

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

A Practitioner's Guide to Securing Your Hadoop Cluster

Your Speakers

- André Araújo
 - Senior Solutions Architect, Cloudera
- Syed Rafice
 - Senior Systems Engineer, Cloudera
- Mubashir Kazia
 - Senior Solutions Architect, Cloudera
- Mark Donsky
 - Director Product Management, Cloudera

Format

- Five sections
- Each section:
 - Introduce a security concept
 - Demo without it in place
 - How to enable
 - Demo with it in place
- **Please hold questions until the end of each section**
- Short break in the middle
- Slides are available from <http://strataconf.com>

Agenda

- Prelude: Network Security – André
- Authentication – André
- Authorization – André
- Wire Encryption – Syed
- Encryption-at-rest – Mubashir
- Data Governance – Mark
- Final Thoughts – Mubashir

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

Prelude: Network Security

Don't Put Your Hadoop Cluster on the Open Internet

- MongoDB ransomware
 - Tens of thousands of open MongoDB instances on the internet
 - With no security turned on
 - The attack: All data deleted or encrypted; ransom note left behind
- It has happened to Hadoop clusters, too



SHODAN

cloudera



Explore

Downloads

Reports



Exploits



Maps



Share Search

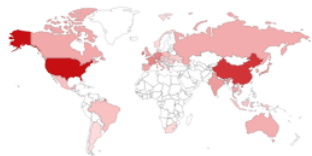


Download Results



Create Report

TOP COUNTRIES



United States	541
China	268
Singapore	46
Ireland	44
France	39

TOP SERVICES

50070	786
8086	269
HTTPS	51
Splunk	29
HTTP	20

TOP ORGANIZATIONS

Amazon.com	305
Google Cloud	68
Microsoft Azure	67

Total results: 1,182

Hadoop Administration

191.234.183.82

Microsoft Azure

Added on 2017-03-07 21:03:55 GMT

Brazil, Campinas

[Details](#)14.0
MB41
Files

Total Blocks 2

Number of Threads 177

Hadoop Administration

104.199.102.238

238.102.199.104.bc.googleusercontent.com

Google Cloud

Added on 2017-03-07 20:52:34 GMT

United States, Mountain View

[Details](#)510.0
MB3,931
Files

Total Blocks 661

Number of Threads 96

Basic Networking Checks

- Make sure your IP address isn't an internet-exposed address
 - These are the private IP address ranges:
 - 10.* (10.0/8)
 - 172.16.* - 172.31.* (172.16/12)
 - 192.168.* (192.168/16)
- Use `nmap` from outside your corporate environment
- If in {AWS, Azure, GCE}, check networking configuration



strataconf.com
#StrataData

PRESENTED BY



Questions?

Authentication

André Araújo

Senior Solutions Architect
Cloudera

Authentication - Agenda

- Intro - identity and authentication
- **DEMO:** Hadoop with no authentication
- Kerberos and LDAP authentication
- Enabling kerberos and LDAP using Cloudera Manager
- **DEMO:** Actual strong authentication in Hadoop
- Questions

Identity

- Before we can talk about authentication, we must understand **identity**
- An object that uniquely identifies a user (usually)
 - Email account, Windows account, passport, driver's license
- In Hadoop, identity largely means **username**
- Using a common source of identity is paramount

Identity Sources

- Individual Linux servers use `/etc/passwd` and `/etc/group`
 - Not scalable and prone to **errors**
- LDAP is the preferred way
 - Integrate at the Linux OS level
 - RedHat SSSD
 - Centrify
 - **All** applications running on the OS can use the same LDAP integration
 - Most enterprises use Active Directory
 - Some enterprises use a Linux-specific LDAP implementation

Identity and Authentication

- So you have an identity database, now what?
- Users and applications must **prove** their identities to each other
- This process is authentication
- Hadoop strong authentication is built around **Kerberos**
- Kerberos is built into Active Directory and this is the most common Hadoop integration

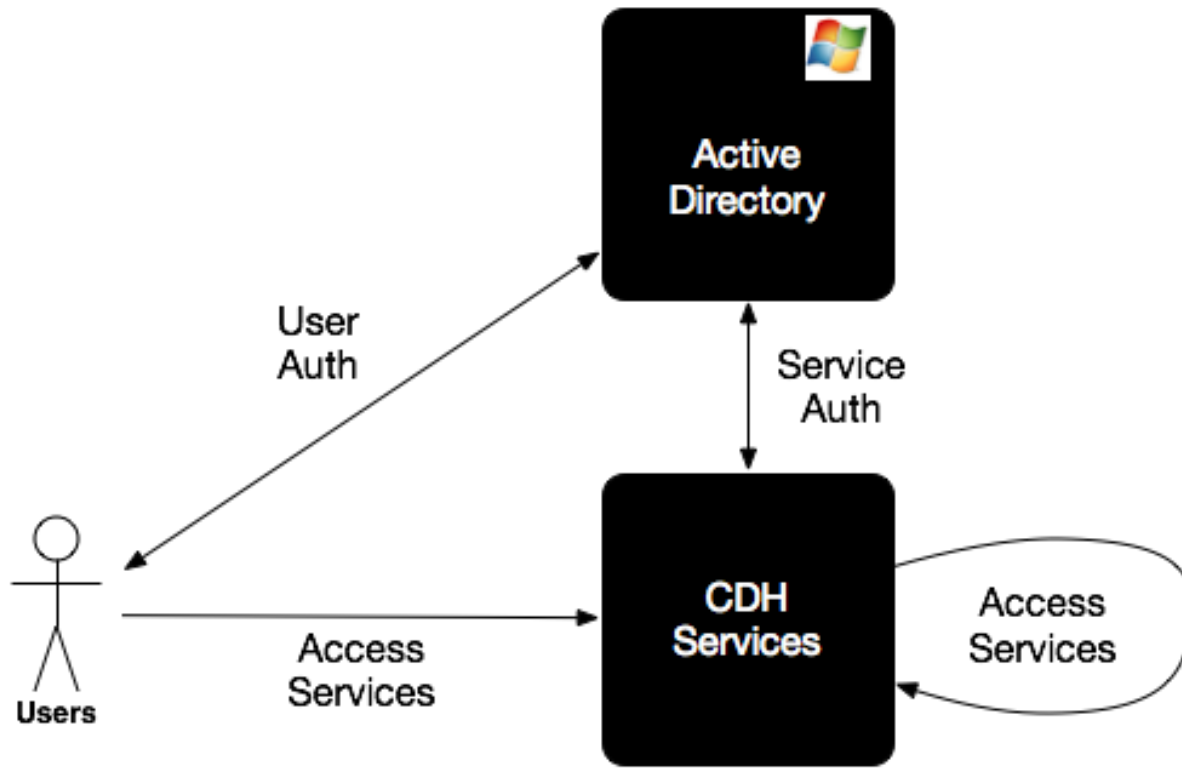
Hadoop Default “Authentication”

- Out of the box, Hadoop “authenticates” users by simply believing whatever username you tell it you are
- This includes telling Hadoop you are the hdfs user, a **superuser**!
- **DEMO:** Let’s see just how bad this is.

