



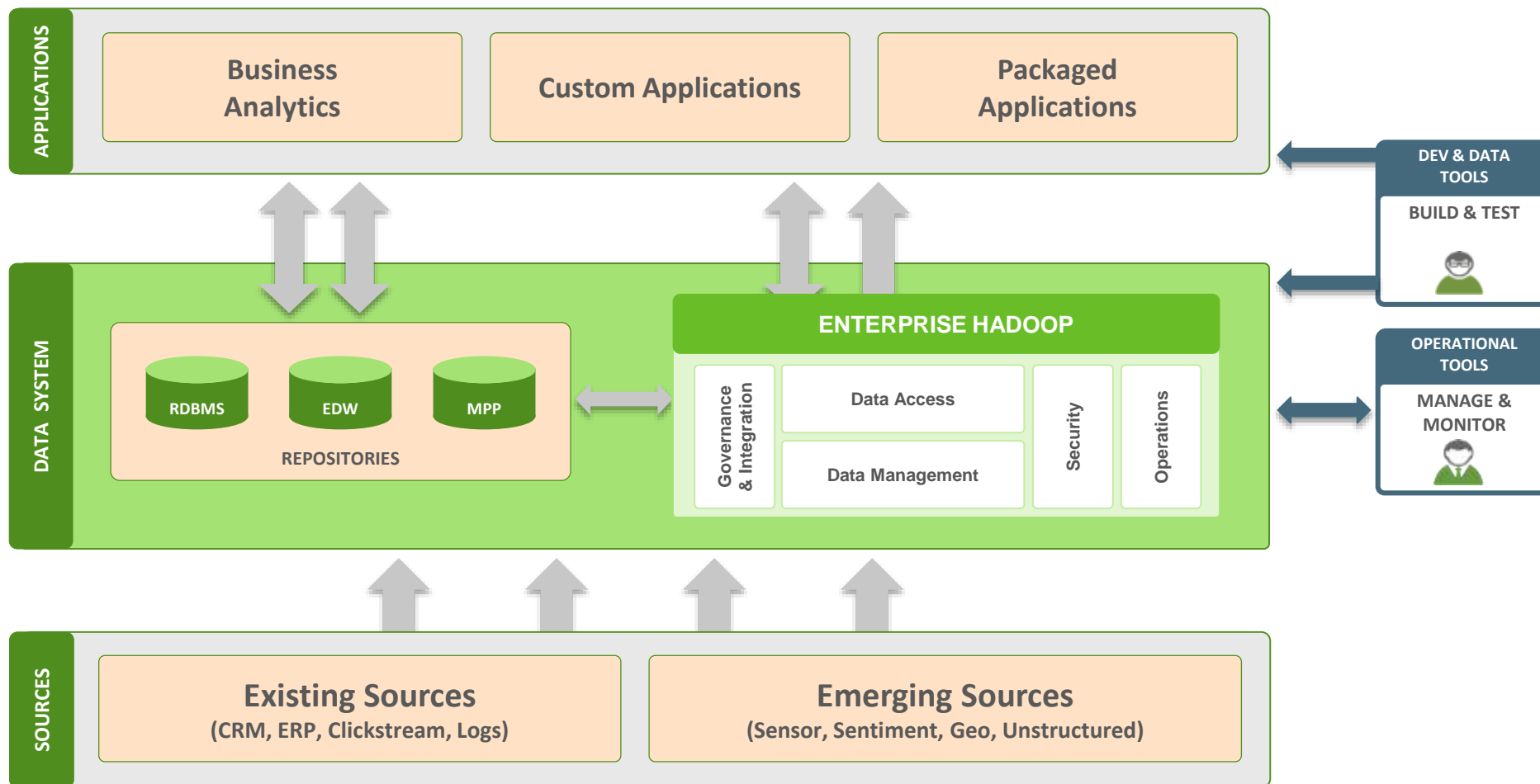
# Securing Hadoop Data Lake

Hortonworks. We do Hadoop.

