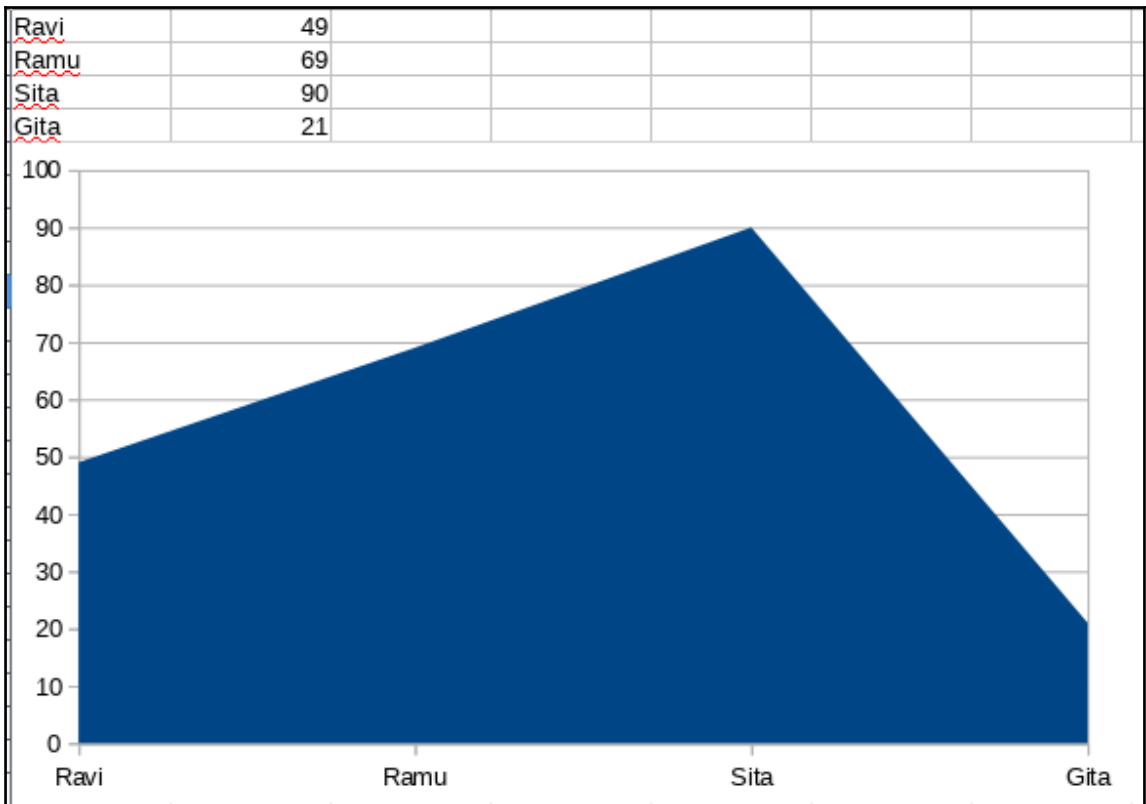


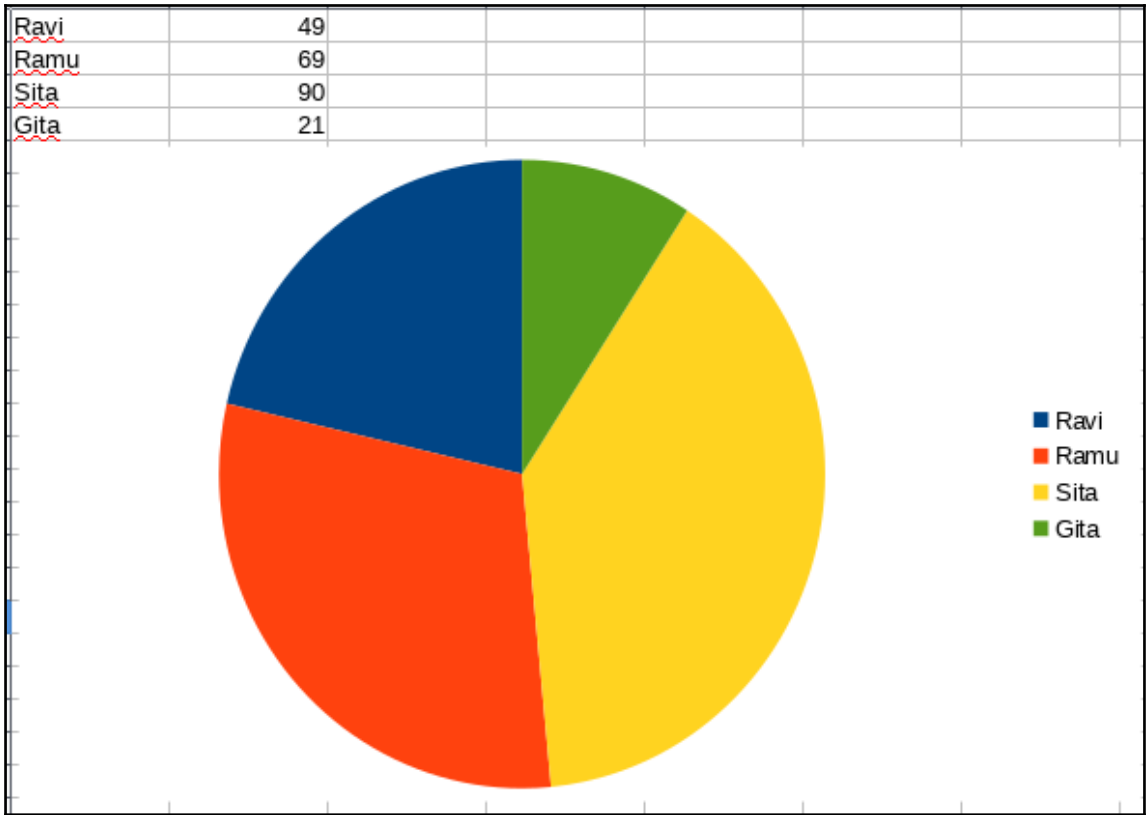




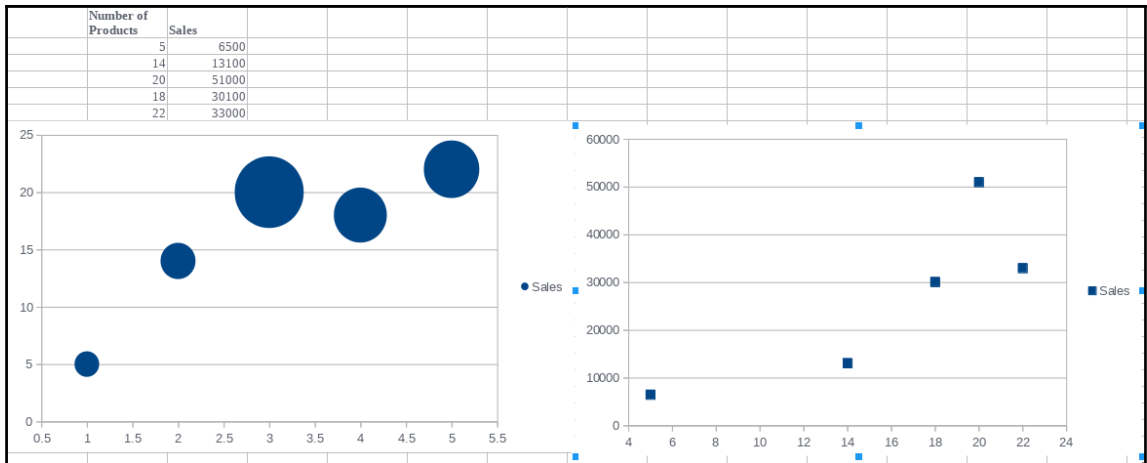
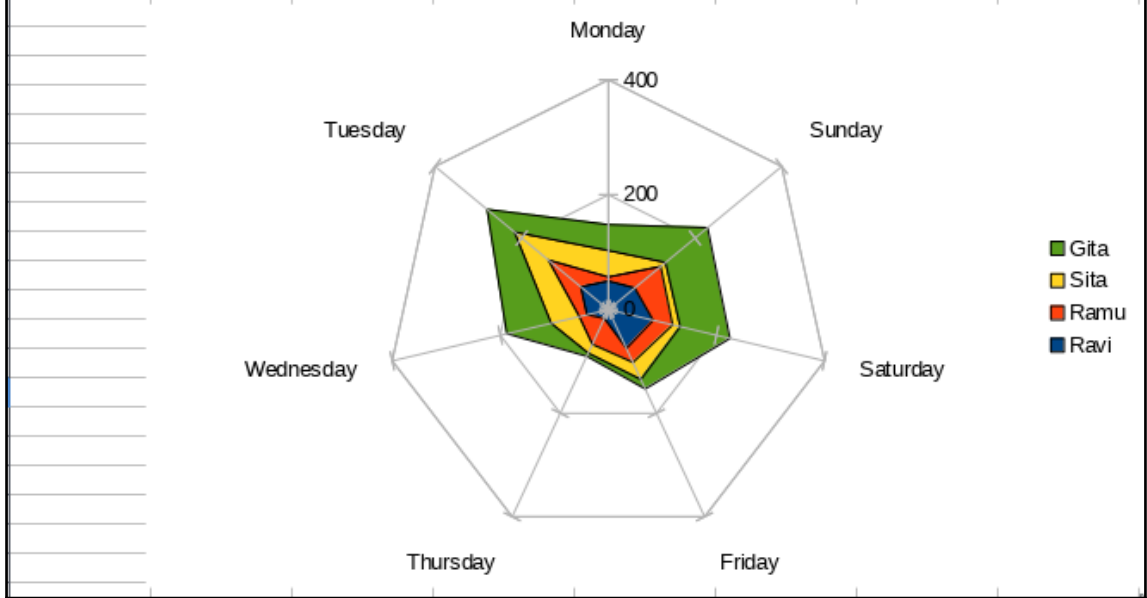
Chapter 9: Designing Data Visualization Solutions

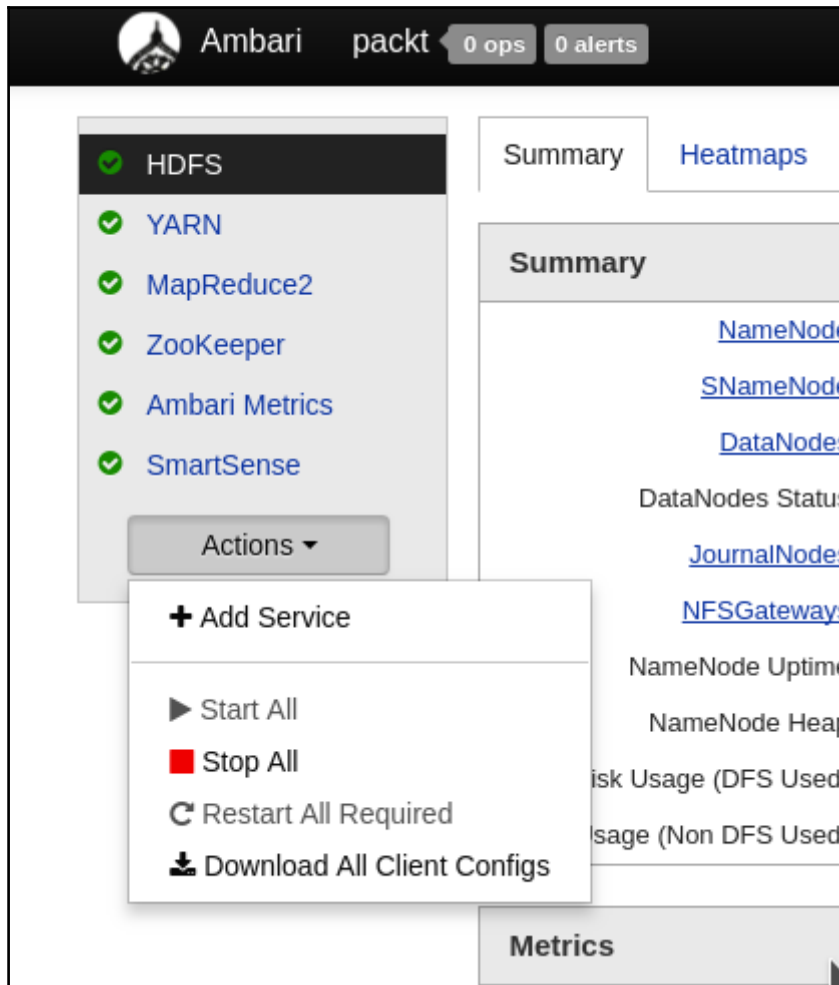






| | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday | Sunday |
|-------------|--------|---------|-----------|----------|--------|----------|--------|
| <u>Ravi</u> | 50 | 64 | 40 | 17 | 75 | 83 | 61 |
| <u>Ramu</u> | 7 | 75 | 14 | 50 | 27 | 36 | 60 |
| <u>Sita</u> | 46 | 78 | 52 | 17 | 33 | 14 | 11 |
| <u>Gita</u> | 45 | 63 | 82 | 6 | 18 | 92 | 97 |