Kerberos

- To enable security in Hadoop, everything starts with Kerberos
- **Every role type of every service has its own unique Kerberos credentials**
- Users must **prove** their identity by obtaining a Kerberos ticket, which is honored by the Hadoop components
- Hadoop components themselves authenticate to each other for intra and inter service communication

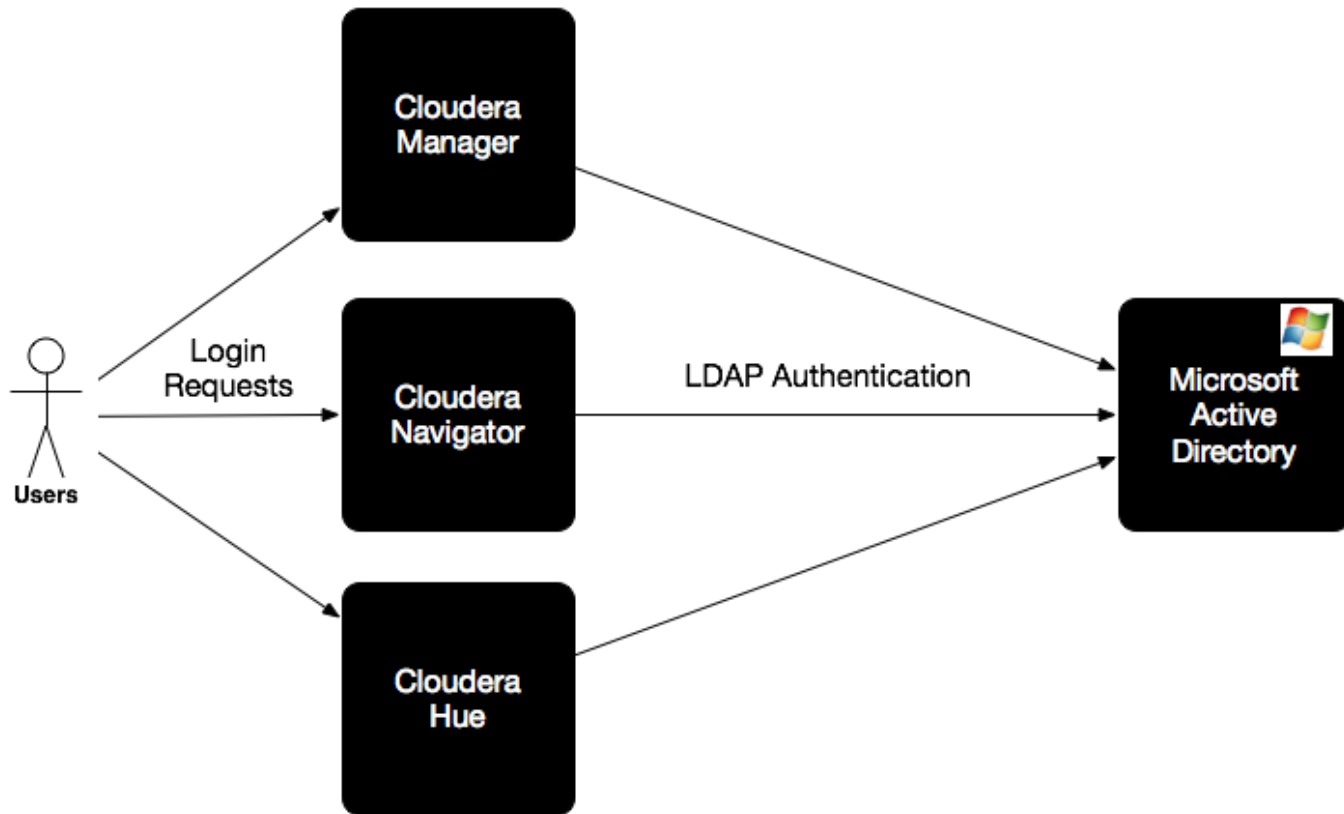
Kerberos Authentication



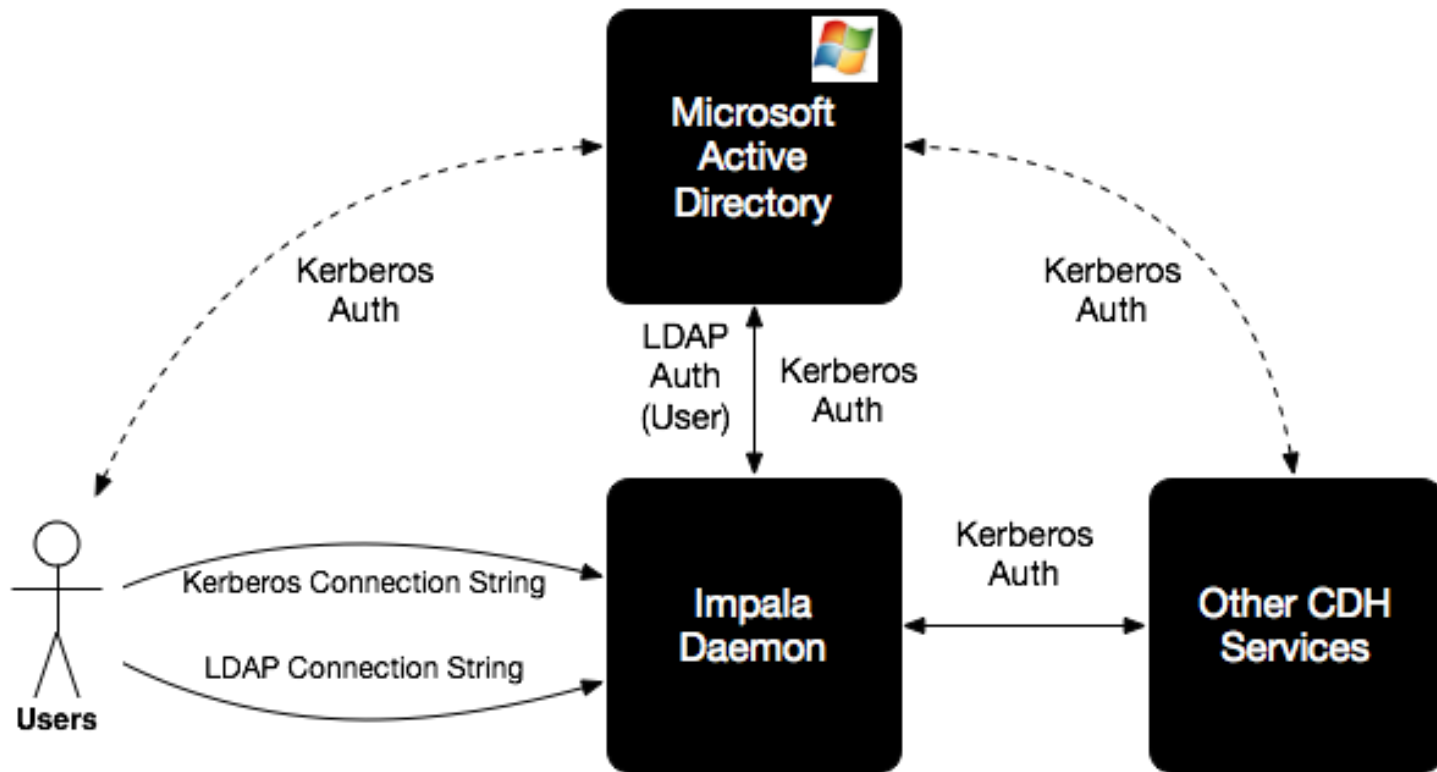
LDAP and SAML

- Beyond just Kerberos, other components such as web consoles and JDBC/ODBC endpoints can authenticate users differently
- **LDAP** authentication is supported for Hive, Impala, Solr, and web-based UIs
- **SAML** (SSO) authentication is supported for Cloudera Manager, Navigator, and Hue
- Some components support both Kerberos and LDAP authentication at the same time
- Generally speaking, LDAP is a much easier authentication mechanism to use for external applications – No Kerberos software and configuration required!
- **...just make sure wire encryption is also enabled to protect passwords**