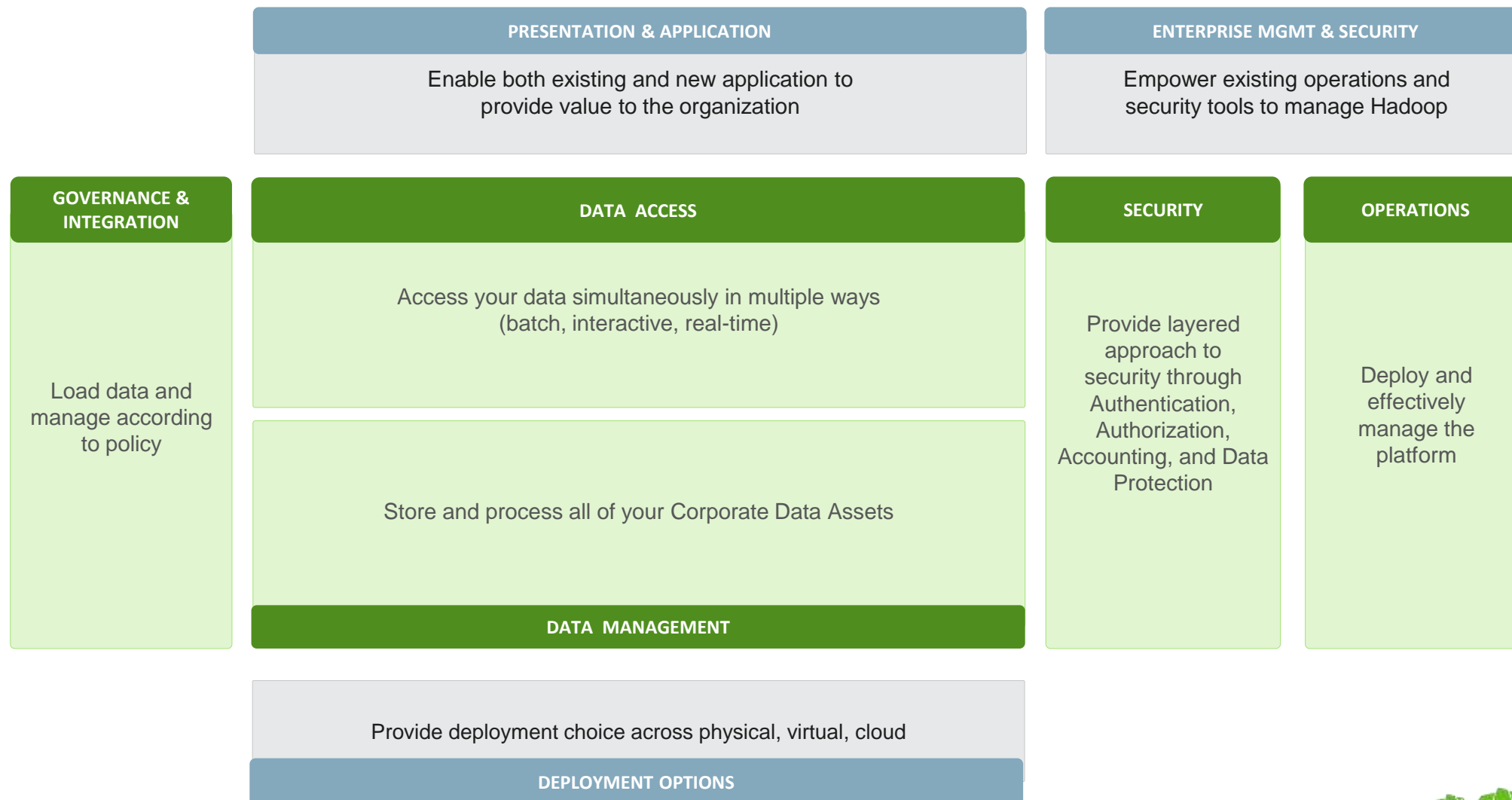
# Agenda

- **Security Approach within Hadoop**
- **Security Pillars**
- **Workshops**
- **Questions**

# A Modern Data Architecture



# Core Capabilities of Enterprise Hadoop



# Security needs are changing

## Administration

Centrally management & consistent security

## Authentication

Authenticate users and systems

## Authorization

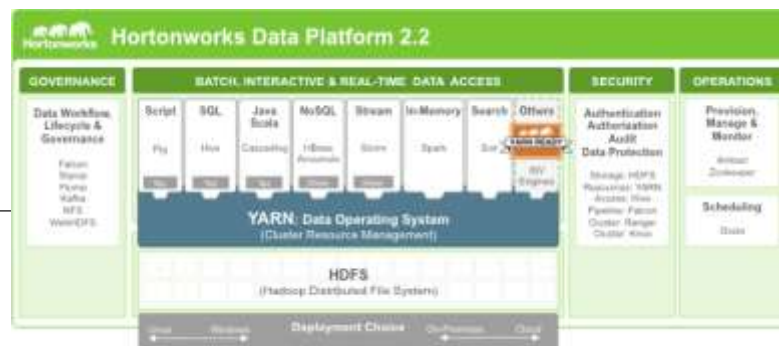
Provision access to data

## Audit

Maintain a record of data access

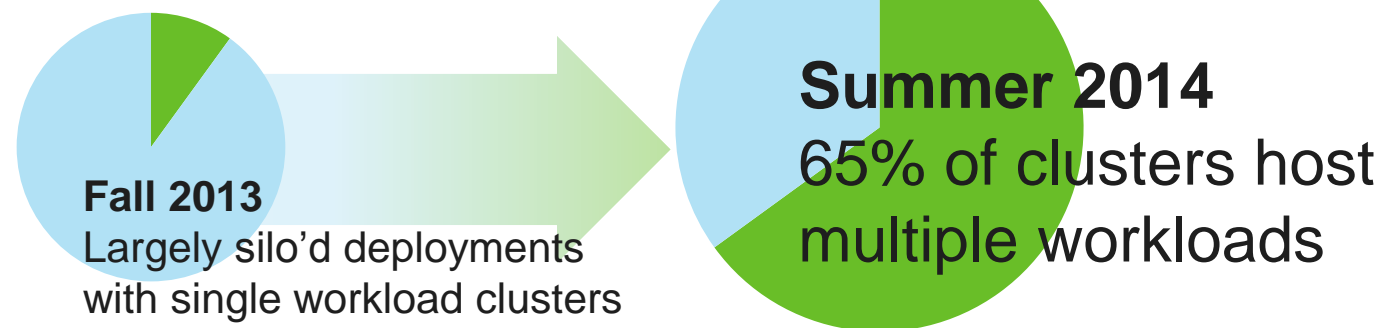
## Data Protection

Protect data at rest and in motion



## Security needs are changing

- YARN unlocks the data lake
- Multi-tenant: Multiple applications for data access
- Changing and complex compliance environment
- Data classification



# Security today in Hadoop with HDP

## Centralized Security Administration

### Authentication

Who am I/prove it?

- Kerberos in native Apache Hadoop
- HTTP/REST API Secured with Apache Knox Gateway

### Authorization

Restrict access to explicit data

- HDFS, Hive and Hbase (Storm and Knox in 2.2)
- Fine grain access control

### Audit

Understand who did what

- Centralized audit reporting
- Policy and access history

### Data Protection

Encrypt data at rest & in motion

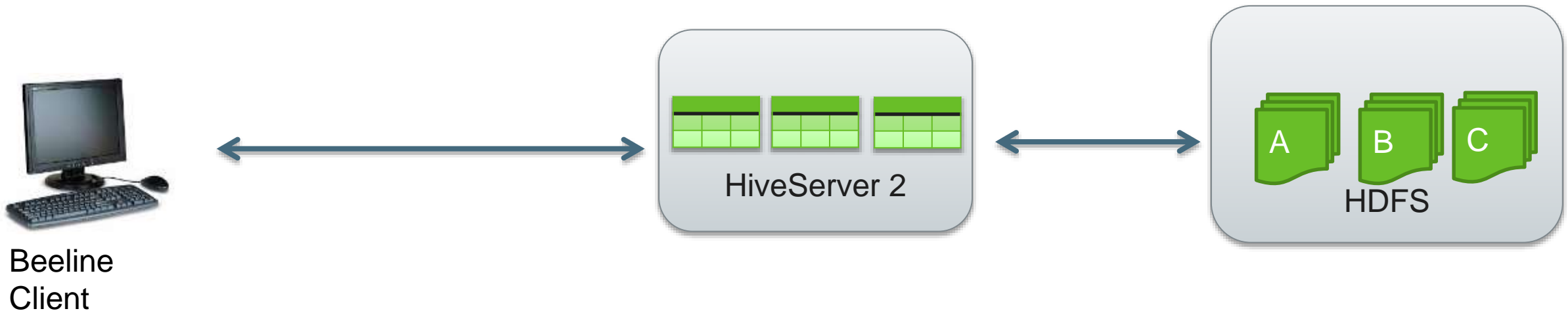
- Wire encryption in Hadoop
- Orchestrated encryption with partner tools