Add Service Wizard

| | | | |
|-------------------------------------|-------------------|-----------------|---|
| <input type="checkbox"/> | Ranger KMS | 0.7.0 | Key Management Server |
| <input checked="" type="checkbox"/> | SmartSense | 1.4.4.2.6.1.5-3 | SmartSense - Hortonworks SmartSense Tool (HST) helps quickly gather configuration, metrics, logs from common HDP services that aids to quickly troubleshoot support cases and receive cluster-specific recommendations. |
| <input type="checkbox"/> | Spark | 1.6.3 | Apache Spark is a fast and general engine for large-scale data processing. |
| <input type="checkbox"/> | Spark2 | 2.2.0 | Apache Spark is a fast and general engine for large-scale data processing |
| <input type="checkbox"/> | Zeppelin Notebook | 0.7.3 | A web-based notebook that enables interactive data analytics. It enables you to make beautiful data-driven, interactive and collaborative documents with SQL, Scala and more. |
| <input checked="" type="checkbox"/> | Druid | 0.10.1 | A fast column-oriented distributed data store. |
| <input type="checkbox"/> | Mahout | 0.9.0 | Project of the Apache Software Foundation to produce free implementations of distributed or otherwise scalable machine learning algorithms focused primarily in the areas of collaborative filtering, clustering and classification |
| <input type="checkbox"/> | Slider | 0.92.0 | A framework for deploying, managing and monitoring existing distributed applications on YARN. |
| <input checked="" type="checkbox"/> | Superset | 0.15.0 | Superset is a data exploration platform designed to be visual, intuitive and interactive. |

Next →

Add Service Wizard

Metrics Collector:

Grafana:

HST Server:

Activity Explorer:

Activity Analyzer:

Druid Broker: +

Druid Overlord: +

Superset: +

Druid Coordinator: +

Druid Router: +

node-3 (12.6 GB, 2 cores)

- ZooKeeper Server
- Metrics Collector
- Druid Broker
- Druid Overlord
- Superset
- Druid Coordinator
- Druid Router

← Back

Next →

Add Service Wizard

ADD SERVICE WIZARD

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on. Hosts that are assigned master components are shown with *.

| Host | | all none | all none | all none | all none |
|--------|---|--|-------------------------------------|--|---|
| node-1 | * | <input type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input type="checkbox"/> Druid Historical | <input type="checkbox"/> Druid MiddleManager |
| node-2 | * | <input type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input type="checkbox"/> Druid Historical | <input type="checkbox"/> Druid MiddleManager |
| node-3 | * | <input checked="" type="checkbox"/> DataNode | <input type="checkbox"/> NFSGateway | <input checked="" type="checkbox"/> Druid Historical | <input checked="" type="checkbox"/> Druid MiddleManager |

Show: 25 1 - 3 of 3 ⏪ ⏩ ⏴ ⏵

← Back
Next →

Add Service Wizard

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Customize Services

We have come up with recommended configurations for the services you selected. Customize them as you see fit.

HDFS ZooKeeper Ambari Metrics SmartSense **Druid** **Superset** 4 Misc

There are 4 configuration changes in 2 services [Show Details](#)

Group: Default (3) [Manage Config Groups](#) Filter...

SUPERSET META DATA STORAGE CONFIG 4 Advanced

SUPERSET META DATA STORAGE

Superset Database name

Superset Database type

MYSQL

Add Service Wizard

ADD SERVICE WIZARD

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Install, Start and Test

Please wait while the selected services are installed and started.

23 % overall

Show: **All (3)** | [In Progress \(3\)](#) | [Warning \(0\)](#) | [Success \(0\)](#) | [Fail \(0\)](#)

| Host | Status | Message |
|--------|--|-------------------------------------|
| node-1 | <div style="width: 33%;"><div style="width: 33%;"></div></div> 33% | Install complete (Waiting to start) |
| node-2 | <div style="width: 33%;"><div style="width: 33%;"></div></div> 33% | Install complete (Waiting to start) |
| node-3 | <div style="width: 5%;"><div style="width: 5%;"></div></div> 5% | Installing Druid Broker |

3 of 3 hosts showing - [Show All](#)
Show: 25 | 1 - 3 of 3

Next -->

Add Service Wizard

ADD SERVICE WIZARD

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Summary

Important: You may also need to restart other services for the newly added services to function properly (for example, HDFS and YARN/MapReduce need to be restarted after adding Oozie). After closing this wizard, please restart all services that have the restart indicator 🔄 next to the service name.

Here is the summary of the install process.

The cluster consists of 3 hosts
 Installed and started services successfully on 2 new hosts
 1 warnings

Complete -->




Coordinator Console

Running Tasks

Show entries

| id | createdTime | queueInsertionTime |
|--|--------------------------|--------------------------|
| index_hadoop_wikiticker_2018-03-16T04:54:38.979Z | 2018-03-16T04:54:39.058Z | 2018-03-16T04:54:39.082Z |

Showing 1 to 1 of 1 entries

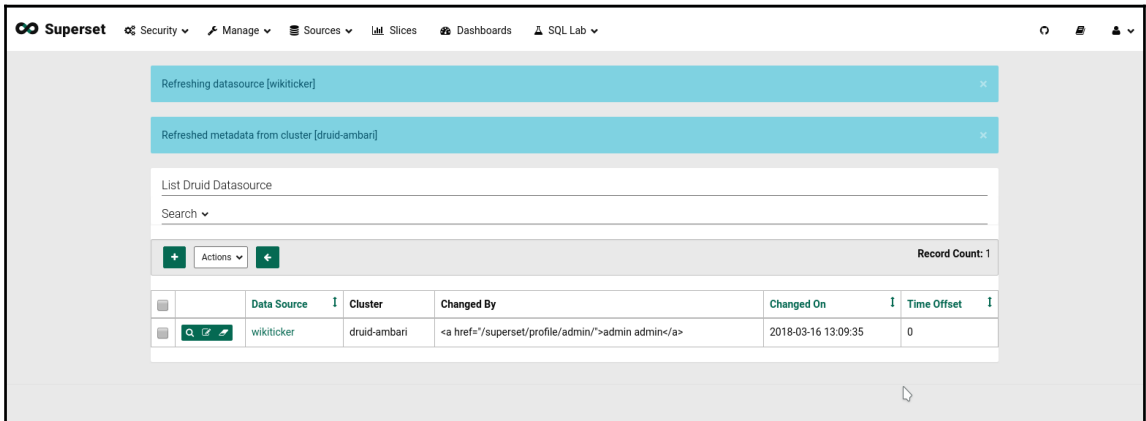
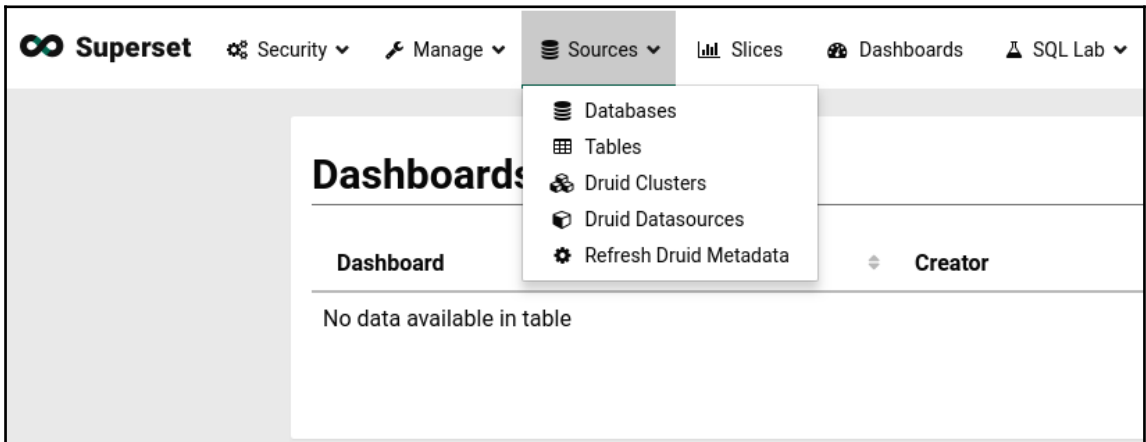
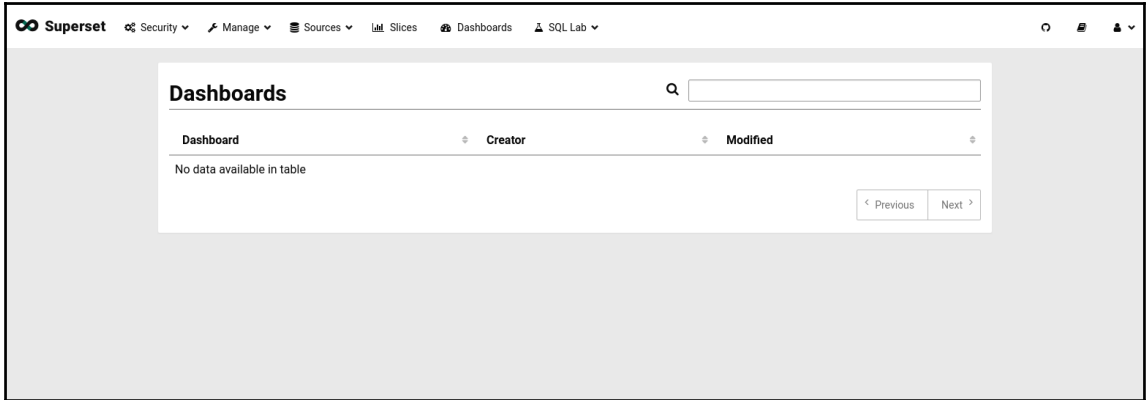
 Superset   [Login](#)

Sign In

Enter your login and password below:

Username:

Password:



Superset Security Manage Sources Slices Dashboards SQL Lab

Query Save as

Datasource & Chart Type

[druid-ambar].[wikiticker]

Table View

Time

Time Granularity: one day Origin:

Since: 7 days ago Until: now

GROUP BY

Group by:

Metrics:

[-] - untitled 0.28 sec json csv Query

No data was returned.

Superset Security Manage Sources Slices Dashboards SQL Lab

List Slice

Search

+ Actions -

Record Count: 0

No records found

+Actions ▾←

No records found

+Actions ▾←

| | | Data Source ↑ | Cluster |
|--------------------------|---|----------------------|----------------|
| <input type="checkbox"/> | | | |
| <input type="checkbox"/> | 🔍 ✎ ✂ | wikiticker | druid-ambari |

⚡ Query ➕ Save as

Datasource & Chart Type

[druid-ambari].[wikiticker] ✕ ✎

Big Number ▼

Time ?

| | |
|---------------------------------|-----------------------|
| Time Granularity ? | Origin ? |
| one day ▼ | ▼ |
| Since ? | Until |
| 5 years ago ▼ | now ▼ |

Metric ?

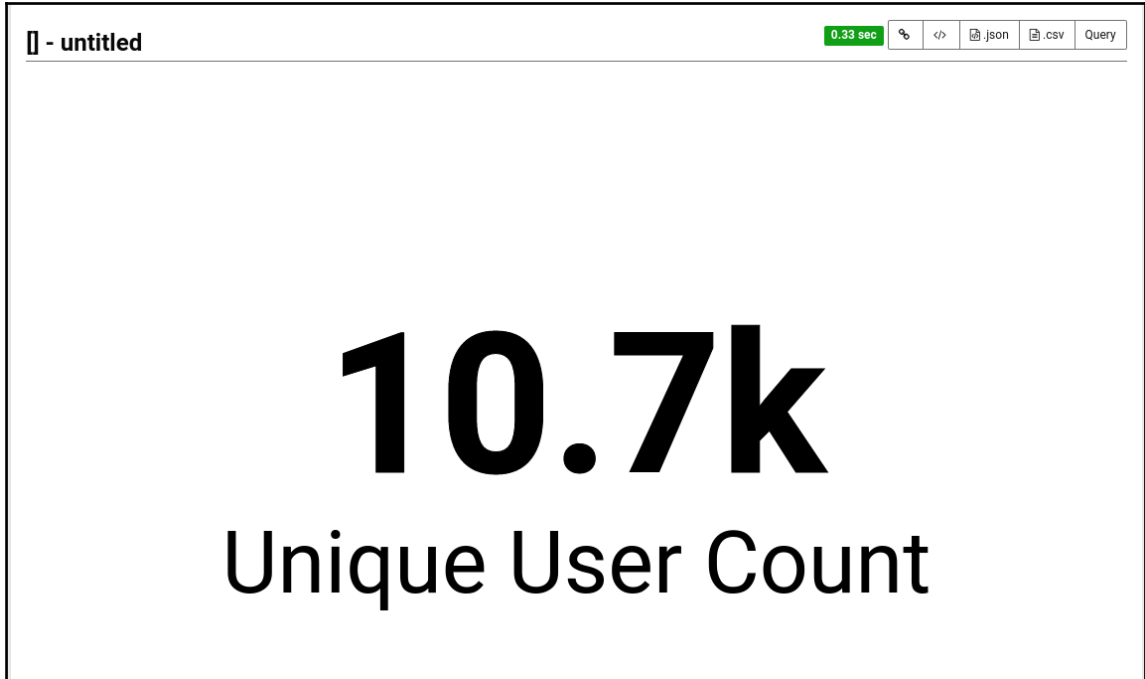
COUNT(DISTINCT user_unique) ▼

Subheader ?