Web UI LDAP Authentication



Impala Dual-mode Authentication



Enabling Kerberos

- Setting up Kerberos for your cluster is no longer a daunting task
- Cloudera Manager and Apache Ambari provide wizards to automate the provisioning of service accounts and the associated keytabs
- Both MIT Kerberos and Active Directory are supported Kerberos KDC types
- Again, most enterprises use Active Directory so let's see what we need to set it up!

Active Directory Prerequisites

- At least one AD domain controller is setup with LDAPS
- An AD account for Cloudera Manager
- A **dedicated OU** in your desired AD domain
- An account that has **create/modify/delete** user privileges on this OU
- This is **not** a domain admin / administrative account!
- While not required, AD **group policies** can be used to further restrict the accounts
- Install **openldap-clients** on the CM server host, **krb5-workstation** on every host
- From here, use the wizard!

Cloudera Manager Kerberos Wizard

Before using the wizard, please ensure that you have performed the following steps:

Set up a working KDC. Cloudera Manager supports MIT KDC and Active Directory.

☒ Yes, I've set up a working KDC.

The KDC should be configured to have non-zero ticket lifetime and renewal lifetime. CDH will not work properly if tickets are not renewable.

☒ Yes, I've checked that the KDC allows renewable tickets.

OpenLdap client libraries should be installed on the Cloudera Manager Server host if you want to use Active Directory. Also, Kerberos client libraries should be installed on ALL hosts.

☒ Yes, I've installed the client libraries.

Cloudera Manager needs an account that has permissions to create other accounts in the KDC.

☒ Yes, I've created a proper account for Cloudera Manager.

KDC Information

Specify information about the KDC. The properties below are used by Cloudera Manager to generate principals for CDH daemons running on the cluster.

KDC Type

- ☐ MIT KDC [C](#)
- ☒ Active Directory

KDC Server Host

kdc

ad.hadoop.com

[C](#)

Kerberos Security Realm

default_realm

HADOOP.COM

Kerberos Encryption Types

aes256-cts

+ - [C](#)

aes128-cts

+ -

rc4-hmac

+ -

Active Directory Suffix

ou=hadoop,DC=hadoop,DC=com

Active Directory Account Prefix

cdh_

[C](#)

Active Directory Domain Controller Override

my-ad-dc1.hadoop.com

[C](#)

Cloudera Manager Kerberos Wizard

KDC Account Manager Credentials

Enter the credentials for the account that has permissions to **create** other users. Cloudera Manager will store it in encrypted form and use it whenever new principals need to be generated.

Username @
Password

Click through the remaining steps

Setting up LDAP Authentication

- CM -> Administration -> Settings
 - Click on category “External Authentication”
- Cloudera Management Services -> Configuration
 - Click on category “External Authentication”
- Hue / Impala / Hive / Solr -> Configuration
 - Search for “LDAP”

Post-Configuration

- Kerberos authentication is setup
- LDAP authentication is setup
- **DEMO:** No more fake authentication!



strataconf.com
#StrataData

PRESENTED BY



Questions?

Authorization

André Araújo

Senior Solutions Architect
Cloudera

Authorization - Agenda

- Authorization – Overview
- **DEMO:** Default Authorization
- Configuration Stronger Authorization
- Apache Sentry
- Record Service
- **DEMO:** Strong Authorization
- Questions

Authorization - Overview

- Authorization dictates what a user is permitted to do
- Happens **after** a user has authenticated to establish identity
- Authorization policies in Hadoop are typically based on:
 - Who the **user** is and what **groups** they belong to
 - Role-based access control (RBAC)
- Many different authorization mechanisms in Hadoop components

Authorization in Hadoop




- HDFS file permissions (POSIX 'rwx rwx rwx' style)
- Yarn job queue permissions
- Sentry (Hive / Impala / Solr / Kafka)
- Cloudera Manager RBAC
- Cloudera Navigator RBAC
- Hue groups
- Hadoop KMS ACLs
- HBase ACLs
- etc.

Default Authorization Examples

- HDFS
 - Default umask is 022, making all new files **world readable**
 - Any authenticated user can execute hadoop shell commands
- YARN
 - Any authenticated user can submit and **kill jobs** for any queue
- Hive metastore
 - Any authenticated user can **modify the metastore** (CREATE/DROP/ALTER/etc.)
- **DEMO:** Let's see just how bad this is.

Configuring HDFS Authorization

- Set default umask to 026
- Setup hadoop-policy.xml (Service Level Authorization)

Default Umask dfs.umaskmode, fs.permissions.umask-mode	HDFS-1 (Service-Wide) 
	<input type="text" value="026"/>
Authorized Groups	HDFS-1 (Service-Wide) 
	<input type="text" value="prod_cdh_users"/>
Authorized Admin Groups	HDFS-1 (Service-Wide) 
	<input type="text" value="prod_cdh_admins"/>