HDP 2.1

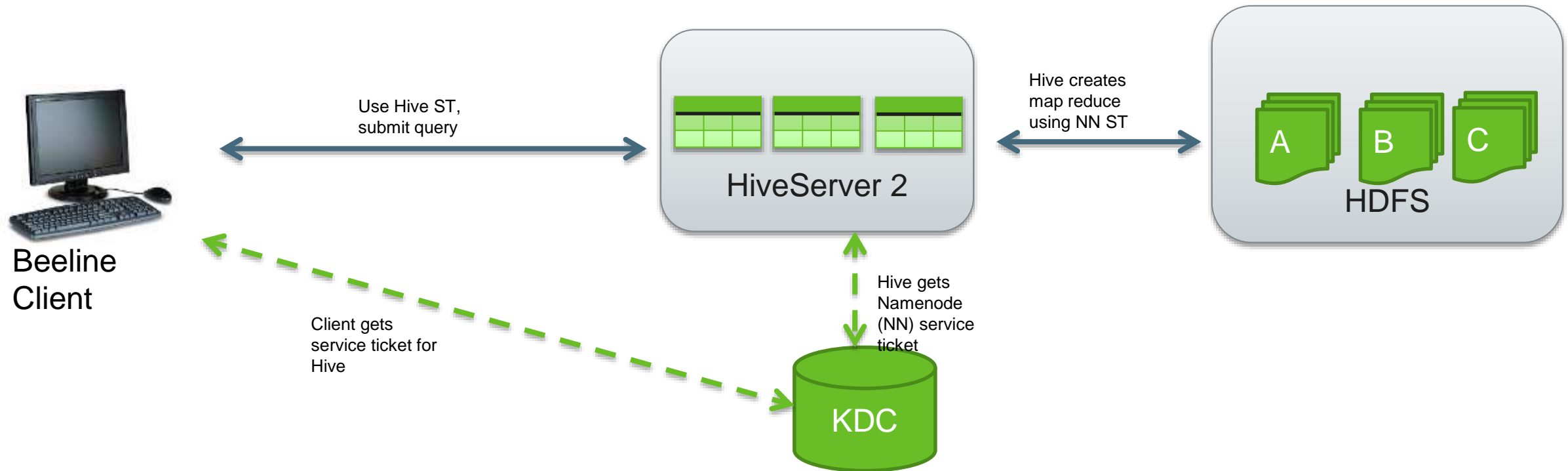


Ranger

# Typical Flow – Hive Access through Beeline client

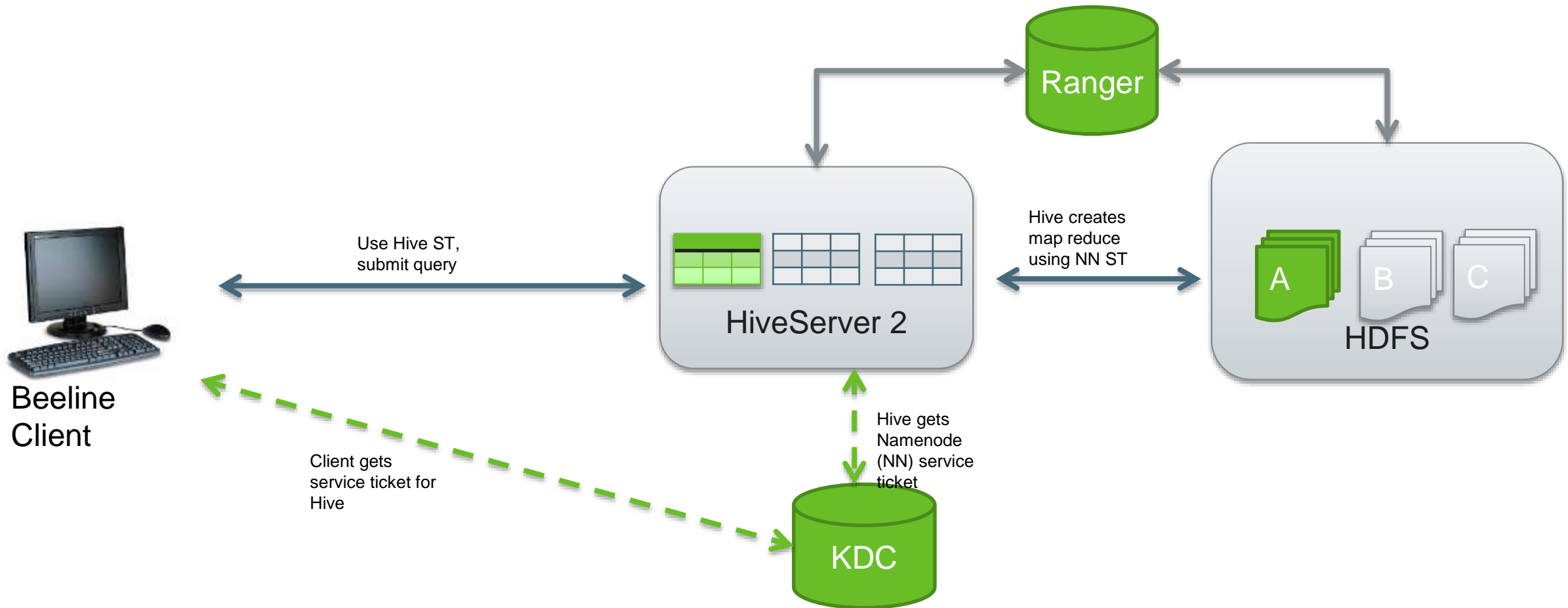