Unique User Count

Number format ?

".3s" | 12.3k ▼



Save a Slice ✕

Save as


Do not add to a dashboard

Add slice to existing dashboard

Add to new dashboard

⚡ Query + Save as

Datasource & Chart Type

[druid-ambari].[wikiticker] ▼ 

Word Cloud ▼

Time ?

| | |
|---------------------------------|-----------------------|
| Time Granularity ? | Origin ? |
| one day ▼ | ▼ |
| Since ? | Until |
| 5 years ago ▼ | now ▼ |

Series ?

regionName ▼

Metric ?

COUNT(*) ▼

Series limit ?

10 ▼

Font Size From ? **Font Size To ?**

20 150

Rotation ?

flat ▼

Filters ?

countryIsoCode ▼

in ▼ US -

regionIsoCode ▼

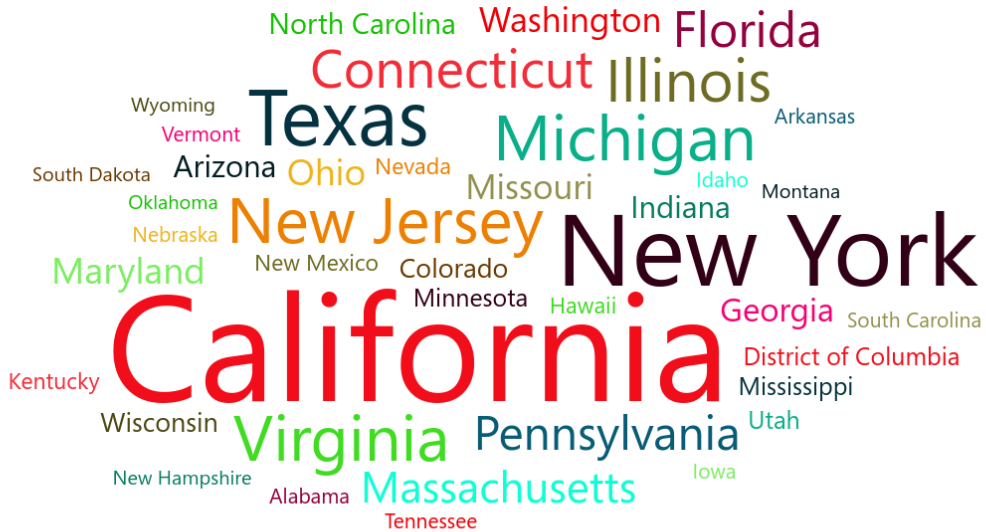
regex ▼ .+ -

+ Add filter

Word Cloud - Top US Regions ☆

0.19 sec

🔗 </> 📄 json 📄 csv Query



Save a Slice ✕

Overwrite slice [Unique Users]

Save as


Do not add to a dashboard

Add slice to existing dashboard

Add to new dashboard

⚡ Query ➕ Save as

Datasource & Chart Type

[druid-ambari].[wikiticker] ▼ 

Sunburst ▼

Time ?

| | |
|---------------------------------|-----------------------|
| Time Granularity ? | Origin ? |
| one day ▼ | ▼ |
| Since ? | Until |
| 5 years ago ▼ | now ▼ |

Hierarchy ?

Primary Metric ?

Secondary Metric ?

Row limit

Filters ?

Save a Slice ✕

Overwrite slice [Word Cloud - Top US Regions]

Save as

Do not add to a dashboard

Add slice to existing dashboard

Add to new dashboard

Sunburst - Top 10 Cities ☆ [📄](#)

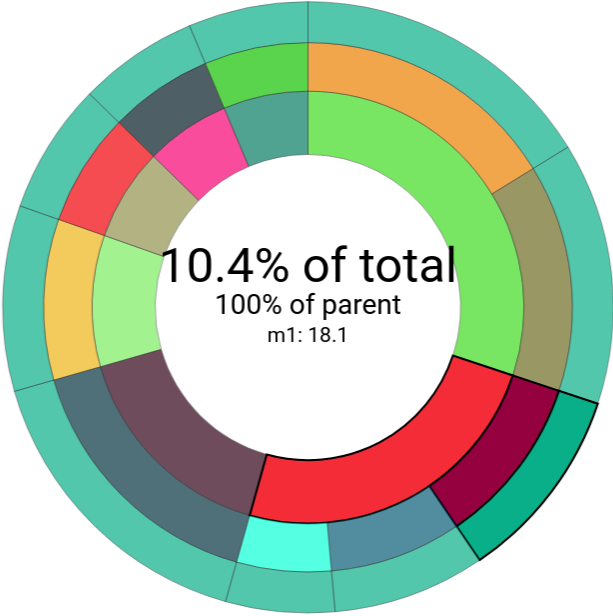
0.17 sec



json


csv

Query



⚡ Query ➕ Save as

Datasource & Chart Type

[druid-ambari].[wikiticker] ▼ 

Directed Force Layout ▼

Time ?

| | |
|---------------------------------|-----------------------|
| Time Granularity ? | Origin ? |
| one day ▼ | ▼ |
| Since ? | Until |
| 5 years ago ▼ | now ▼ |

Source / Target ?

✕ channel ✕ namespace

Metric ?

COUNT(DISTINCT user_unique) ▼

Row limit

50 ▼

Save a Slice ✕

Overwrite slice [Sunburst - Top 10 Cities]

Save as

Do not add to a dashboard

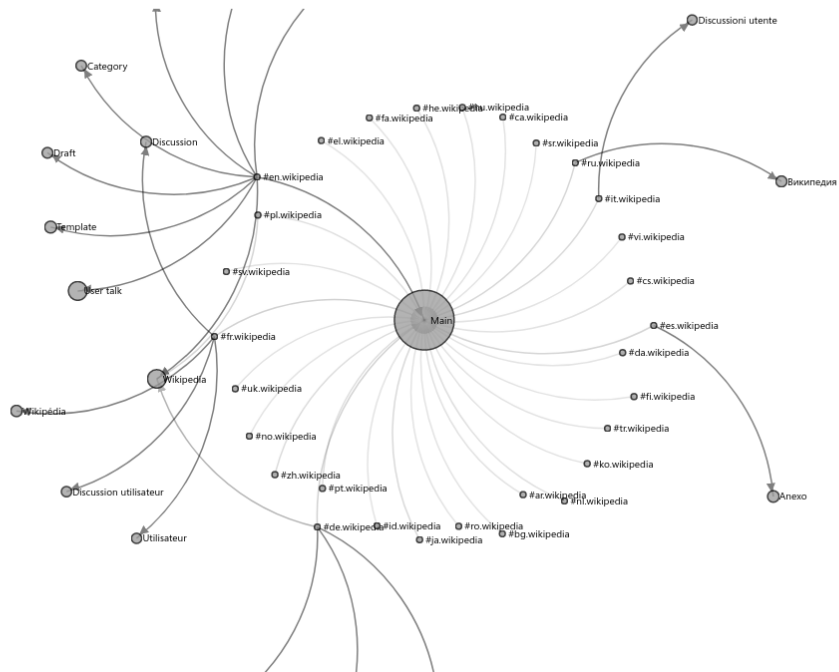
Add slice to existing dashboard

Add to new dashboard

DFL - Top 50 Channels & Namespaces

0.28 sec

[🔍](#) [</>](#) [.json](#) [.csv](#) [Query](#)



⚡ Query ➕ Save as

Datasource & Chart Type

[druid-ambari].[wikiticker] ▼ ✎

Sankey ▼

Time ?

| | |
|---------------------------------|-----------------------|
| Time Granularity ? | Origin ? |
| one day ▼ | ▼ |
| Since ? | Until |
| 5 years ago ▼ | now ▼ |

Source / Target ?

Metric ?

Row limit

Filters ?

Save a Slice ×

Save as

Do not add to a dashboard

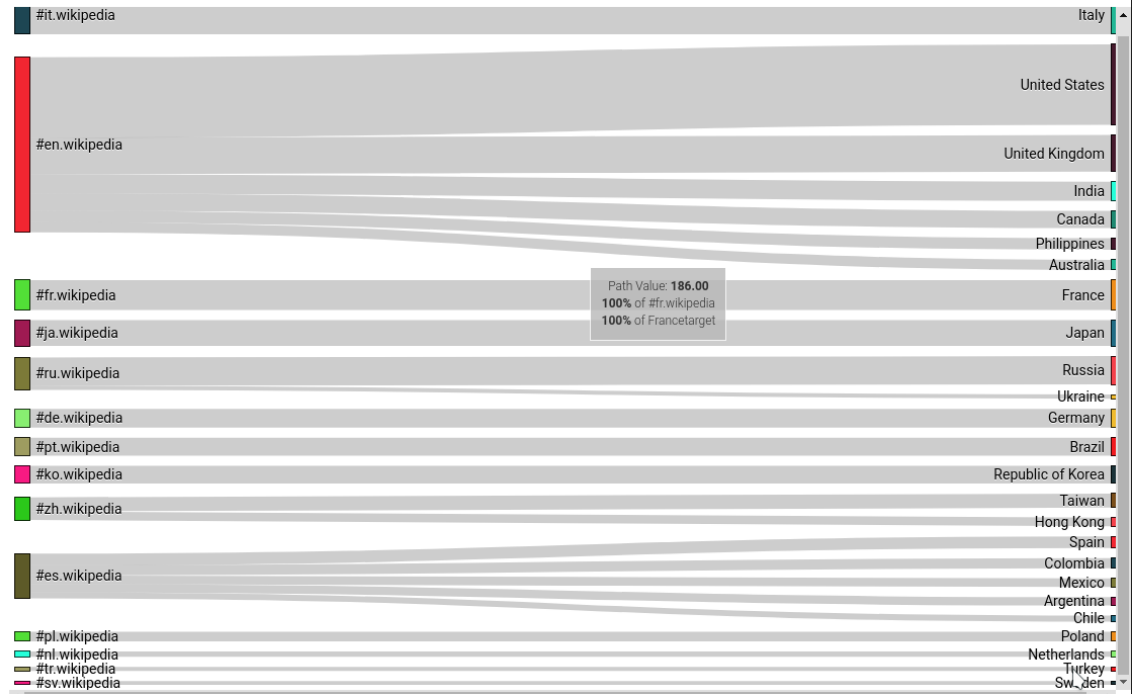
Add slice to existing dashboard

Add to new dashboard

Top 25 Countries / Channels Contribution

0.24 sec

🔗 </> json .csv Query



Superset Security Manage Sources Slices Dashboards SQL Lab

Add Dashboard

| | |
|--|--|
| Title | My Dashboard 1 |
| Slug | dash1 |
| To get a readable URL for your dashboard | |
| Slices | <ul style="list-style-type: none"> <input type="checkbox"/> Sunburst - Top 10 Cities <input type="checkbox"/> DFL - Top 50 Channels & Namespaces <input type="checkbox"/> Top 25 Countries / Channels Contribution <input type="checkbox"/> Word Cloud - Top US Regions <input type="checkbox"/> Unique Users |