Configuring Yarn Authorization

- Setup the YARN admin ACL

Admin ACL

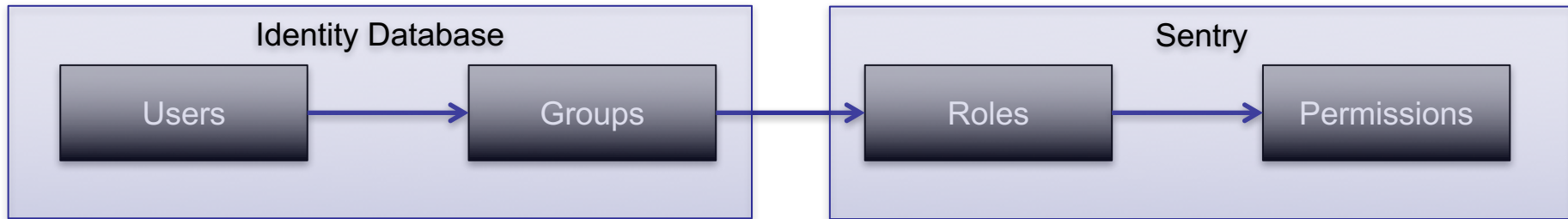
yarn.admin.acl

YARN-1 (Service-Wide) C

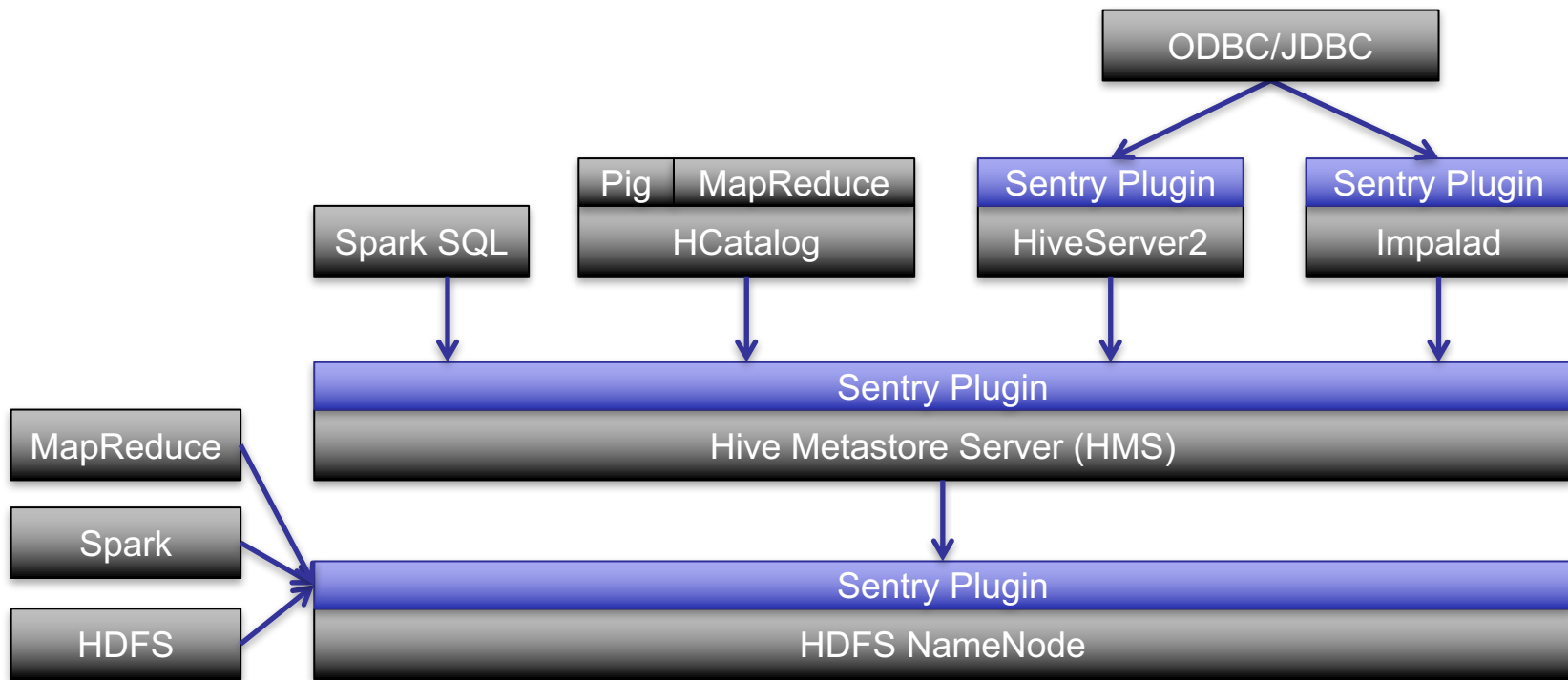
yarn prod_cdh_admins

Apache Sentry

- Provides **centralized RBAC** for several components
 - **Hive / Impala:** Databases, tables, views, columns
 - **Solr:** Collections, documents, indexes
 - **Kafka:** Cluster, topic, consumer group



Apache Sentry (Cont.)

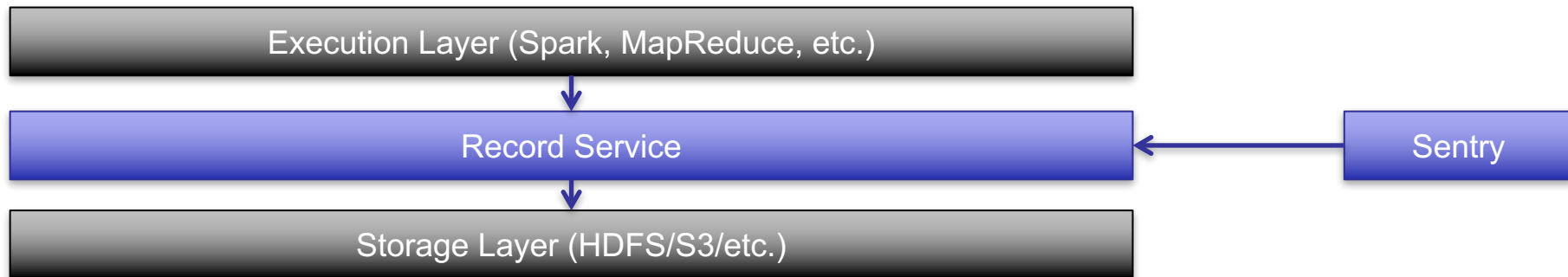


Configuring Sentry

- Cloudera Manager -> Add Service -> Sentry
- Hive
 - Set Sentry service
 - Disable HiveServer2 impersonation
- Impala
 - Set Sentry Service
- HDFS
 - Enable Sentry HDFS Synchronization
 - Enable extended ACLs
 - Specify path prefixes

Record Service

- Currently in public beta
- Layer that enforces fine-grained Sentry permissions
 - Without RS, file based access is all-or-nothing
 - View / Column permissions for Spark and MapReduce
- Opens up possibilities for dynamic data masking and tokenization
- Can swap out entire storage layer.



Post Configuration

- HDFS setup with a better umask and service level authorization
- YARN setup with restrictive admin ACLs
- Hive, Impala, and HDFS setup with Sentry integration
- **DEMO:** No more default authorization holes!

Authorization - Summary

- HDFS file permissions (POSIX 'rwx rwx rwx' style)
- Yarn job queue permissions
- Sentry (Hive / Impala / Solr / Kafka)
- Cloudera Manager RBAC
- Cloudera Navigator RBAC
- Hue groups
- Hadoop KMS ACLs
- HBase ACLs
- etc.