# Typical Flow – Authenticate through Kerberos

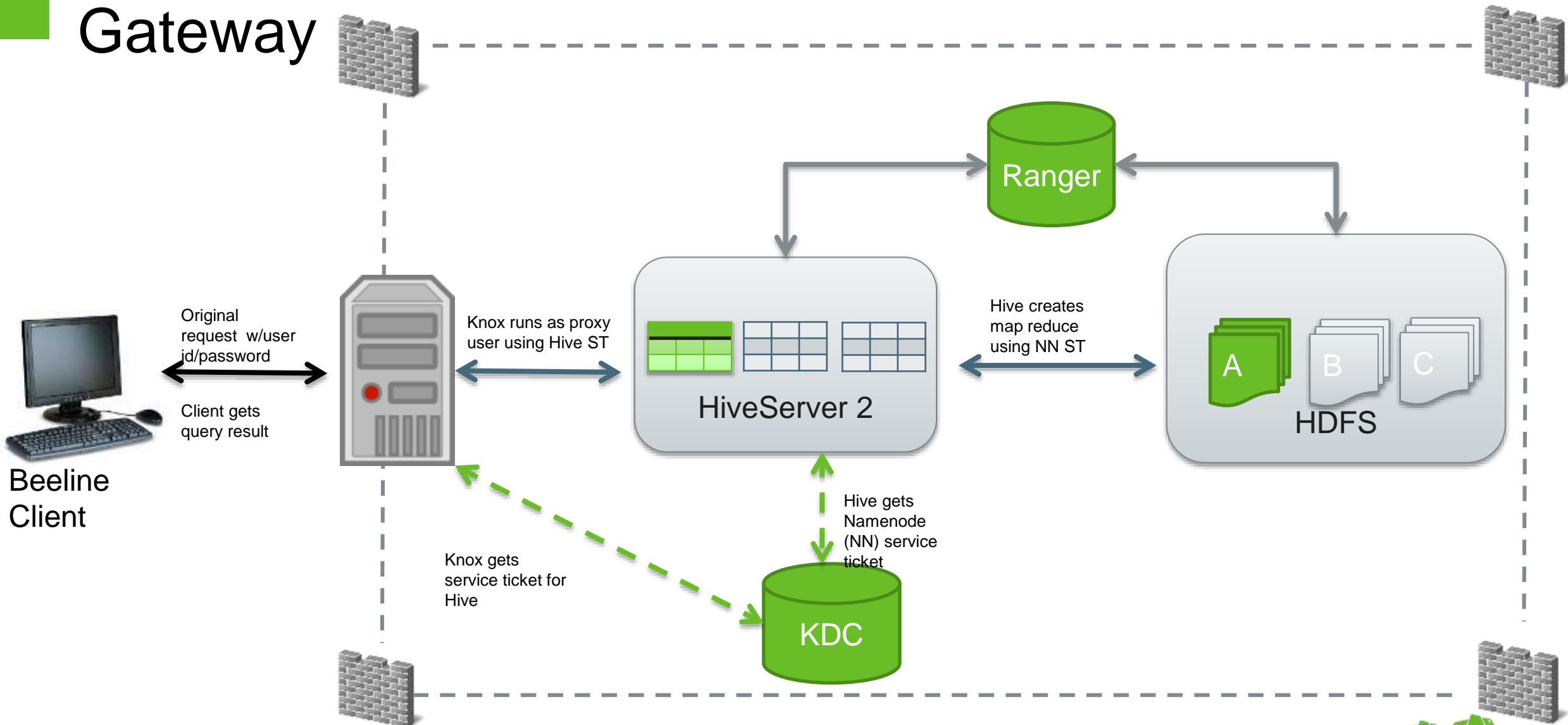




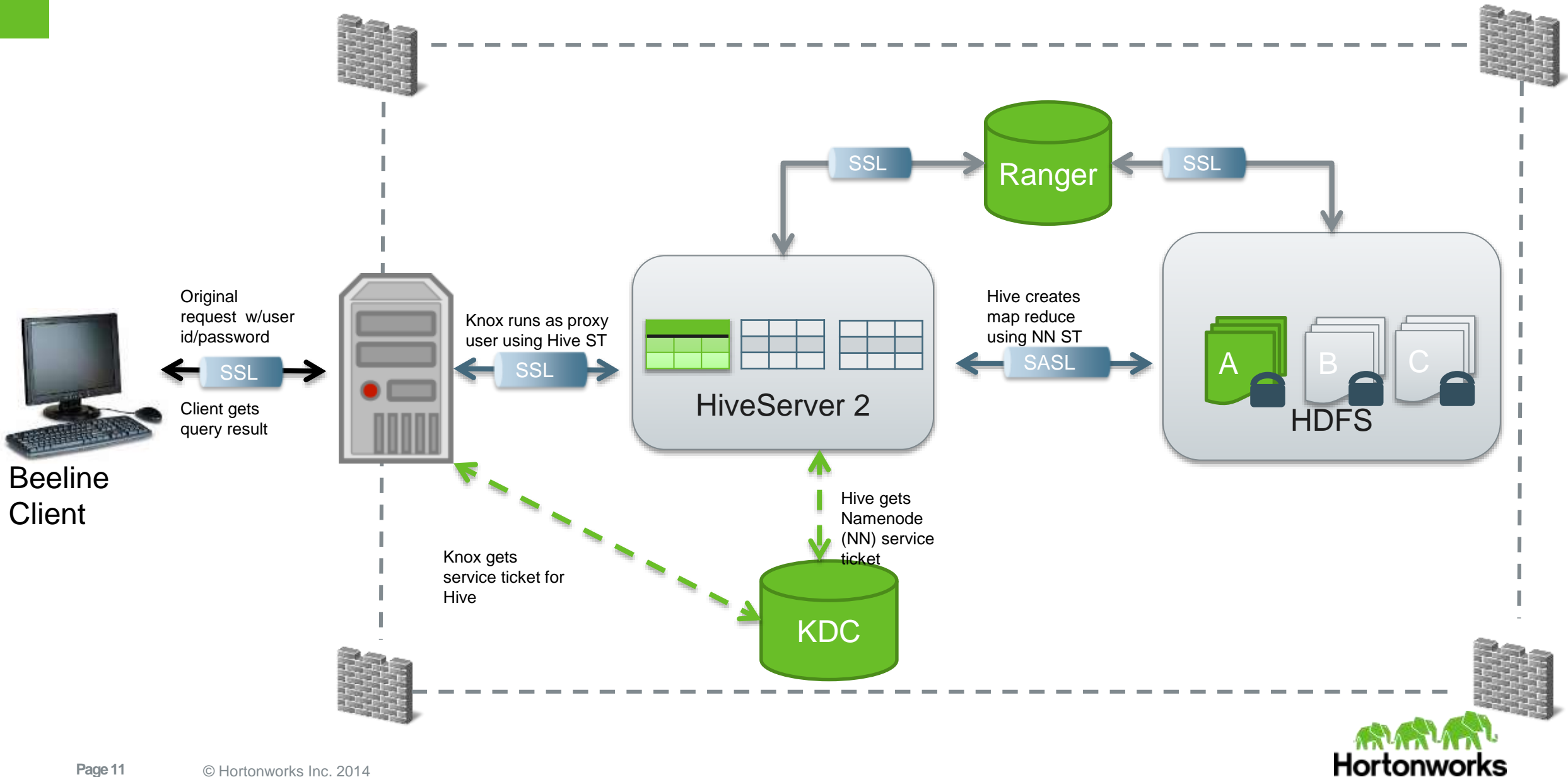
# Typical Flow – Add Authorization through Ranger(XA Secure)