Superset Security Manage Sources Slices Dashboards SQL Lab

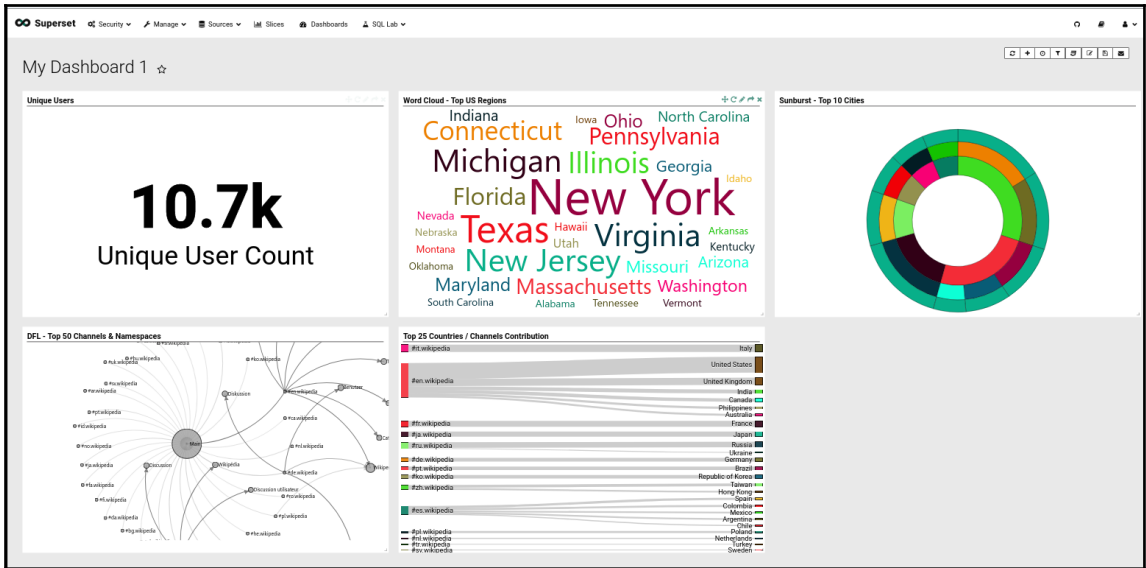
Added Row

List Dashboard

Search

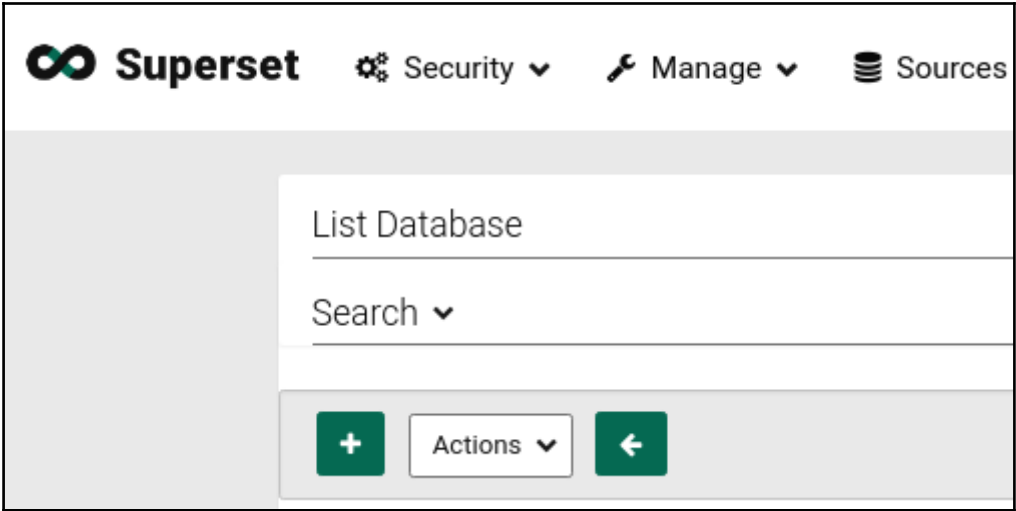
+ Actions ←

| | Dashboard | Creator |
|--------------------------|----------------|-------------|
| <input type="checkbox"/> | My Dashboard 1 | admin admin |



Sources ▾ **Slices**

- Databases
- Tables
- Druid Clusters
- Druid Datasources
- Refresh Druid Metadata



| | |
|--|---|
| Database | <input type="text" value="employees"/> |
| SQLAlchemy URI | <input type="text" value="mysql+pymysql://superset:XXXXXXXXXX@master:3306/employees"/> <p>Refer to the SQLAlchemy docs for more information on how to structure your URI.</p> <input type="button" value="Test Connection"/> |
| Cache Timeout | <input type="text" value="Cache Timeout"/> |
| Extra | <pre>{ "metadata_params": {}, "engine_params": {} }</pre> <p>JSON string containing extra configuration elements. The <code>engine_params</code> object gets unpacked into the <code>sqlalchemy.create_engine</code> call, while the <code>metadata_params</code> gets unpacked into the <code>sqlalchemy.MetaData</code> call.</p> |
| Expose in SQL Lab | Expose this DB in SQL Lab <input checked="" type="checkbox"/> |
| Allow Run Sync | Allow users to run synchronous queries, this is the default and should work well for queries that can be executed within a web request scope (<~1 minute) <input checked="" type="checkbox"/> |
| Allow Run Async | Allow users to run queries, against an async backend. This assumes that you have a Celery worker setup as well as a results backend. <input type="checkbox"/> |
| Allow CREATE TABLE AS | Allow CREATE TABLE AS option in SQL Lab <input type="checkbox"/> |
| Allow DML | Allow users to run non-SELECT statements (UPDATE, DELETE, CREATE, ...) in SQL Lab <input type="checkbox"/> |
| CTAS Schema | <input type="text" value="CTAS Schema"/> <p>When allowing CREATE TABLE AS option in SQL Lab, this option forces the table to be created in this schema</p> |
| <input type="button" value="Save"/> <input type="button" value="←"/> | |

Superset Security Manage Sources Slices Dashboards SQL Lab

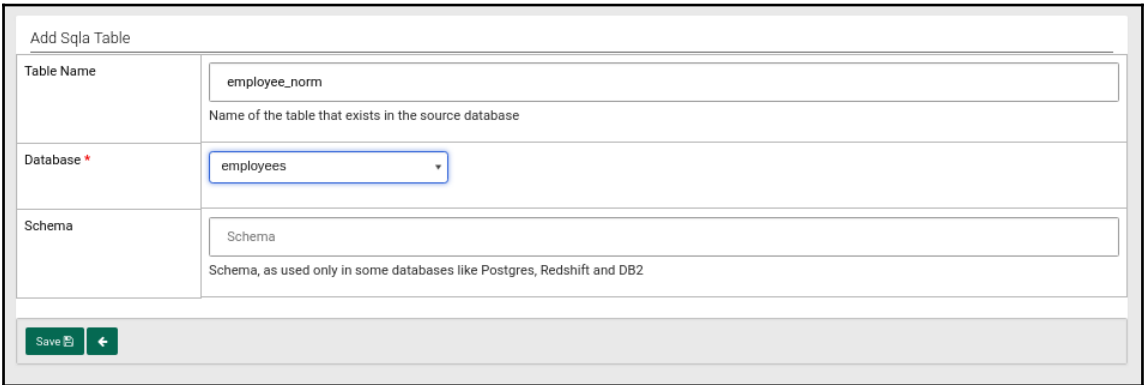
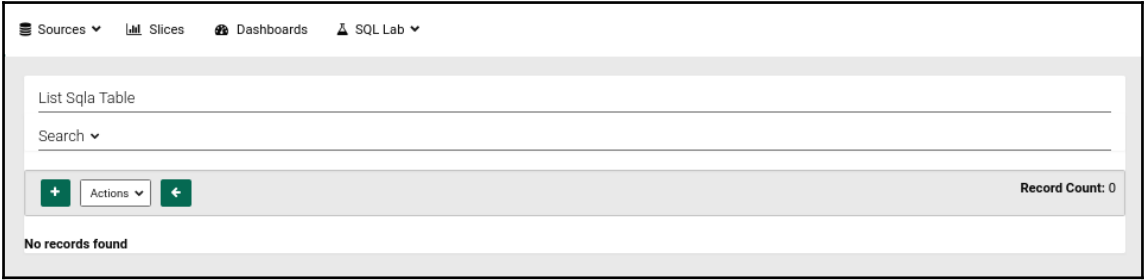
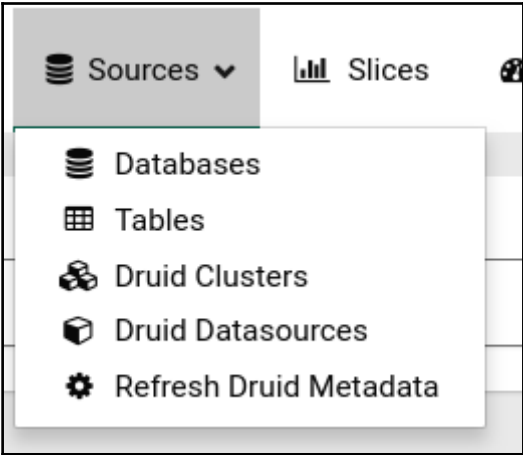
Added Row

List Database

Search

Record Count: 2

| | Database | Backend | Allow Run Sync | Allow Run Async | Allow DML | Creator | Last Changed |
|--------------------------|---|-----------|----------------|-----------------|-----------|---------|---------------------------------|
| <input type="checkbox"/> | <input type="button" value="🔍"/> <input type="button" value="🔗"/> | employees | mysql | True | False | False | admin admin 2018-03-24 08:34:15 |
| <input type="checkbox"/> | <input type="button" value="🔍"/> <input type="button" value="🔗"/> | main | mysql | True | False | False | 2018-03-15 13:19:14 |



The table was created. As part of this two phase configuration process, you should now click the edit button by the new table to configure it. ✕

Added Row ✕

List Sqla Table

Search ▾

+
Actions ▾
←
Record Count: 1

| <input type="checkbox"/> | Table | Database | Is Featured | Changed By | Last Changed |
|--------------------------|---|-----------|-------------|-------------|---------------------|
| <input type="checkbox"/> | 🔍 📄 🗑️ employee_norm | employees | False | admin admin | 2018-03-24 08:38:08 |

Edit Sqla Table

Detail List Table Column List Sql Metric

+
←
Record Count: 10

| | Column | Type | Groupable | Filterable | Count Distinct | Sum | Min | Max | Is temporal |
|---|------------|---------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|--------------------------|--------------------------|-------------------------------------|
| 🔍 📄 🗑️ | emp_no | INTEGER(11) | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| 🔍 📄 🗑️ | birth_date | DATE | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| 🔍 📄 🗑️ | full_name | VARCHAR(31) | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| 🔍 📄 🗑️ | gender | ENUM('M','F') | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| 🔍 📄 🗑️ | hire_date | DATE | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| 🔍 📄 🗑️ | salary | INTEGER(11) | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| 🔍 📄 🗑️ | from_date | DATE | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| 🔍 📄 🗑️ | to_date | DATE | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| 🔍 📄 🗑️ | dept_name | VARCHAR(40) | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| 🔍 📄 🗑️ | title | VARCHAR(50) | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |

⚡ Query ➕ Save as

Datasource & Chart Type

[employees].[employee_norm] ✎

Distribution - NVD3 - Pie Chart ▼

Time ?

Time Column ? birth_date ▼ Time Grain ? Time Column ▼

Since ? 100 years ago ▼ Until ? now ▼

Metrics ?

✕ sum__salary

Group by ?

✕ dept_name

Series limit ?

50 ▼

Label Type ?

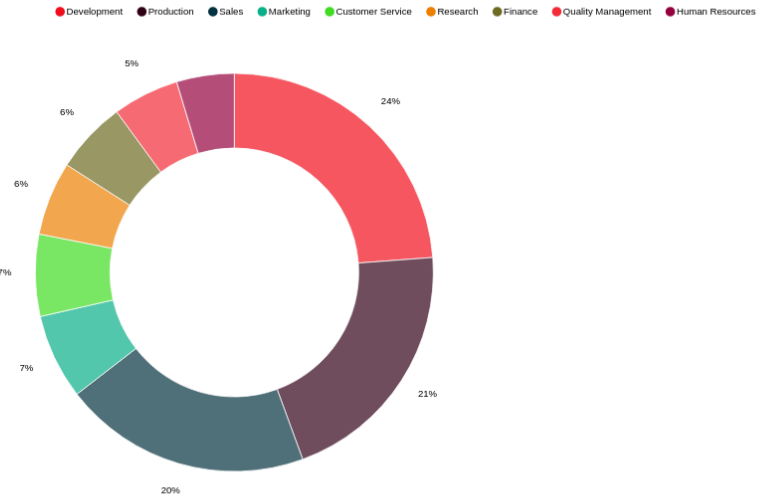
Percentage ▼

Donut ? Legend ?

Put labels outside ?