strataconf.com
#StrataData

PRESENTED BY



Questions

Encryption of Data in Transit

Syed Rafice

Senior Systems Engineer
Cloudera

Agenda

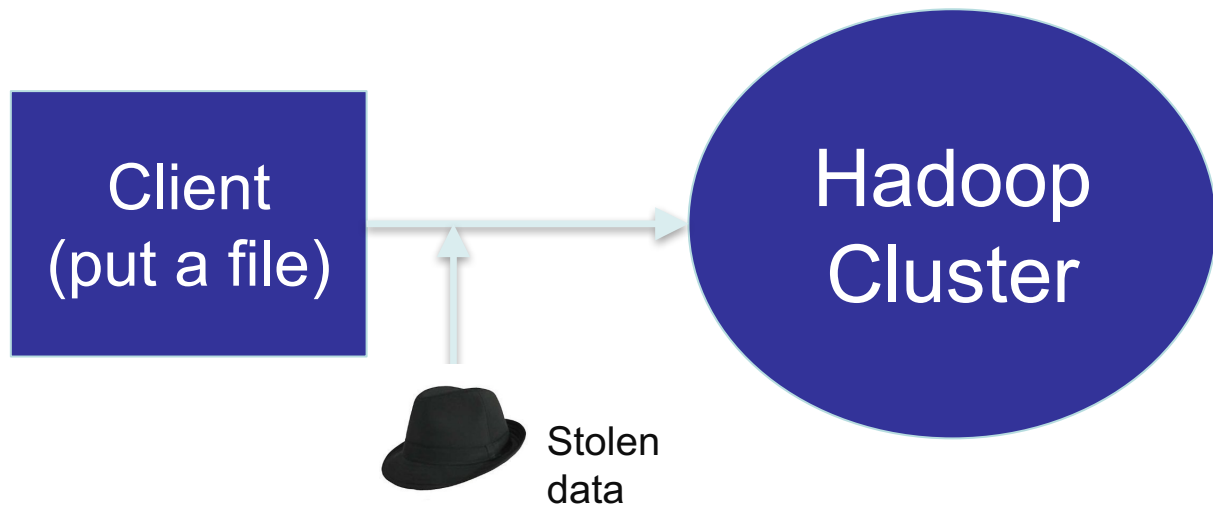
- Why encryption of data on the wire is important
- Technologies used in Hadoop
 - SASL “Privacy”
 - TLS
- For each:
 - Demo without
 - Discussion
 - Enabling in Cloudera Manager
 - Demo with it enabled

Why Encrypt Data in Transit?

- Networking configuration (firewalls) can mitigate some risk
- Attackers may already be inside your network
- Data and credentials (usernames and passwords) have to go into and out of the cluster
- Regulations around transmitting sensitive information
- Let's see this for real using wireshark

Demo

- Transfer data into a cluster
- Simple file transfer: “hadoop fs -put”
- Attacker sees file contents go over the wire



Two Encryption Technologies




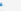
- SASL “confidentiality” or “privacy” mode
 - Protects core hadoop
- TLS – Transport Layer Security
 - Used for “everything else”

SASL

- Simple Authentication and Security Layer
- Not a protocol, but a framework for passing authentication steps between a client and server
- Pluggable with different authentication types
 - GSS-API for Kerberos (Generic Security Services)
- Can provide transport security
 - “auth-int” – integrity protection: signed message digests
 - “auth-conf” – confidentiality: encryption

SASL Encryption - Setup

- First, enable Kerberos
- HDFS:
 - Hadoop RPC Protection
 - Datanode Data Transfer Protection
 - Enable Data Transfer Encryption
 - Data Transfer Encryption Algorithm
 - Data Transfer Cipher Suite Key Strength

Hadoop RPC Protection hadoop.rpc.protection	HDFS-1 (Service-Wide)  <input type="radio"/> authentication <input type="radio"/> integrity <input checked="" type="radio"/> privacy
DataNode Data Transfer Protection dfs.data.transfer.protection	HDFS-1 (Service-Wide)  <input type="radio"/> Authentication <input type="radio"/> Integrity <input checked="" type="radio"/> Privacy
Enable Data Transfer Encryption dfs.encrypt.data.transfer	HDFS-1 (Service-Wide) <input checked="" type="checkbox"/> 
Data Transfer Encryption Algorithm dfs.encrypt.data.transfer.algorithm	HDFS-1 (Service-Wide)  <input type="radio"/> 3des <input type="radio"/> rc4 <input checked="" type="radio"/> AES/CTR/NoPadding
Data Transfer Cipher Suite Key Strength dfs.encrypt.data.transfer.cipher.key.bitlength	HDFS-1 (Service-Wide) <input type="radio"/> 128 <input type="radio"/> 192 <input checked="" type="radio"/> 256

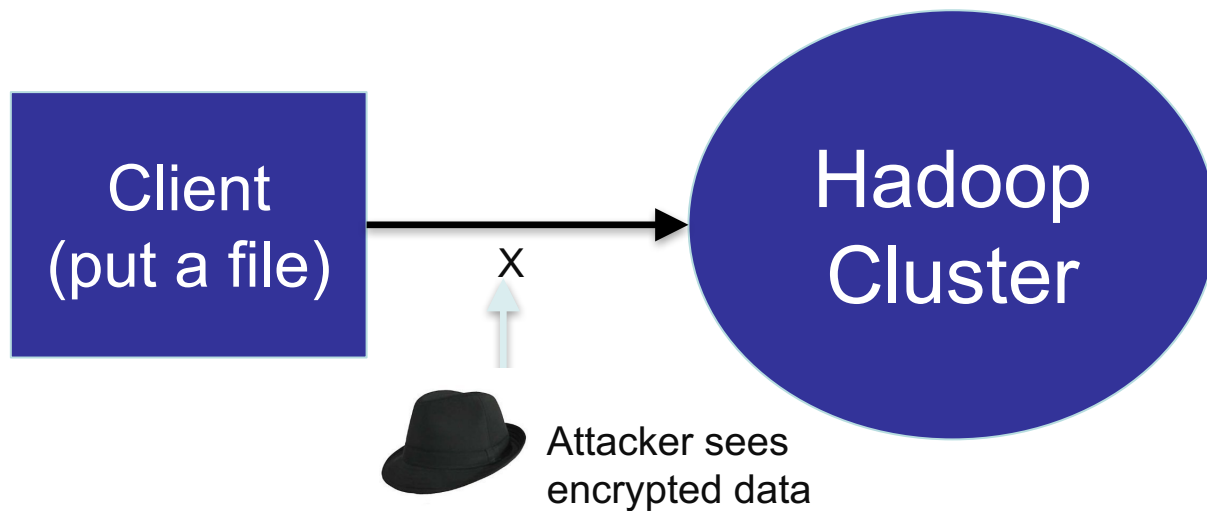
SASL Encryption - Setup

- Hbase
 - HBase Thrift Authentication
 - Hbase Transport Security

HBase Thrift Authentication hbase.thrift.security.qop	HBASE-1 (Service-Wide) C <ul style="list-style-type: none"><input type="radio"/> none<input type="radio"/> auth<input type="radio"/> auth-int<input checked="" type="radio"/> auth-conf
HBase Transport Security hbase.rpc.protection	HBASE-1 (Service-Wide) C <ul style="list-style-type: none"><input type="radio"/> authentication<input type="radio"/> integrity<input checked="" type="radio"/> privacy