# Typical Flow – Firewall, Route through Knox Gateway



# Typical Flow – Add Wire and File Encryption



# Security Features

HDP Security	
Authentication	
Kerberos Support	✓
Perimeter Security – For services and rest API	✓
Authorizations	
Fine grained access control	HDFS, Hbase and Hive, Storm and Knox (next release)
Role base access control	✓
Column level	✓
Permission Support	Create, Drop, Index, lock, user
Auditing	
Resource access auditing	Extensive Auditing
Policy auditing	✓

# Security Features

HDP Security	
<b>Data Protection</b>	
Wire Encryption	✓
Volume Encryption	✓
File/Column Encryption	HDFS TDE & Partners
<b>Reporting</b>	
Global view of policies and audit data	✓
<b>Manage</b>	
User/ Group mapping	✓
Global policy manager, Web UI	✓
Delegated administration	✓

# Authorization and Auditing

Apache Ranger

# Authorization and Audit

## Authorization

Fine grain access control

- HDFS – Folder, File
- Hive – Database, Table, Column
- HBase – Table, Column Family, Column

Flexibility  
in defining  
policies

## Audit

Extensive user access auditing in  
HDFS, Hive and HBase

- IP Address
- Resource type/ resource
- Timestamp
- Access granted or denied

Control  
access into  
system

# Central Security Administration

## HDP Advanced Security

- Delivers a 'single pane of glass' for the security administrator
- Centralizes administration of security policy
- Ensures consistent coverage across the entire Hadoop stack

The screenshot shows the 'Create Policy' page in the Hortonworks Policy Manager. The breadcrumb trail is 'Manage Repository > hivedev Policies > Edit Policy'. The 'Policy Details' section includes:

- Select Database Name \***: A text box containing 'xademo'.
- Table**: A dropdown menu set to 'Table'.
- Select Column Name**: A multi-select box containing 'x customer\_details', 'x phone\_number', 'x plan', 'x date', 'x status', 'x balance', and 'x region'.
- include**: Two toggle switches, both set to 'include'.
- Audit Logging**: A toggle switch set to 'ON'.

The screenshot shows the 'Manage Repository' page in the Hortonworks Policy Manager. The breadcrumb trail is 'Manage Repository'. The page displays three repository cards:

- HDFS**: Contains 'hadoopdev'.
- HIVE**: Contains 'hivedev'.
- HBASE**: Contains 'hbasedev'.

Each card has a green plus icon for adding and a red minus icon for removing. To the right of the repository cards, there are two permission matrices. Each matrix has columns for 'Drop', 'Alter', 'Index', 'Lock', 'All', and 'Admin', and rows for each repository. The 'All' and 'Admin' columns have checkboxes that are currently unchecked. A red 'x' icon is visible to the right of each matrix.



# Setup Authorization Policies

The screenshot shows the Hortonworks Policy Manager interface. The top navigation bar includes 'Hortonworks', 'Policy Manager', 'Users/Groups', 'Analytics', and 'Audit'. Below this, a breadcrumb trail shows 'Manage Repository' > 'hadoopdev Policies' > 'Edit Policy'. The main section is titled 'Create Policy' and contains two sub-sections: 'Policy Details' and 'User and Group Permissions'.

**Policy Details:**

- Resource Path \***: A text field containing '/demo/data/Website/Website-Logs'.
- Description**: An empty text area.
- Recursive**: A toggle switch set to 'YES'.
- Audit Logging**: A toggle switch set to 'ON'.






**User and Group Permissions:**

Select Group	Read	Write	Execute	Admin	
<input type="text" value="IT"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="button" value="x"/>
<input type="text" value="Network"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="button" value="x"/>
<input type="button" value="+"/>					