Department Salary Breakup


5.01 sec % </> json .csv Query





⚡ Query 💾 Save as



Datasource & Chart Type



[employees].[employee_norm] 


Time Series - Line Chart 

Time

Time Column  from_date 


Time Grain  year 



Since  100 years ago 



Until now 

Metrics

✖ avg_salary

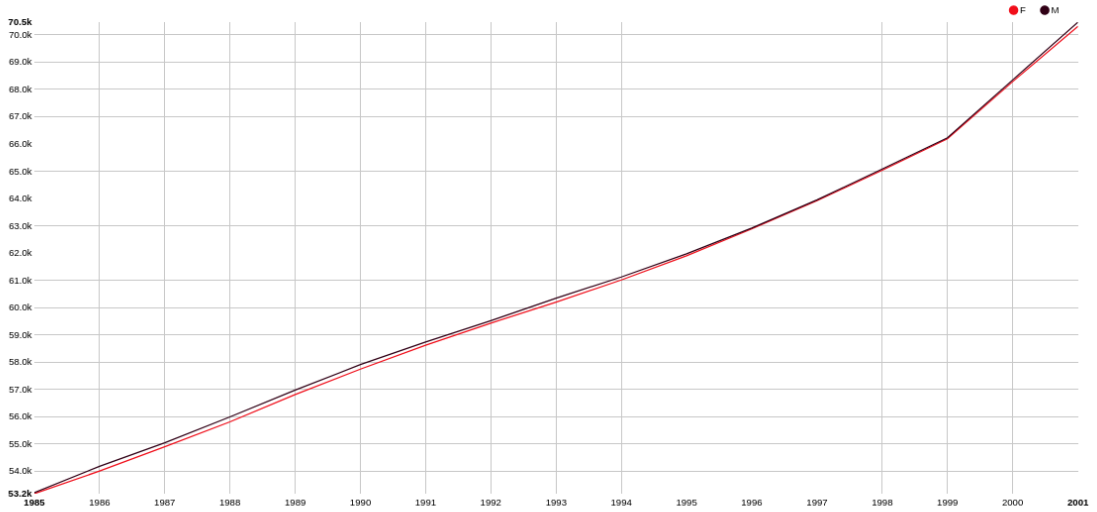
Group by  ✖ gender

Series limit  50 

Sort By  


Salary Diversity


11.67 sec [🔍](#) [🏠](#) [📄 json](#) [📄 csv](#) [Query](#)






⚡ Query ➕ Save as



Datasource & Chart Type


[employees].[employee_norm] 



Time Series - Percent Change 


Time 

Time Column  **Time Grain **


from_date  year 

Since  **Until**



100 years ago  now 



Metrics 

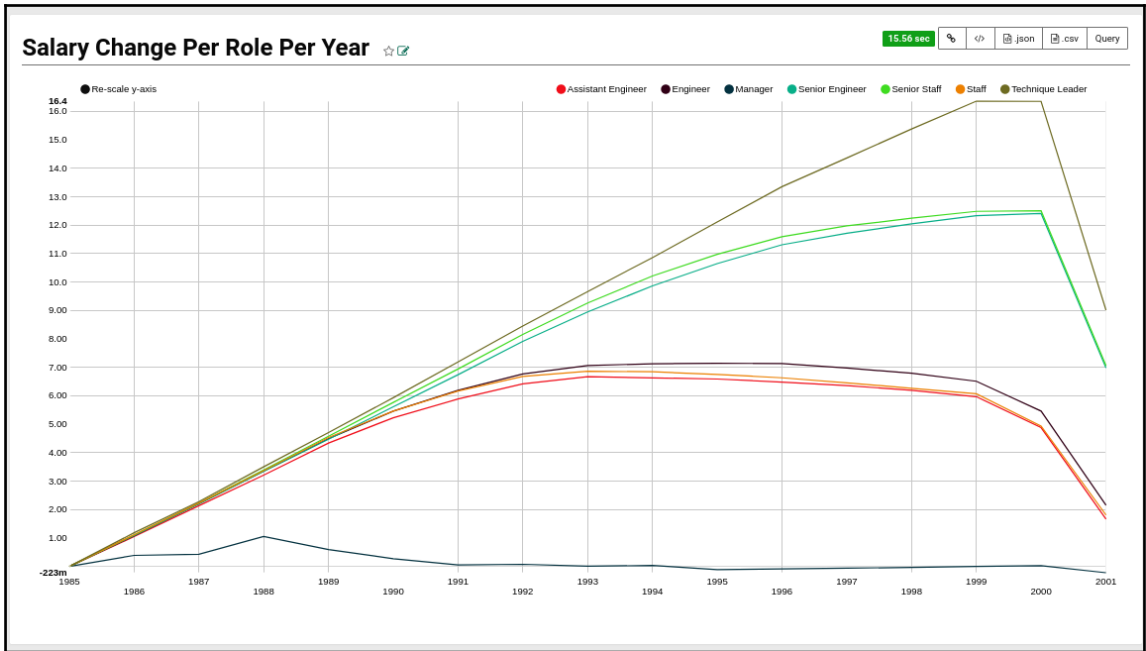
✖ sum__salary

Group by 

✖ title

Series limit  **Sort By **

50  

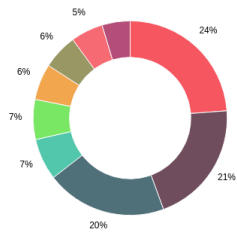


| | |
|--------|--|
| Title | My Dashboard 2 |
| Slug | dash2 To get a readable URL for your dashboard |
| Slices | <input checked="" type="checkbox"/> Department Salary Breakup <input checked="" type="checkbox"/> Salary Diversity <input checked="" type="checkbox"/> Salary Change Per Role Per Year |

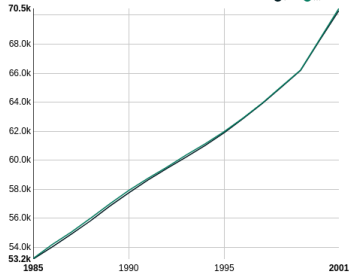
My Dashboard 2 ☆

Department Salary Breakup

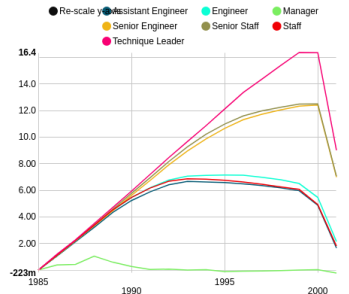
- Development
- Customer Service
- Human Resources
- Production
- Research
- Sales
- Marketing
- Finance
- Quality Management