Demo 2

- Put a file into HDFS again
- But this time with SASL encryption turned on



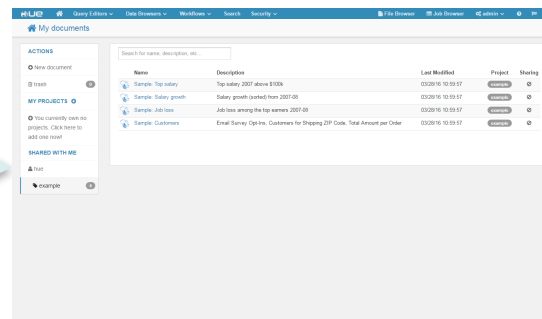
TLS

- Transport Layer Security
 - The successor to SSL – Secure Sockets Layer
 - The term SSL was deprecated 15 years ago, but we still use it
 - TLS is what's behind https:// web pages
- Let's what happens with no TLS (an http connection to Hue)

Web Browser (http)



Stolen admin
credentials



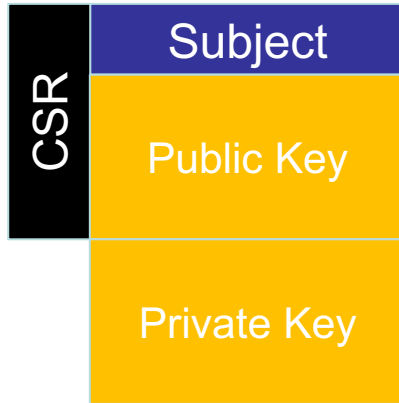
TLS - Certificates

- TLS relies on certificates for authentication
- You'll need one certificate per machine
- Certificates:
 - Cryptographically prove that you are who you say you are
 - Are issued by a "Certificate Authority" (CA)
 - Have a "subject", an "issuer" and a "validity period"
 - Many other attributes, like "Extended Key Usage"
 - Let's look at an https site

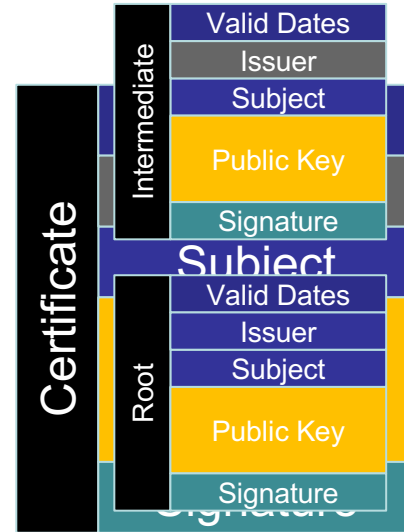
TLS – Certificate Authorities

- “Homemade” CA using openssl
 - Suitable for test/dev clusters only
- Internal Certificate Authority
 - A CA that is trusted widely inside your organization, but not outside
 - Commonly created with Active Directory Certificate Services
 - Web browsers need to trust it as well
- External Certificate Authority
 - A widely known CA like VeriSign, GeoTrust, Symantec, etc
 - Costs \$\$\$ per certificate

You



Certificate Authority



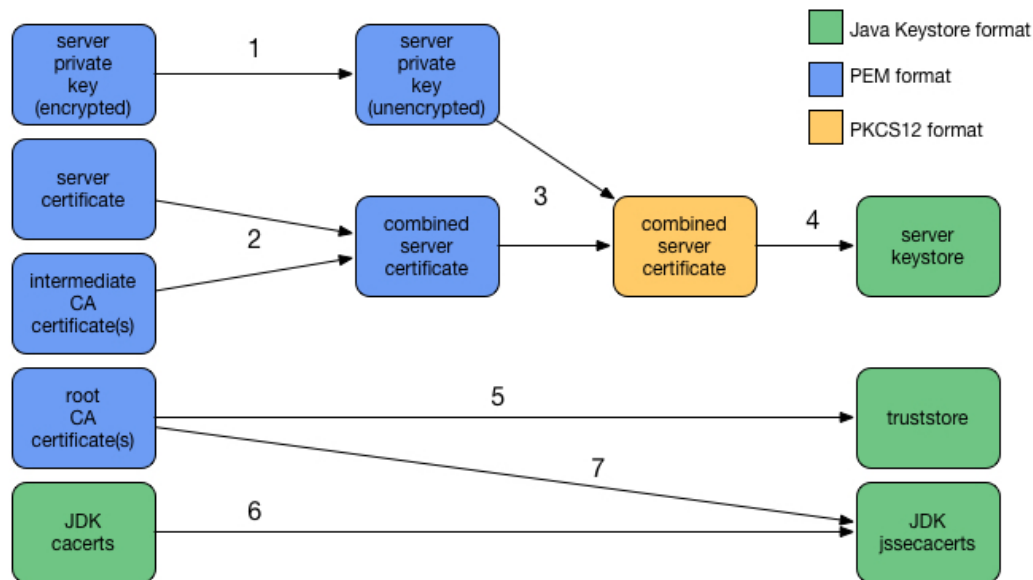
TLS – Certificate File Formats

- Two different formats for storing certificates and keys
- PEM
 - “Privacy Enhanced Mail” (yes, really)
 - Used by openssl; programs written in python and C++
- JKS
 - Java KeyStore
 - Used by programs written in Java
- The Hadoop ecosystem uses both
- Therefore you must translate private keys and certificates into both formats

TLS – Key Stores and Trust Stores



- Keystore
 - Used by the server side of a TLS client-server connection
 - JKS: Contains private keys and the hosts's certificate; Password protected
 - PEM: typically one certificate file and one password-protected private key file
- Truststore
 - Used by the client side of a TLS client-server connection
 - Contains certificates that the client trusts: the Certificate Authorities
 - JKS: Password protected, but only for an integrity check
 - PEM: Same concept, but no password
 - There is a system-wide certificate store for both PEM and JKS formats.

TLS – Key Stores and Trust Stores



- 1 - openssl rsa -in `hostname` -f.key.temp -out `hostname` -f.key
- 2 - cat server.pem int-CA.pem > `hostname` -f.pem
- 3 - openssl pkcs12 -export -in `hostname` -f.pem -inkey `hostname` -f.key -out `hostname` -f.pfx
- 4 - keytool -importkeystore -srcstoretype PKCS12 -srckeystore `hostname` -f.pfx -destkeystore `hostname` -f.jks
- 5 - keytool -importcert -file root-CA.pem -alias root-CA -keystore truststore.jks
- 6 - cp \$JAVA_HOME/jre/lib/security/cacerts \$JAVA_HOME/jre/lib/security/jssecacerts
- 7 - keytool -importcert -file root-CA.pem -alias root-CA -keystore \$JAVA_HOME/jre/lib/security/jssecacerts

TLS – Securing Cloudera Manager

- CM Web UI -  <https://>
- CM Agent -> CM Server communication – 3 “Levels” of TLS use
 - Level 1: Encrypted but no certificate verification. Akin to clicking on  ~~<https://>~~
 - Level 2: Agent verifies the server’s certificate
 - Level 3: Agent and Server verify each other’s certificate. This is called TLS mutual authentication: each side is confident that it’s talking to the other
 - Note: TLS level 3 requires that certificates are suitable for both “TLS Web Server Authentication” and “TLS Web Client Authentication”
 - Very Sensitive Information goes over this channel
 - Like Kerberos Keytabs. Therefore, set up TLS in CM first before Kerberos