file level  
access  
control,  
flexible  
definition

Control  
permissions

# Monitor through Auditing


 **Policy Manager**  **Users/Groups**  **Analytics**  **Audit**  **admin**

Big Data

Admin

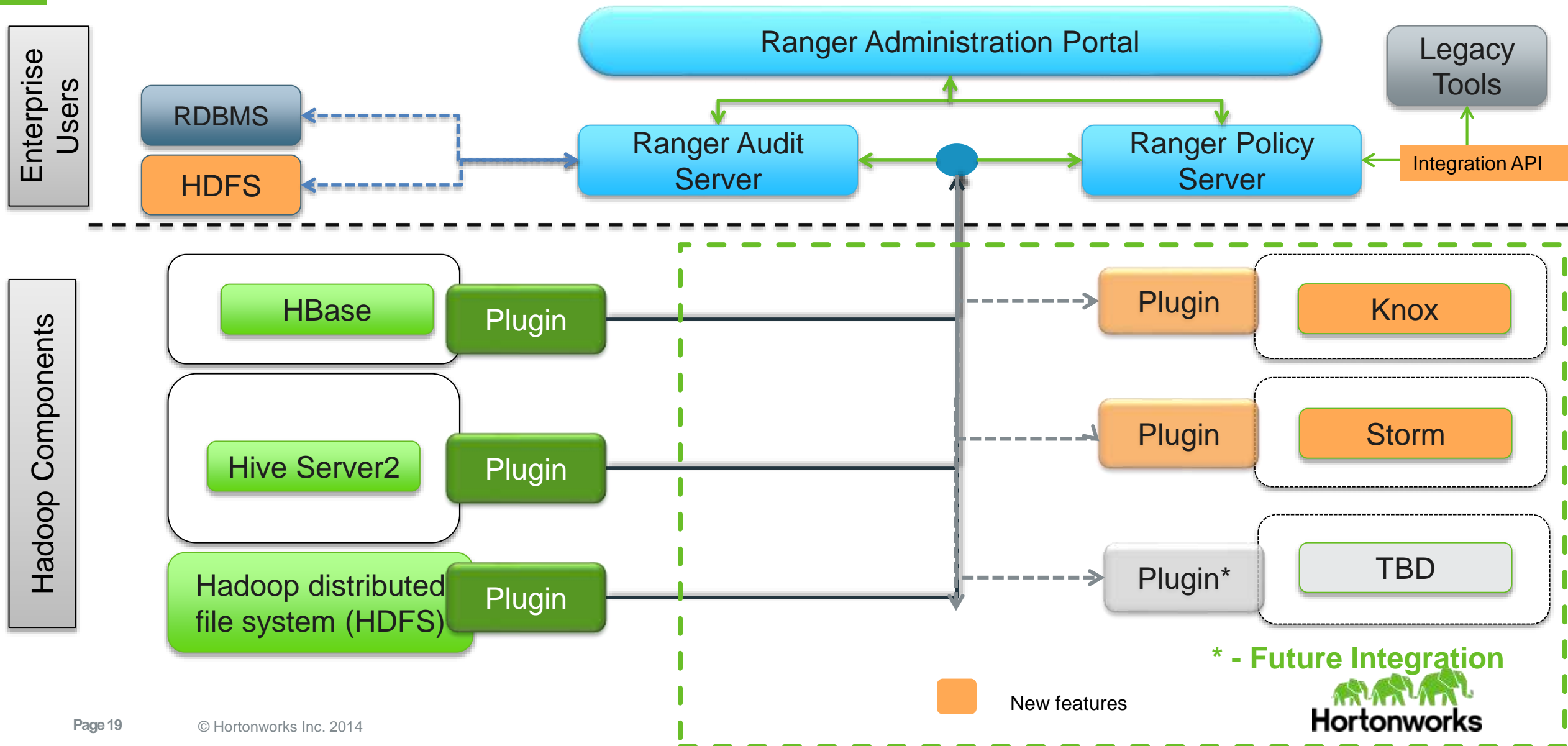
Login Sessions

Agents

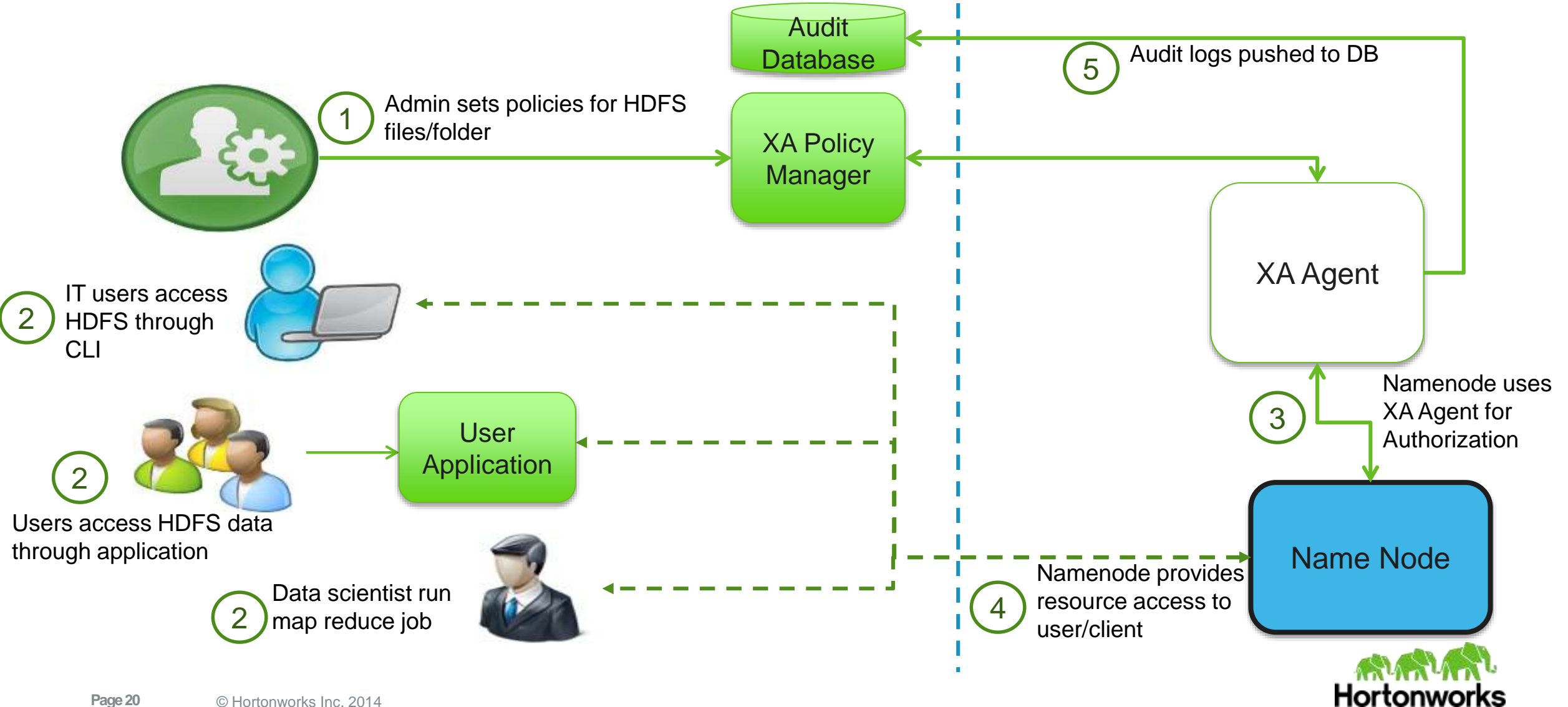
Last Updated Time : 06/20/2014 05:32:12 PM 

Event Time	User	Repository	Resource Name	Access Type	Result	Access Enforcer	Client IP
		Name / Type					
06/20/2014 05:31:53 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:31:30 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:30:30 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:29:30 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:28:53 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:28:30 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:27:30 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15
06/20/2014 05:26:31 PM	mapred	hadoopdev HDFS	/mr-history/tmp	READ_EXECUTE	Allowed	hadoop-acl	/10.0.2.15