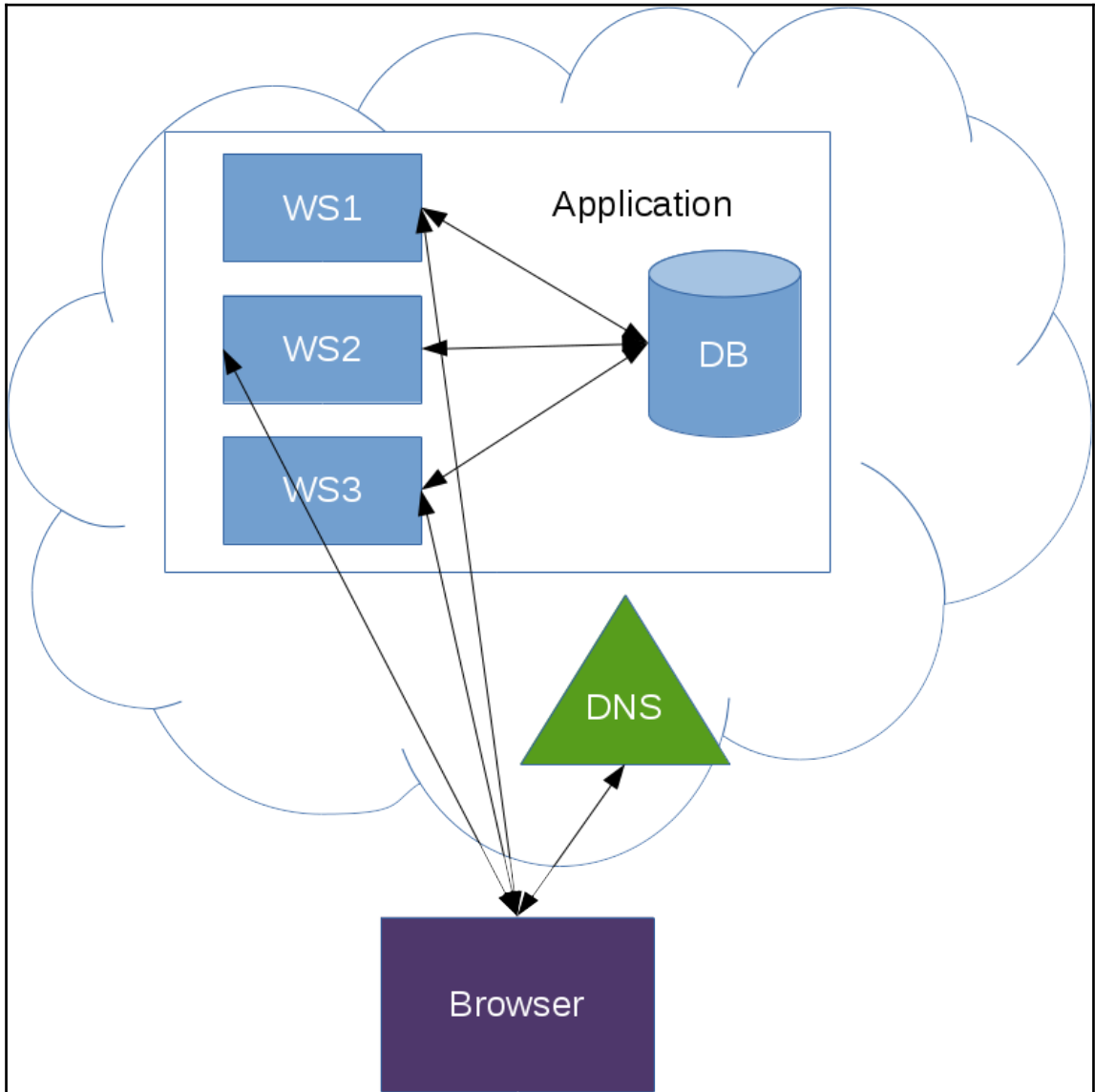
Salary Diversity

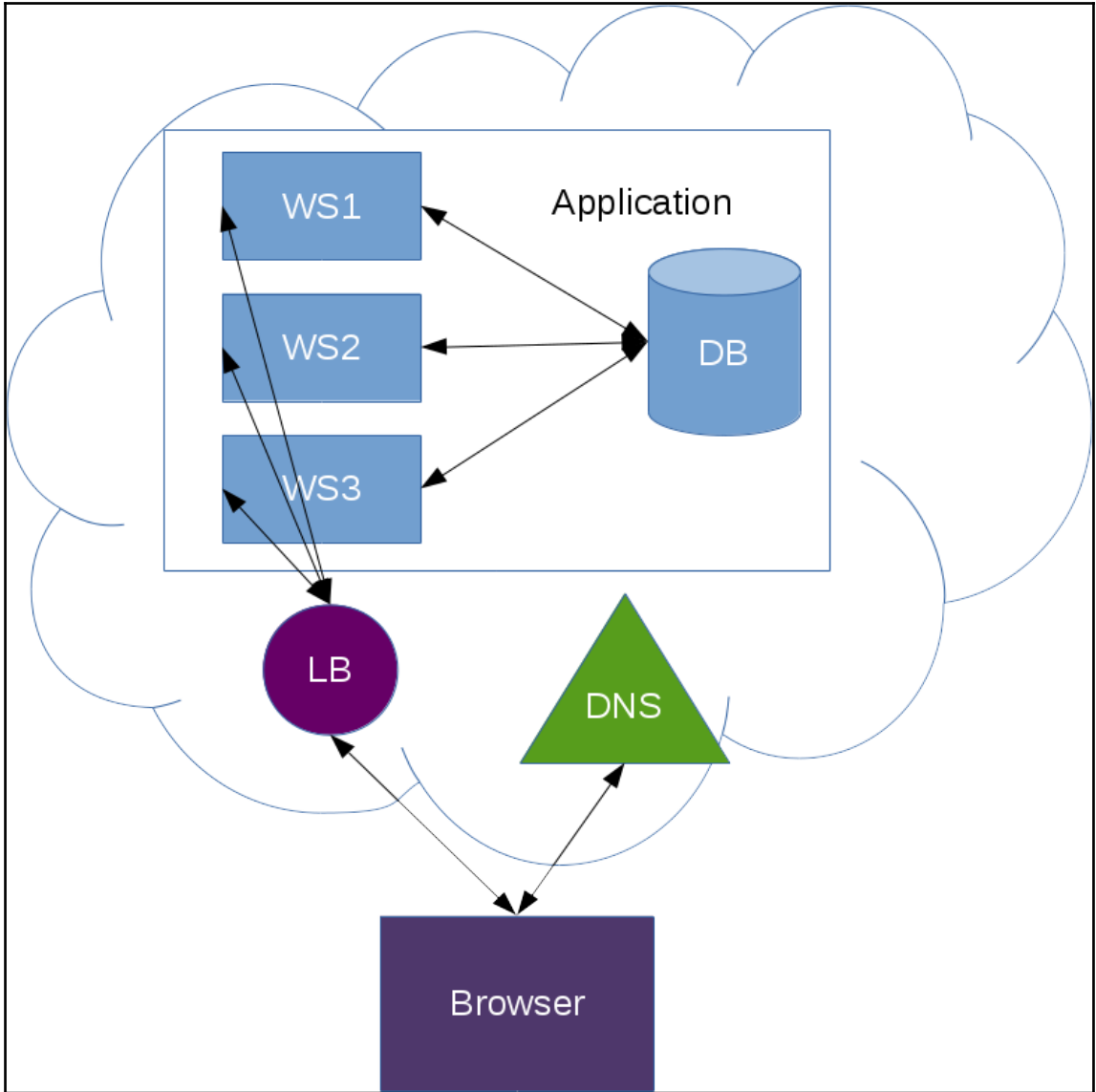


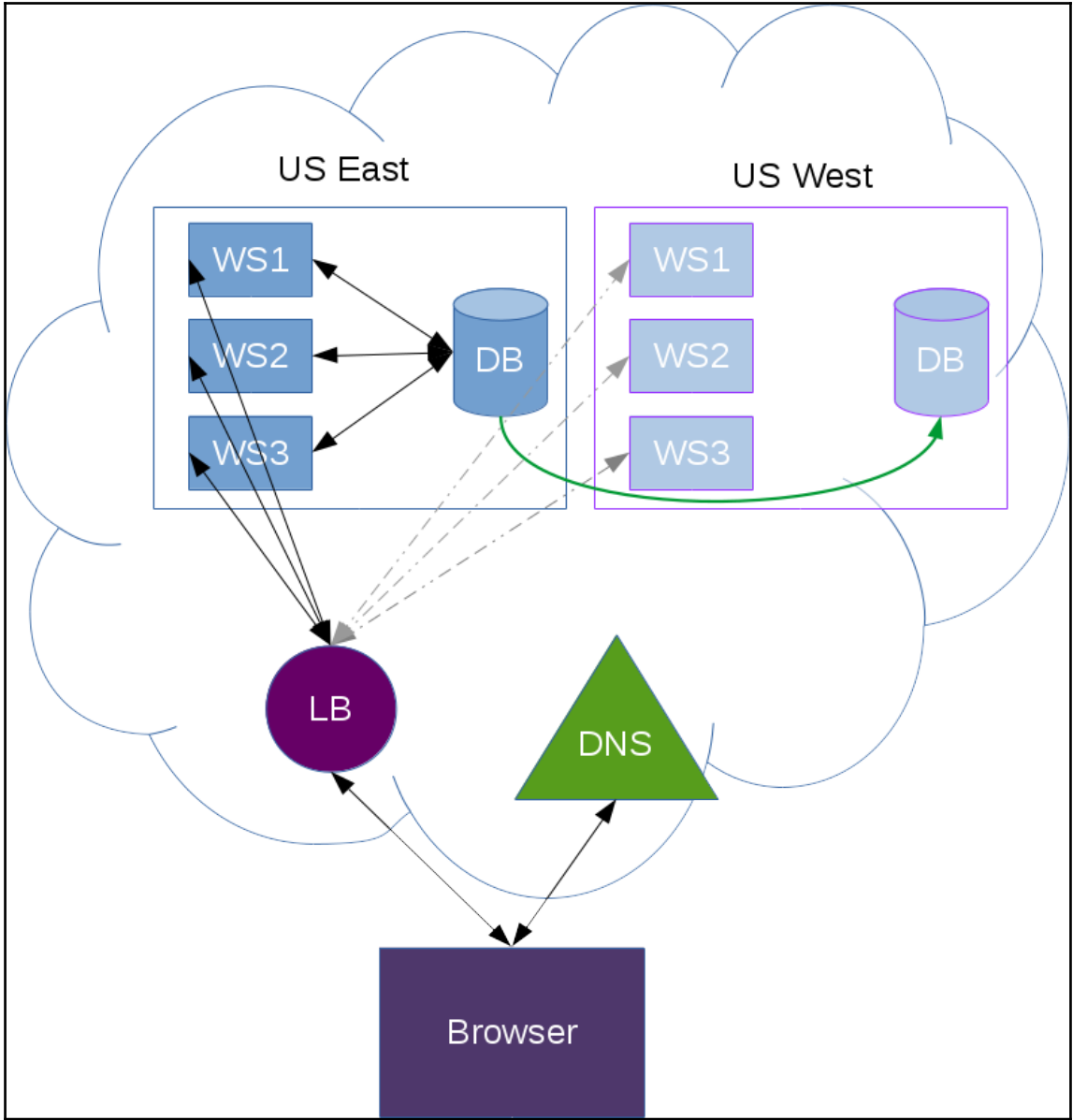
Salary Change Per Role Per Year

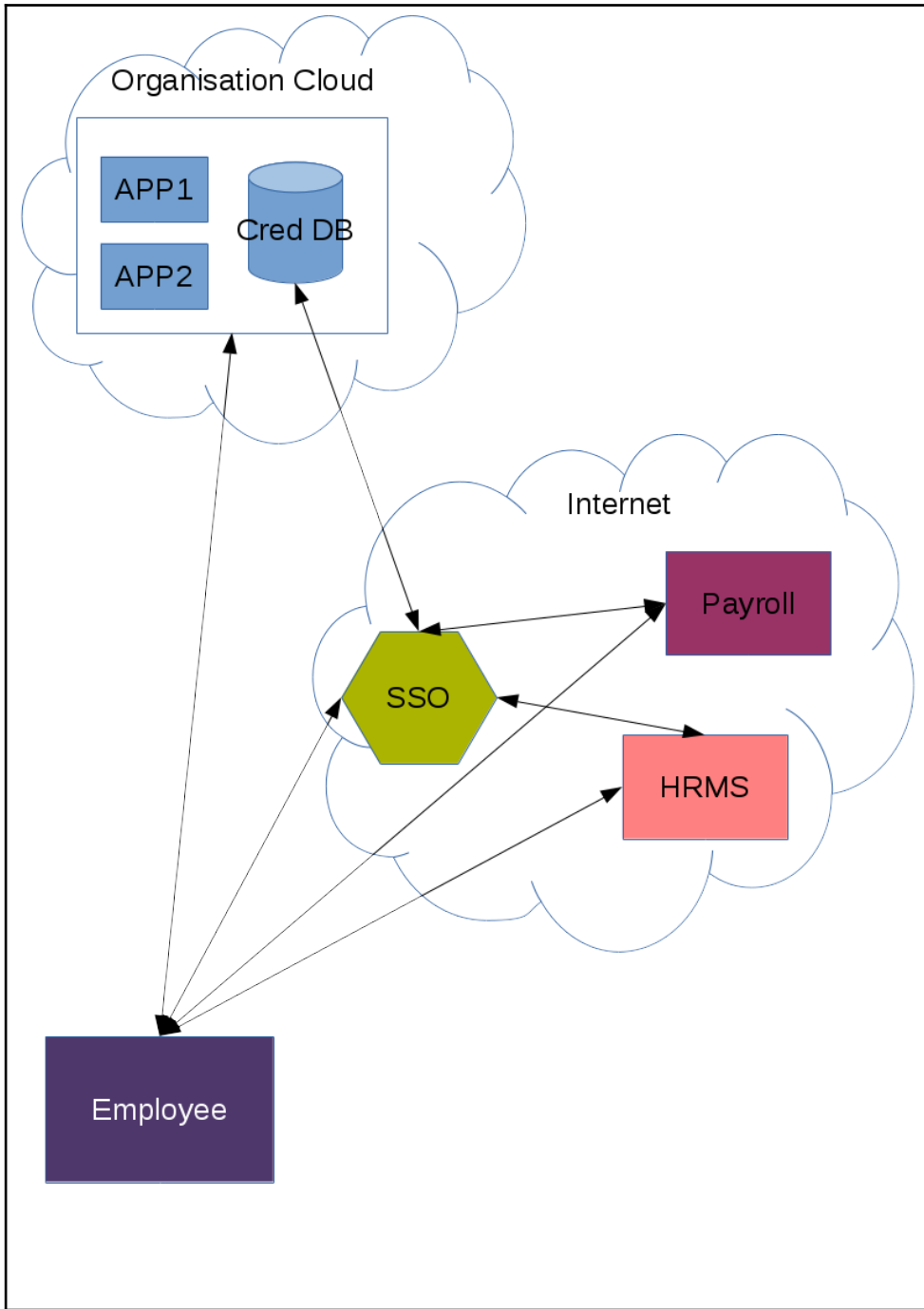


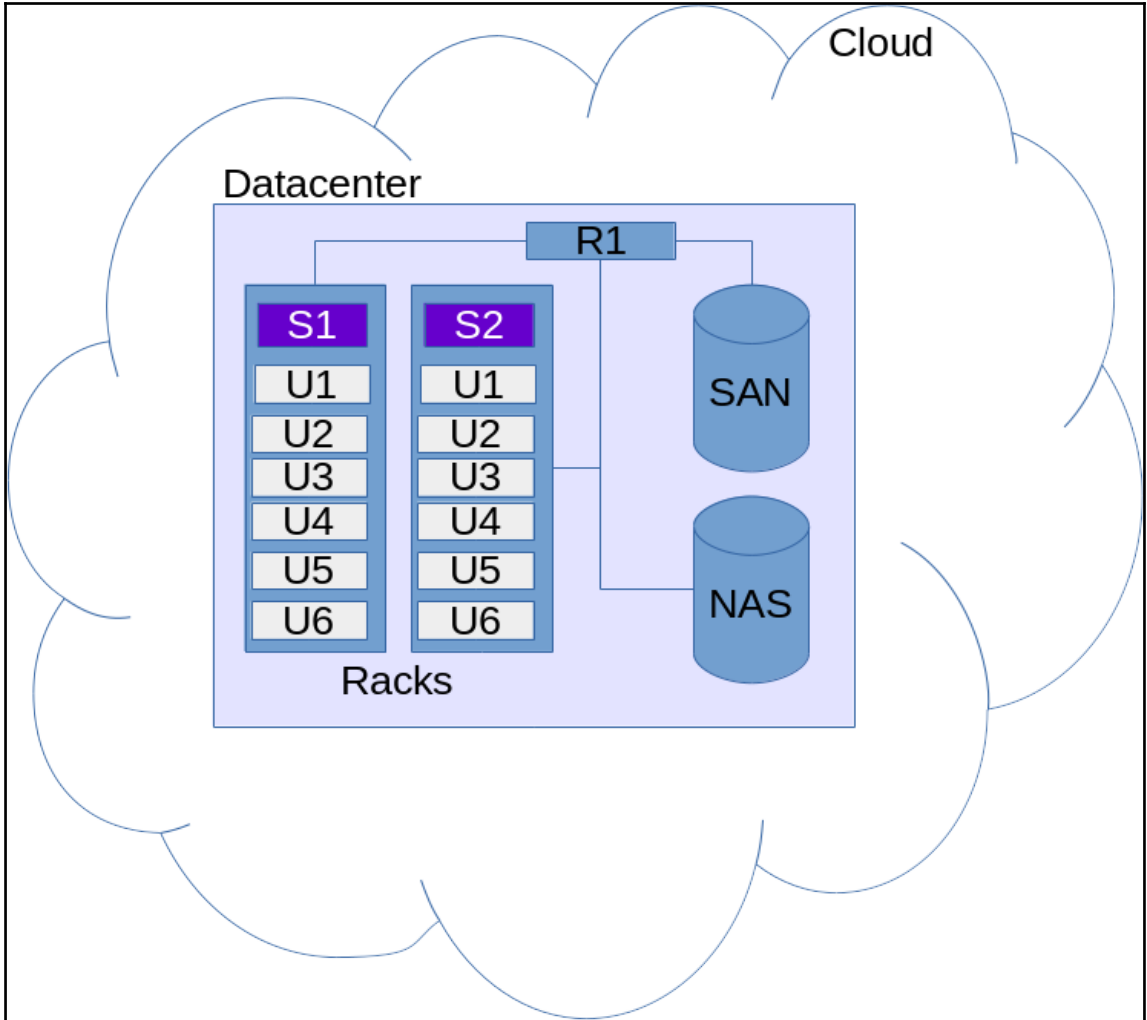
Chapter 10: Developing Applications Using the Cloud











The screenshot shows the Google Cloud Platform interface for the 'Clusters' page. The top navigation bar includes the Google Cloud logo, the text 'Google Cloud Platform', and a dropdown menu for 'My First Project'. A search bar is located on the right side of the navigation bar. The main content area is titled 'Clusters' and contains a card for 'Cloud Dataproc Clusters'. The card text reads: 'Google Cloud Dataproc lets you provision Apache Hadoop clusters and connect to underlying analytic data stores. Create your first cluster to get started.' Below this text is a button labeled 'Enabling' with three dots to its right, indicating a loading or progress state. The left sidebar contains navigation icons for home, clusters, and a menu icon.