Cloudera Manager TLS

Use TLS Encryption for Admin Console	<input checked="" type="checkbox"/>	← CM Web UI
Requires Server Restart		
Use TLS Encryption for Agents	<input checked="" type="checkbox"/>	← TLS Level 1
Requires Server Restart		
Use TLS Authentication of Agents to Server	<input checked="" type="checkbox"/>	← TLS Level 3
Requires Server Restart		
Cloudera Manager TLS/SSL Server JKS Keystore File Location	<input type="text" value="/opt/cloudera/security/jks/keystore.jks"/>	
Requires Server Restart		
Cloudera Manager TLS/SSL Server JKS Keystore File Password	<input type="password" value="....."/>	
Requires Server Restart		
Cloudera Manager TLS/SSL Certificate Trust Store File	<input type="text" value="/opt/cloudera/security/jks/truststore.jks"/>	
Requires Server Restart		
Cloudera Manager TLS/SSL Certificate Trust Store Password	<input type="password" value="....."/>	
Requires Server Restart		

The CM Agent Settings

- Agent `/etc/cloudera-scm-agent/config.ini`

`use_tls=1` ← TLS Level 1

`verify_cert_file=` full path to CA certificate.pem file ← TLS Level 2

`client_key_file=` full path to private key.pem file

`client_keypw_file=` full path to file containing password for key

`client_cert_file=` full path to certificate.pem file

} TLS Level 3

TLS for CM-Managed Services

- CM requires that all files (jks and pem) are in the same location on each machine
- For each service (HDFS, Hue, Hbase, Hive, Impala, ...)
 - Search the configuration for “TLS”
 - Check the “enable” boxes
 - Provide keystore, truststore, and passwords

Hive Example

Enable TLS/SSL for HiveServer2 hive.server2.enable.SSL, hive.server2.use.SSL	HIVE-1 (Service-Wide) <input checked="" type="checkbox"/> ↕
HiveServer2 TLS/SSL Server	HIVE-1 (Service-Wide) ↕
JKS Keystore File Location hive.server2.keystore.path	<input type="text" value="/opt/cloudera/security/jks/keystore.jks"/> ⓘ
HiveServer2 TLS/SSL Server	HIVE-1 (Service-Wide)
JKS Keystore File Password hive.server2.keystore.password	<input type="password" value="....."/> ⓘ
HiveServer2 TLS/SSL Certificate	HIVE-1 (Service-Wide) ↕
Trust Store File	<input type="text" value="/opt/cloudera/security/jks/truststore.jks"/>
HiveServer2 TLS/SSL Certificate	HIVE-1 (Service-Wide)
Trust Store Password	<input type="password" value="....."/> ⓘ

TLS - Troubleshooting

- To examine certificates
 - `openssl x509 -in <cert>.pem -noout -text`
 - `keytool -list -v -keystore <keystore>.jks`
- To attempt a TLS connection as a client
 - `openssl s_client -connect <host>:<port>`
 - This tells you all sorts of interesting TLS things

Demo - TLS

- Let's try to attack an https connection to Hue
- Note that this is only one example, TLS protects many, many things in hadoop

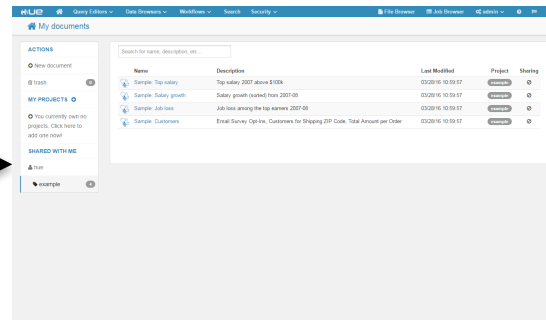
Web Browser (https)



X



Attacker sees
encrypted data



Conclusions

- You need to encrypt information on the wire
- Technologies used are SASL encryption and TLS
- TLS requires certificate setup



strataconf.com
#StrataData

PRESENTED BY



Questions?

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

HDFS Encryption at Rest

Mubashir Kazia

Senior Solutions Architect
Cloudera

Agenda

- Why Encrypt Data
- Demo
- HDFS Encryption
- Demo
- Questions

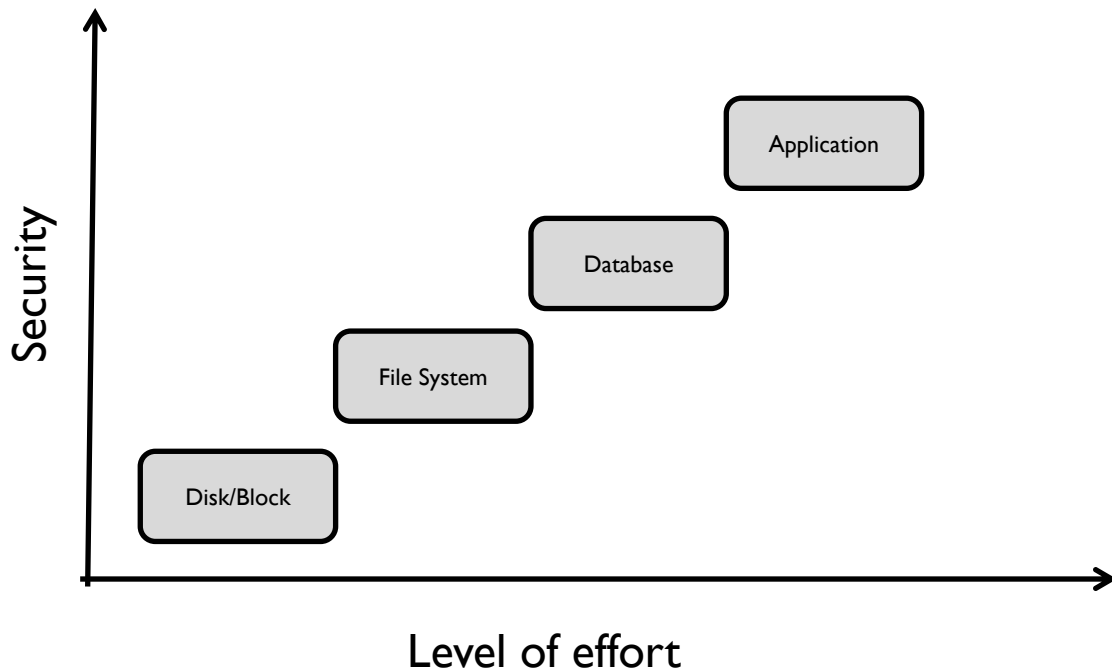
Why store encrypted data?