# Authorization and Auditing w/ Ranger



# Simplified Workflow - HDFS



# Ranger Investments for HDP 2.2

- **New Components Coverage**
  - Storm Authorization & Auditing
  - Knox Authorization & Auditing
- **Deeper Integration with HDP**
  - Windows Support
  - Integration with Hive Auth API, support grant/revoke commands
  - Support grant/revoke commands in Hbase
- **Enterprise Readiness**
  - Rest APIs for policy manager
  - Store Audit logs locally in HDFS
  - Support Oracle DB
  - Ambari support, as part of Ambari 2.0 release

# REST API Security through Knox

Securely share Hadoop Cluster

# Hadoop REST API with Knox

Service	Direct URL	Knox URL
WebHDFS	<a href="http://namenode-host:50070/webhdfs">http://namenode-host:50070/webhdfs</a>	<a href="https://knox-host:8443/webhdfs">https://knox-host:8443/webhdfs</a>
WebHCat	<a href="http://webhcat-host:50111/templeton">http://webhcat-host:50111/templeton</a>	<a href="https://knox-host:8443/templeton">https://knox-host:8443/templeton</a>
Oozie	<a href="http://ooziehost:11000/oozie">http://ooziehost:11000/oozie</a>	<a href="https://knox-host:8443/oozie">https://knox-host:8443/oozie</a>
HBase	<a href="http://hbasehost:60080">http://hbasehost:60080</a>	<a href="https://knox-host:8443/hbase">https://knox-host:8443/hbase</a>
Hive	<a href="http://hivehost:10001/cliservice">http://hivehost:10001/cliservice</a>	<a href="https://knox-host:8443/hive">https://knox-host:8443/hive</a>
YARN	<a href="http://yarn-host:yarn-port/ws">http://yarn-host:yarn-port/ws</a>	<a href="https://knox-host:8443/resourcemanager">https://knox-host:8443/resourcemanager</a>

Masters could  
be on many  
different hosts

One hosts,  
one port

SSL config  
at one host

Consistent  
paths

# Why Knox?

## Enhanced Security

- Protect network details
- SSL for non-SSL services
- WebApp vulnerability filter

## Centralized Control

- Central REST API auditing
- Service-level authorization
- Alternative to SSH “edge node”

## Simplified Access

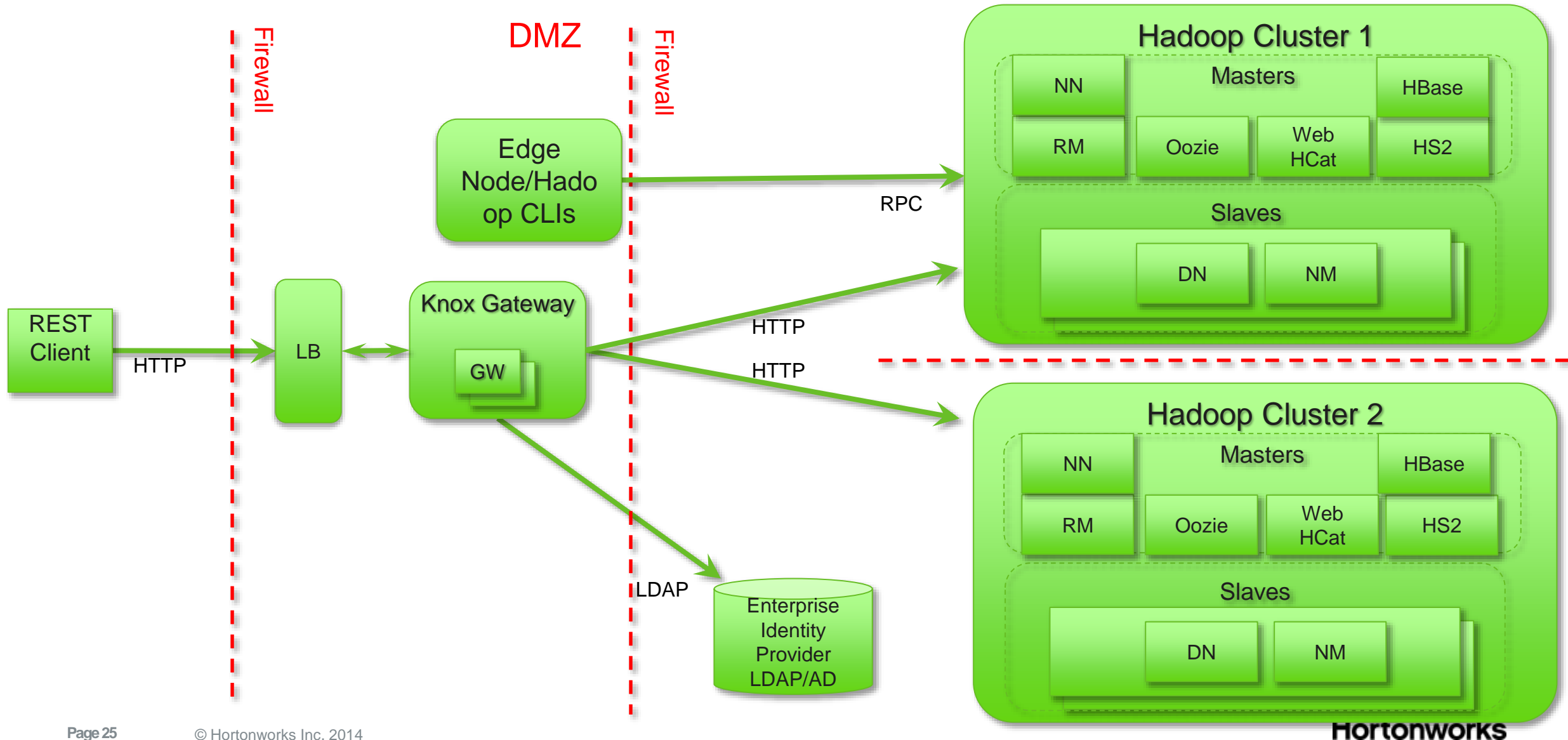
- Kerberos encapsulation
- Extends API reach
- Single access point
- Multi-cluster support
- Single SSL certificate

## Enterprise Integration

- LDAP integration
- Active Directory integration
- SSO integration
- Apache Shiro extensibility
- Custom extensibility



# Hadoop REST API Security: Drill-Down



# What's New in Knox with HDP 2.2

- Use Ambari for Install/start/stop/configuration
- Knox support for HDFS HA
- Support for YARN REST API
- Support for SSL to Hadoop Cluster Services (WebHDFS, HBase, Hive & Oozie)
- Knox Management REST API
- Integration with Ranger (fka XA Secure) to for Knox Service Level Authorization