Google Cloud Platform My First Project

Create a cluster

Name ?
packt

Region ? asia-south1 **Zone** ? asia-south1-a

Cluster mode ?
Single Node (1 master, 0 workers)

Single Node
Contains all Hadoop master and worker daemons and all job drivers

Machine type ?
4 vCPUs 15 GB memory [Customize](#)
[Upgrade your account](#) to create instances with up to 96 cores

Primary disk size (minimum 10 GB) ?
500 GB

YARN cores ? 4 **YARN memory** ? 12.0 GB

[Bucket, network, version, initialization, & access options](#)

[Create](#) [Cancel](#)

Secure | <https://ssh.cloud.google.com/projects/molten-amulet-189901/zones/asia-south1-a/instances/packt-m?authuser...>
⋮ ⚙

```

[.-text [-ignoreCrc] <src> ...]
[.-touchz <path> ...]
[.-truncate [-w] <length> <path> ...]
[.-usage [cmd ...]]

Generic options supported are
.-conf <configuration file> specify an application configuration file
.-D <property=value> use value for given property
.-fs <file:///hdfs://namenode:port> specify default filesystem URL to use, overrides 'fs.defaultFS' property from
configurations.
.-jt <local|resourceManager:port> specify a ResourceManager
.-files <comma separated list of files> specify comma separated files to be copied to the map reduce cluster
.-libjars <comma separated list of jars> specify comma separated jar files to include in the classpath.
.-archives <comma separated list of archives> specify comma separated archives to be unarchived on the compute m
achines.

The general command line syntax is
command [genericOptions] [commandOptions]

hop five_in@packt-m:/opt$ hadoop fs -ls
18/02/17 07:04:07 INFO gcs.GoogleHadoopFileSystemBase: GHFS version: 1.6.3-hadoop2
ls: `.`: No such file or directory
hop five_in@packt-m:/opt$ cd
hop five_in@packt-m:~$ hadoop fs -ls /
18/02/17 07:04:15 INFO gcs.GoogleHadoopFileSystemBase: GHFS version: 1.6.3-hadoop2
Found 2 items
drwxrwxrwt - mapred hadoop 0 2018-02-17 07:01 /tmp
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user
hop five_in@packt-m:~$ hadoop fs -mkdir /user/packt
18/02/17 07:04:37 INFO gcs.GoogleHadoopFileSystemBase: GHFS version: 1.6.3-hadoop2
hop five_in@packt-m:~$ hadoop fs -ls /user
18/02/17 07:04:47 INFO gcs.GoogleHadoopFileSystemBase: GHFS version: 1.6.3-hadoop2
Found 9 items
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/hbase
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/hdfs
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/hive
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/mapred
drwxr-xr-x - hop_five_in hadoop 0 2018-02-17 07:04 /user/packt
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/pig
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/spark
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/yarn
drwxrwxrwt - hdfs hadoop 0 2018-02-17 07:01 /user/zookeeper
hop five_in@packt-m:~$

```

Google Cloud Platform
My First Project
🔍

Cluster details
REFRESH
DELETE
VIEW LOGS

packt
Overview Jobs VM Instances Configuration

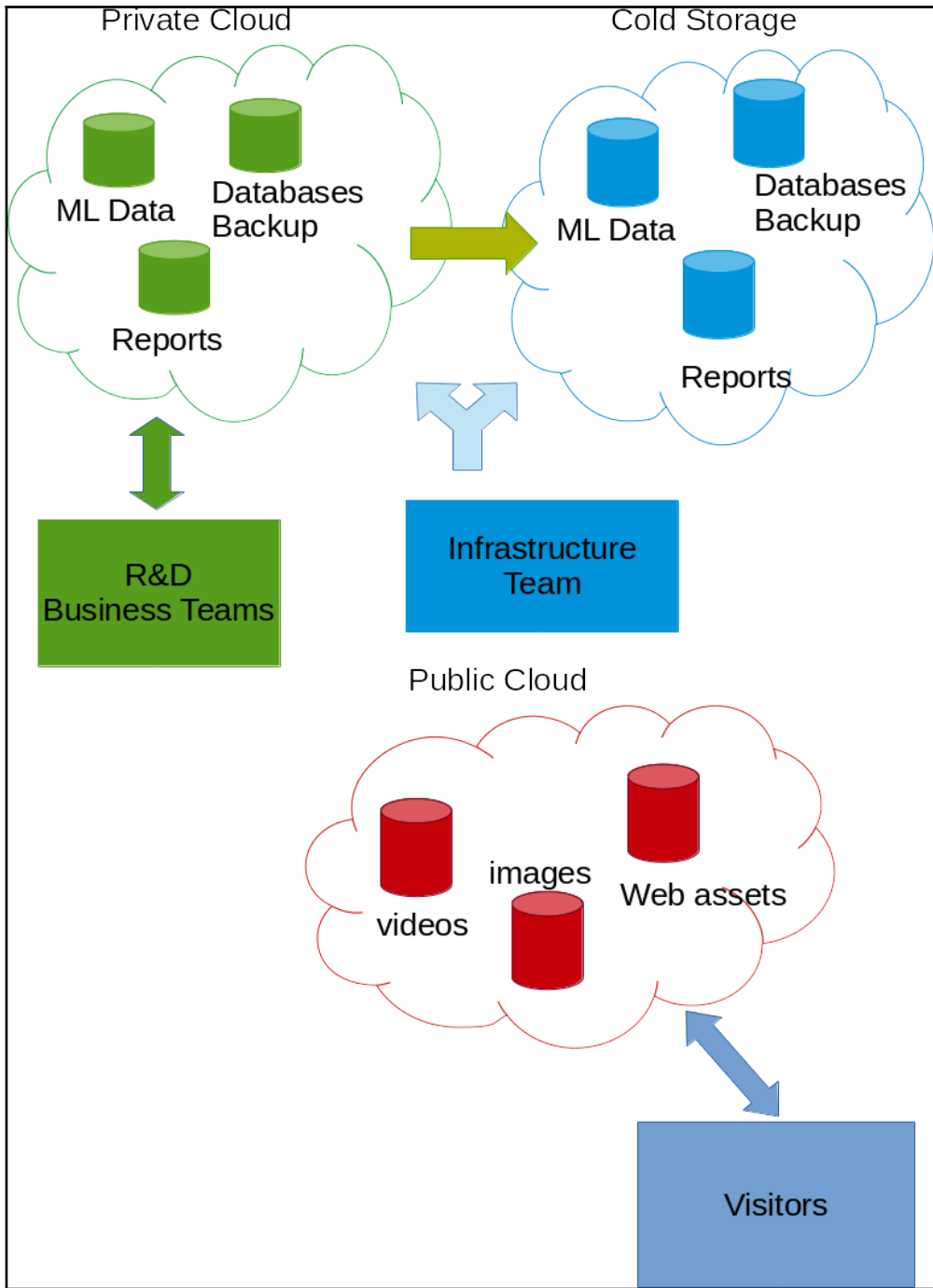
Edit

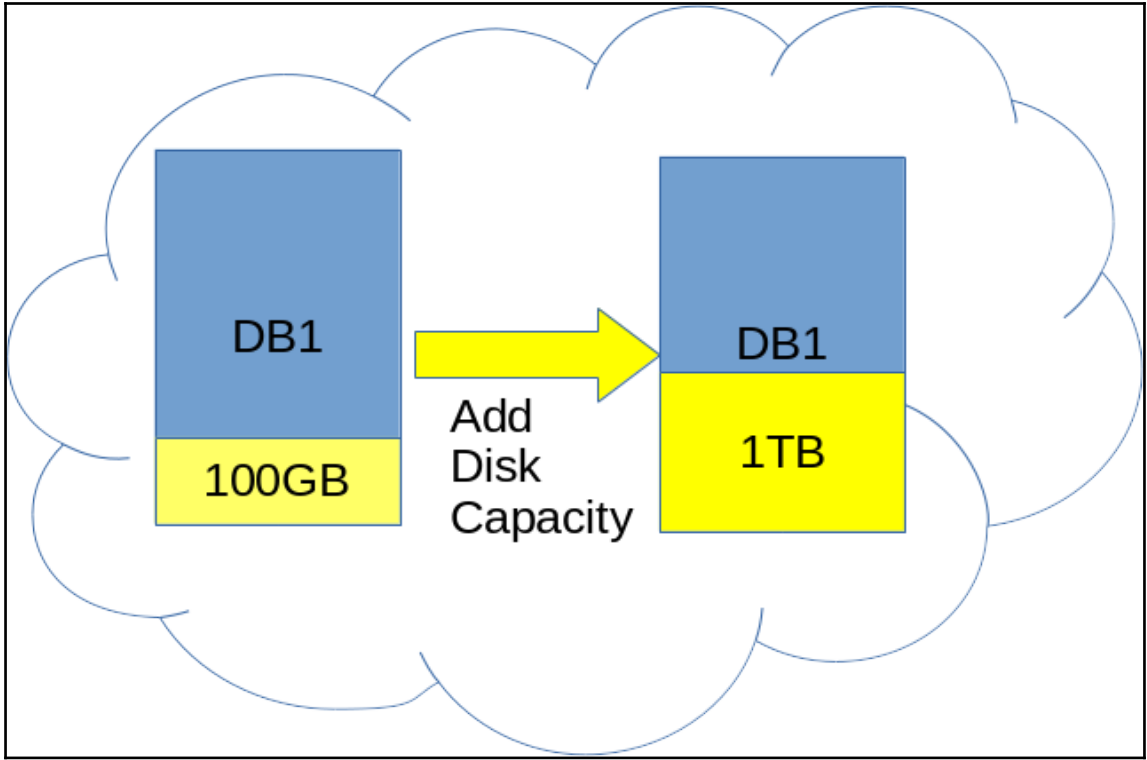
| | |
|------------------------------|---|
| Name | packt |
| Region | asia-south1 |
| Zone | asia-south1-a |
| Master node | Single Node (1 master, 0 work) |
| Machine type | n1-standard-4 (4 vCPU, 15.0 GB) |
| Primary disk size | 500 GB |
| Cloud Storage staging bucket | dataproc-8322bc7f-2ff5-4e3f-930c-300000000000 |
| Subnetwork | default |
| Network tags | None |
| Internal IP only | No |

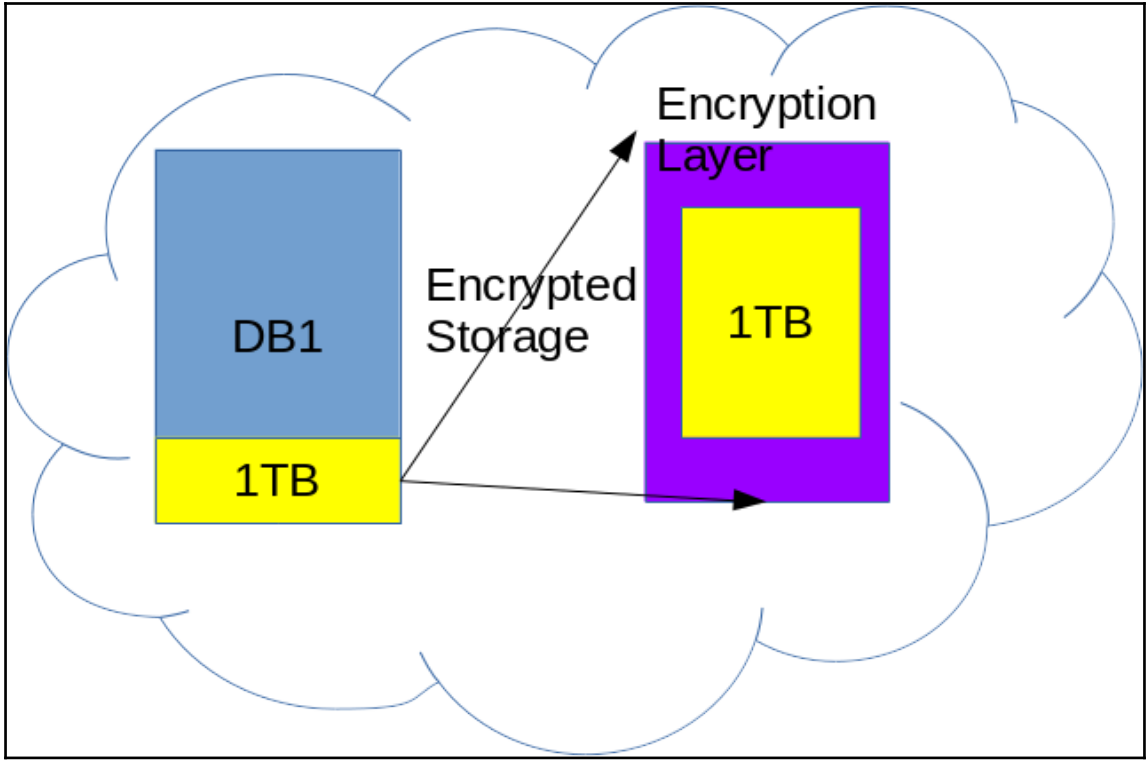
Confirm deletion

Are you sure you wish to delete cluster packt?

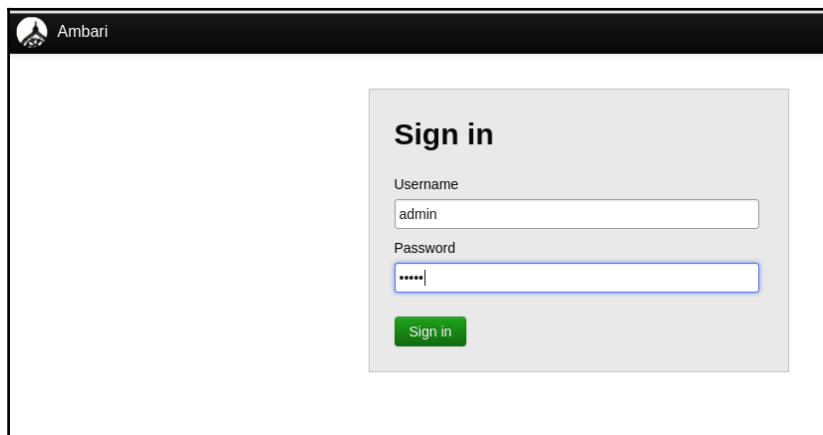
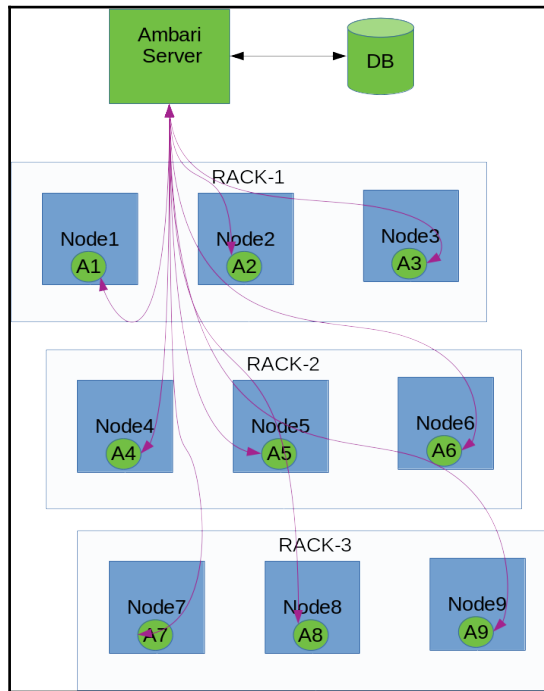
CANCEL OK

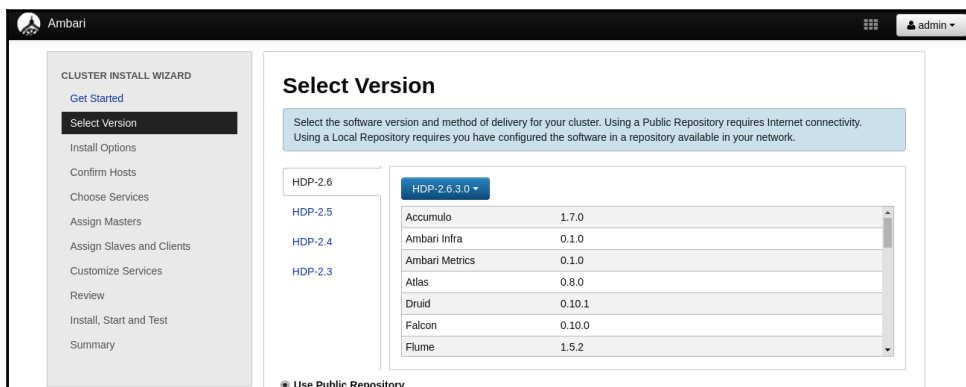
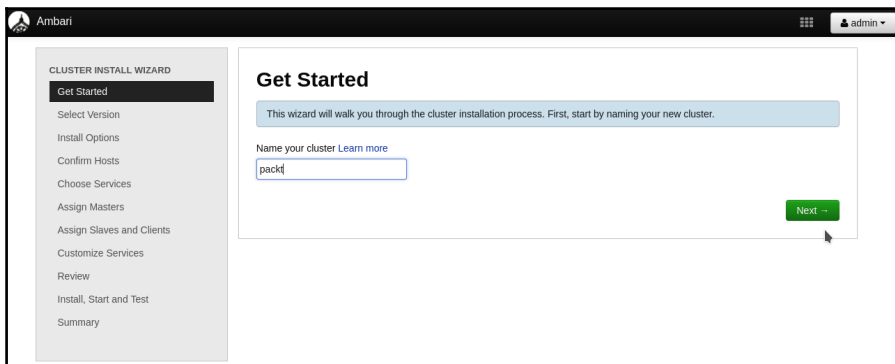
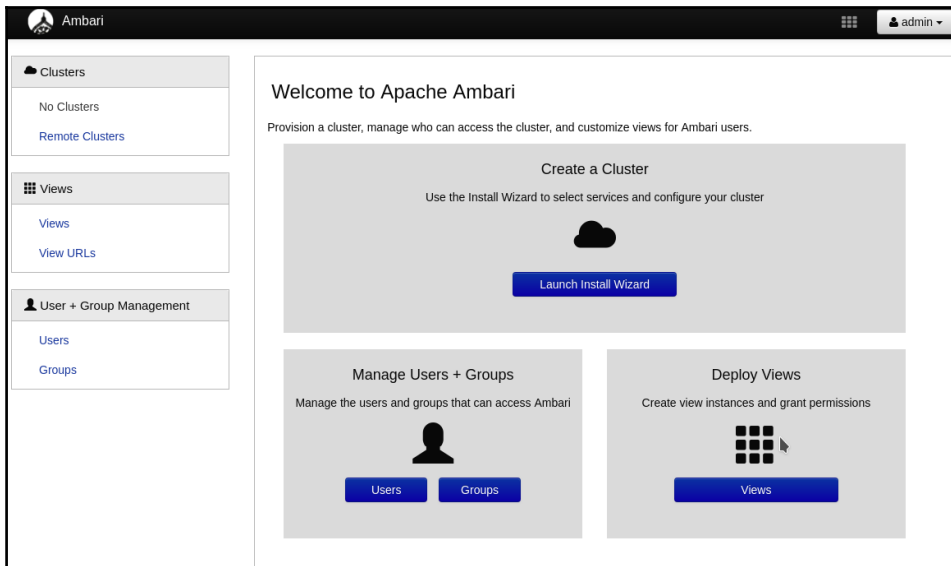






Chapter 11: Production Hadoop Cluster Deployment





Ambari admin

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options**
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Install Options

Enter the list of hosts to be included in the cluster and provide your SSH key.

Target Hosts
Enter a list of hosts using the Fully Qualified Domain Name (FQDN), one per line. Or use [Pattern Expressions](#)

```
node-1
node-2
node-3
```

Host Registration Information

- Provide your **SSH Private Key** to automatically register hosts
 - Choose file: No file chosen
 - 90fC01085019cXpHs4ecXpXK1t04EVBnVhMPH0NEXZC0dAwV2qYTAfQ
-----END RSA PRIVATE KEY-----
 - SSH User Account:
 - SSH Port Number:
- Perform **manual registration** on hosts and do not use SSH

[-- Back](#) [Register and Confirm --](#)

Ambari admin

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts**
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Confirm Hosts

Registering your hosts.
Please confirm the host list and remove any hosts that you do not want to include in the cluster.

Remove Selected Show: All (3) | Installing (0) | Registering (0) | Success (3) | Fail (0)

| Host | Progress | Status | Action |
|--------|----------------------------------|---------|------------------------|
| node-1 | <div style="width: 100%;"></div> | Success | Remove |
| node-2 | <div style="width: 100%;"></div> | Success | Remove |
| node-3 | <div style="width: 100%;"></div> | Success | Remove |

Show: 25 1 - 3 of 3 [←](#) [→](#)

Some warnings were encountered while performing checks against the 3 registered hosts above [Click here to see the warnings.](#)

[-- Back](#) [Next --](#)

Ambari admin

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts
- Choose Services**
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Choose Services

Choose which services you want to install on your cluster.

| Service | Version | Description |
|---|----------|---|
| <input checked="" type="checkbox"/> HDFS | 2.7.3 | Apache Hadoop Distributed File System |
| <input checked="" type="checkbox"/> YARN + MapReduce2 | 2.7.3 | Apache Hadoop NextGen MapReduce (YARN) |
| <input type="checkbox"/> Tez | 0.7.0 | Tez is the next generation Hadoop Query Processing framework written on top of YARN. |
| <input type="checkbox"/> Hive | 1.2.1000 | Data warehouse system for ad-hoc queries & analysis of large datasets and table & storage management service |
| <input type="checkbox"/> HBase | 1.1.2 | A Non-relational distributed database, plus Phoenix, a high performance SQL layer for low latency applications. |
| <input type="checkbox"/> Pig | 0.16.0 | Scripting platform for analyzing large datasets |
| <input type="checkbox"/> Sqoop | 1.4.6 | Tool for transferring bulk data between Apache Hadoop and structured data stores such as relational databases |

Ambari admin

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts
- Choose Services
- Assign Masters**
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Assign Masters

Assign master components to hosts you want to run them on.

NameNode: node-1

SNameNode: node-2

ResourceManager: node-2

App Timeline Server: node-2

History Server: node-2

ZooKeeper Server: node-1

ZooKeeper Server: node-3

ZooKeeper Server: node-2

Grafana: node-1

Metrics Collector: node-3

HST Server: node-1

Activity Explorer: node-1

node-1 (12.6 GB, 2 cores)

- NameNode
- ZooKeeper Server
- Grafana
- HST Server
- Activity Explorer
- Activity Analyzer

node-2 (12.6 GB, 2 cores)

- SNameNode
- ResourceManager
- App Timeline Server
- History Server
- ZooKeeper Server

node-3 (12.6 GB, 2 cores)

- ZooKeeper Server
- Metrics Collector

Ambari admin

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on. Hosts that are assigned master components are shown with *.

*"Client" will install HDFS Client, YARN Client, MapReduce2 Client and ZooKeeper Client.

| Host | all none | all none | all none | all none |
|----------|--------------------------|--------------------------|--------------------------|-------------------------------------|
| node-1 * | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| node-2 * | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| node-3 * | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

Services: HDFS, YARN, MapReduce2, ZooKeeper, Ambari Metrics, SmartSense, Misc

Group: Default (3) Manage Config Groups Filter...

Basic Data Capture Operations Gateway Activity Analysis **Advanced**

Activity Explorer

UI port: 9060

Anonymous access allowed: No

Password for user 'admin':

Activity Analyzer Activity Storage

Back Next

Ambari admin

Customize Services

We have come up with recommended configurations for the services you selected. Customize them as you see fit.

HDFS YARN MapReduce2 ZooKeeper Ambari Metrics SmartSense Misc

Group: Default (3) Manage Config Groups Filter...

Basic Data Capture Operations Gateway Activity Analysis **Advanced**

Activity Explorer

UI port: 9060

Anonymous access allowed: No

Password for user 'admin':

Activity Analyzer Activity Storage

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Review

Please review the configuration before installation

Admin Name : admin

Cluster Name : packt

Total Hosts : 3 (3 new)

Repositories:

```

debian7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/debian7/2.x/updates/2.6.3.0
debian7 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/debian7
redhat-ppc7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos7-ppc/2.x/updates/2.6.3.0
redhat-ppc7 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/centos7-ppc
redhat6 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos6/2.x/updates/2.6.3.0
redhat6 (HDP-UTILS-1.1.0.21):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.21/repos/centos6
redhat7 (HDP-2.6):
http://public-repo-1.hortonworks.com/HDP/centos7/2.x/updates/2.6.3.0
          
```

← Back
Print
Deploy →

Ambari
adm

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Install, Start and Test

Please wait while the selected services are installed and started.

100 % overall

| Host | Status | Message |
|--------|--|---------|
| node-1 | <div style="background-color: #4CAF50; width: 100%; height: 10px;"></div> 100% | Success |
| node-2 | <div style="background-color: #4CAF50; width: 100%; height: 10px;"></div> 100% | Success |
| node-3 | <div style="background-color: #4CAF50; width: 100%; height: 10px;"></div> 100% | Success |

3 of 3 hosts showing · [Show All](#) Show: 25 | 1 - 3 of 3

Successfully installed and started the services.

Next →

Ambari

admin

CLUSTER INSTALL WIZARD

- Get Started
- Select Version
- Install Options
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary**

Summary

Here is the summary of the install process.

The cluster consists of 3 hosts
 Installed and started services successfully on 3 new hosts
 Master services installed
 NameNode installed on node-1.c.coastal-airlock-197705.internal
 SNameNode installed on node-2.c.coastal-airlock-197705.internal
 ResourceManager installed on node-2.c.coastal-airlock-197705.internal
 History Server installed on node-2.c.coastal-airlock-197705.internal
 All services started
 All tests passed
 Install and start completed in 8 minutes and 11 seconds

Complete ->

Ambari packt 0 ops 0 alerts

Dashboard Services Hosts Alerts Admin admin

- ✓ HDFS
- ✓ YARN
- ✓ MapReduce2
- ✓ ZooKeeper
- ✓ Ambari Metrics
- ✓ SmartSense

Actions

Metrics Heatmaps Config History

Metric Actions Last 1 hour

| | | | | |
|----------------------------|---------------------------------|---|-------------------------|------------------------------|
| HDFS Disk Usage | DataNodes Live | HDFS Links NameNode Secondary NameNode 1 DataNodes More... | Memory Usage | Network Usage |
| CPU Usage | Cluster Load | NameNode Heap | NameNode RPC | NameNode CPU WIO |
| NameNode Uptime | ResourceManager Heap | ResourceManager Uptime | YARN Memory | NodeManagers Live |

YARN Links
 ResourceManager
 1 NodeManagers
 More...