# Workshop: Enabling Security

# Let's Begin

- We will use HDP Sandbox with FreeIPA Software Installed
  - FreeIPA is an integrated security information management solution combining Linux (Fedora), 389 Directory Server, MIT Kerberos, NTP, DNS, Dogtag (Certificate System). It consists of a web interface and command-line administration tools
  - In the workshop we use FreeIPA for User Identity Management
  - **Note:** Steps outlined in the workshop are applicable for other identity management solutions such as Active Directory

# Authentication

## 1. Create end users and groups in FreeIPA

- End Users will query HDP via Hue, Beeline & JDBC/ODBC clients

## 2. Enable Kerberos for the HDP Cluster

- Hadoop now authenticates all access to the cluster

## 3. Integrate Hue with FreeIPA

- Users are validated against FreeIPA

## 4. Configure Linux to use FreeIPA as central store of posix data using nslcd

- Enables Hadoop to determine user groups without requiring a local linux user account

## We have now set Authentication

- A user can open a shell, authenticate using kinit and submit hadoop commands or alternatively log into HUE to access Hadoop.

# Enable Perimeter Security

## 1. KNOX Is Available on Sandbox

- Enables Perimeter Security. Enables single point of cluster access using Hadoop REST APIs, JDBC and ODBC calls

## 2. Configure KNOX to authenticate against FreeIPA

## 3. Configure WebHDFS & Hiveserver2 to support JDBC/ODBC access over HTTP

## 4. Use Excel to access Hive via KNOX

- Note, Knox eliminates the need to secure Kerberos ticket on the client machine for user authentication

## We have now set Perimeter Security

- Users can now access the cluster via the Gateway services

# Authorization & Audit

## 1. Install Apache Ranger

- Comprehensive authorization and audit tool for Hadoop

## 2. Sync users between Apache Ranger and FreeIPA

- Note, end users are only required to be maintained in **one** enterprise identity management system

## 3. Configure HDFS & Hive to use Apache Ranger

- In this workshop we will only show steps as it relates to hive authorization. Similar capabilities are available for other HDP components.

## 4. Define HDFS & Hive Access Policy For Users

- User “hive” is a special user and must be assigned universal access

## 5. Log into Hue as the end user and note the authorization policies being enforced

- Review Audit Information

## We have now set Authorization & Audit

- All user access to a Hive is governed & audited by policies maintained in Apache Ranger.

# Encryption

## 1. Wire Level Encryption

- Follow instruction here [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.0.6.0/bk\\_reference/content/ch\\_wire6.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.0.6.0/bk_reference/content/ch_wire6.html)

## 2. Volume Level Encryption

- Leverage LUKS. Sample script provided

## 3. Column level encryption & data masking

- Collaborate with our key security partners




# Resources

# Security Page

Apache Ranger

hortonworks.com/hadoop/ranger/

Apps Hortonworks EC2 test environme Run-HDP-Installer Run-Ambari-S



ENTERPRISE

WHAT IS HADOOP? • APACHE RANGER

Apache Ranger

Comprehensive security for Enterprise Hadoop

Ranger delivers a comprehensive approach to central security policy administration across the Hadoop ecosystem, including authentication, authorization, accounting and auditing.

It already extends baseline features for coordinated enforcement of security policies consistently against additional Hadoop ecosystem components like Hive, Pig, Tez, Storm, Solr, Spark, and more. It truly represents a major step towards a comprehensive approach to security – all completely as open source.


Hortonworks Investment Themes for Ranger

Investment Theme	Planned Enhancements
Extension of support	Additional investments extend administration of authorization and auditing to more

Apache Knox Gateway

hortonworks.com/hadoop/knox-gateway/

Apps Hortonworks EC2 test environme Run-HDP-Installer Run-Ambari-Setup EC2 test environme VPN



ENTERPRISE HADOOP SOLUTIONS PRODUCTS LABS TRAINING GET STARTED PARTNERS ABOUT BLOG

WHAT IS HADOOP? • APACHE KNOX GATEWAY

Apache Knox Gateway

A single point of secure access for Hadoop clusters


The Knox Gateway ("Knox") is a system that provides a single point of authentication and access for Apache™ Hadoop® services in a cluster. The goal of the project is to simplify Hadoop security for users who access the cluster data and execute jobs, and for operators who control access and manage the cluster. Knox runs as a server (or cluster of servers) that serve one or more Hadoop clusters.


What Knox Gateway Does

Knox Gateways provides security for multiple Hadoop clusters, with these advantages:


- **Provide perimeter security** to make Hadoop security setup easier
- **Support authentication and token verification** security scenarios
- **Deliver users a single cluster end-point** that aggregates capabilities for data and jobs
- **Enable integration** with enterprise and cloud identity management environments

Resources


 [Discover HDP2.1: Apache Storm for Stream Data](#)



APACHE TOP-LEVEL PROJECT SINCE  
February 2013  
HORTONWORKS COMMITTEES  
12  
PROJECT PAGE  
<http://knox.incubator.apache.org/>



Try Knox Gateway with Sandbox  
Hortonworks Sandbox is a self-



Hortonworks

# Hortonworks Security Investment Plans

## HDP + XA

### Comprehensive Security for Enterprise Hadoop

Goals:

#### Comprehensive Security

Meet all security requirements across Authentication, Authorization, Audit & Data Protection for all HDP components

#### Central Administration

Provide one location for administering security policies and audit reporting for entire platform

#### Consistent Integration

Integrate with other security & identity management systems, for compliance with IT policies

...all **IN** Hadoop

#### Investment themes

##### Previous Phases

- ✓ Kerberos Authentication
- ✓ HDFS, Hive & Hbase authorization
- ✓ Wire Encryption for data in motion
- ✓ Knox for perimeter security
- ✓ Basic Audit in HDFS & MR
- ✓ SQL Style Hive Authorization
- ✓ ACLs for HDFS

Delivered

##### XA Secure Phase

- Centralized Security Admin for HDFS, Hive & HBase
- Centralized Audit Reporting
- Delegated Policy Administration

Delivered XA  
Secure

##### Future Phases

- Encryption in HDFS, Hive & Hbase
- Centralized security administration of entire Hadoop platform
- Centralized auditing of entire platform
- Expand Authentication & SSO integration choices
- Tag based global policies (e.g. Policy for PII)

# Q&A