- Customers often are mandated to protect data at rest
 - PCI
 - HIPAA
 - National Security
 - Company confidential
- Encryption of data at rest helps mitigate certain security threats
 - Rogue administrators (insider threat)
 - Compromised accounts (masquerade attacks)
 - Lost/stolen hard drives

Demo

- How to access HDFS data from Linux storage bypassing HDFS authorization

Options for encrypting data



Architectural Concepts

- Encryption Zones
- Keys
- Key Management Server

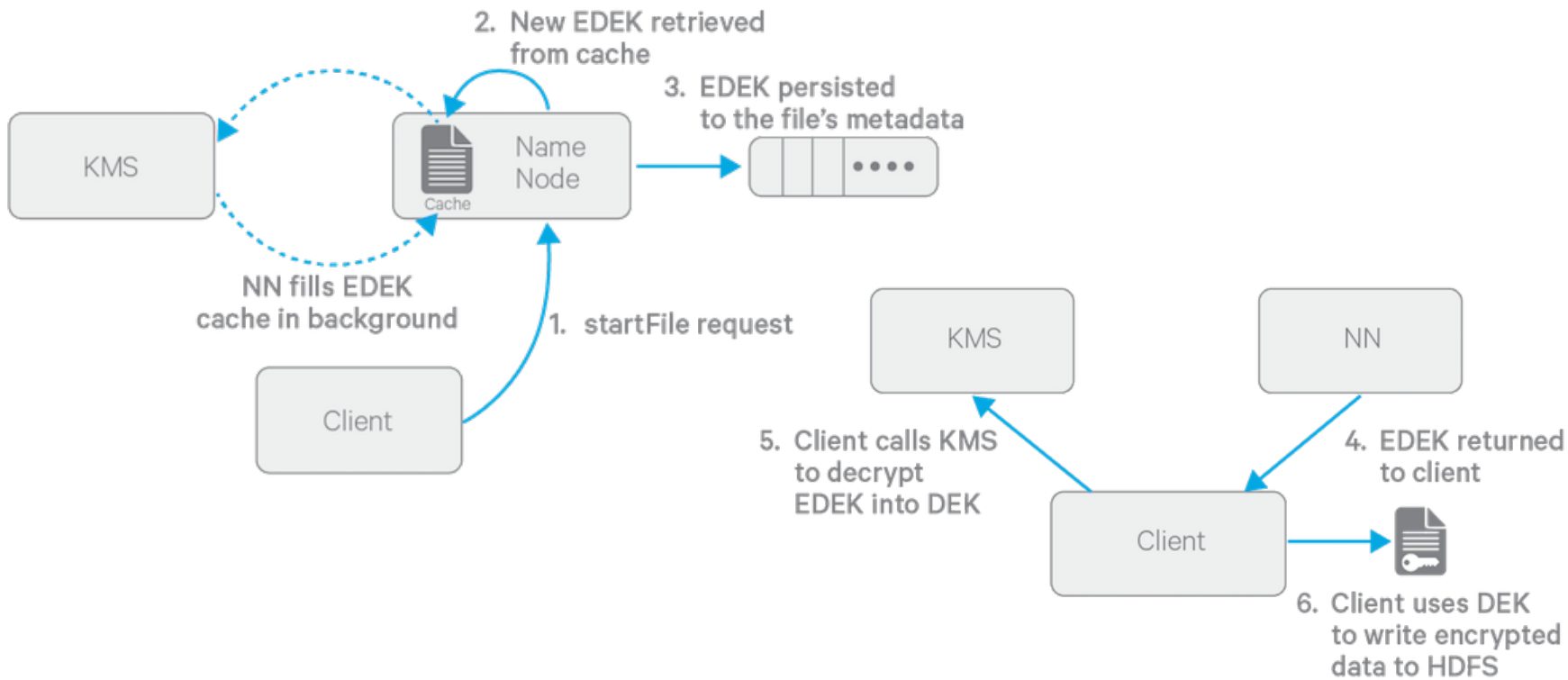
Encryption Zones

- An HDFS directory in which the contents (including subdirs) are encrypted on write and decrypted on read.
- An EZ begins life as an empty directory
- Rename/Move in/out of an EZ are prohibited
- Encryption is transparent to application with no code changes

EZ Keys, Data Encryption Keys, and Encrypted Data Encryption Keys



Key Handling



Key Management Server (KMS)

- KMS sits between client and key server
 - E.g. Cloudera Navigator Key Trustee
- Provides a unified API and scalability
- REST API
- Does not actually store keys (backend does that), but does cache them
- ACLs on per-key basis

HDFS Encryption Configuration

- `hadoop key create <keyname> -size <keySize>`
- `hdfs dfs -mkdir <path>`
- `hdfs crypto -createZone -keyName <keyname> -path <path>`

KMS Per-User ACL Configuration

- White lists (check for inclusion) and black lists (check for exclusion)
- `etc/hadoop/kms-acls.xml`
 - `hadoop.kms.acl.CREATE`
 - `hadoop.kms.blacklist.CREATE`
 - ... `DELETE`, `ROLLOVER`, `GET`, `GET_KEYS`, `GET_METADATA`, `GENERATE_EEK`, `DECRYPT_EEK`
 - `hadoop.kms.acl.<keyname>.<operation>`
 - `MANAGEMENT`, `GENERATE_EEK`, `DECRYPT_EEK`, `READ`, `ALL`

Best practices

- Enable authentication (Kerberos)
- Enable TLS/SSL
- Use KMS acls to setup KMS roles, blacklist HDFS admins and grant per key access
- Do not use the KMS with default JCEKS backing store
- Use hardware that offers AES-NI instruction set
 - Install openssl-devel so Hadoop can use Openssl crypto codec
- Make sure you have enough entropy on all the nodes
 - Run rngd or haveged

Best practices

- Do not run KMS on master or worker nodes
- Run multiple instances of KMS for high availability and load balancing
- Harden KMS instance and use internal firewall so only KMS and ssh etc. ports are reachable from known subnets
- Make secure backups of KMS

HDFS Encryption - Summary

- Good performance (4-10% hit) with AES-NI
- No mods to existing applications
- Prevents attacks at the filesystem and below
- Data is encrypted all the way to the client
- Key management is independent of HDFS
- Can prevent HDFS admin from accessing secure data

Demo

- Accessing HDFS encrypted data from Linux storage

User	Group	Role
hdfs_admin	cdh_admin	HDFS Admin
kms_admin	cdh_admin	KMS Admin
alice	cdh_user	User with DECRYPT_EEK access to key1
bob	cdh_user	User with DECRYPT_EEK access to key2



strataconf.com
#StrataData

PRESENTED BY



Questions?

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera


Hadoop Data Governance

Mark Donsky


Director, Product Management
Cloudera

Data Governance


Frequently Asked Questions



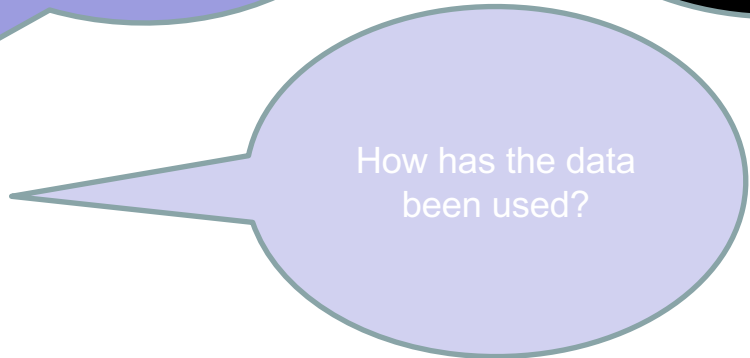
What data do I have?



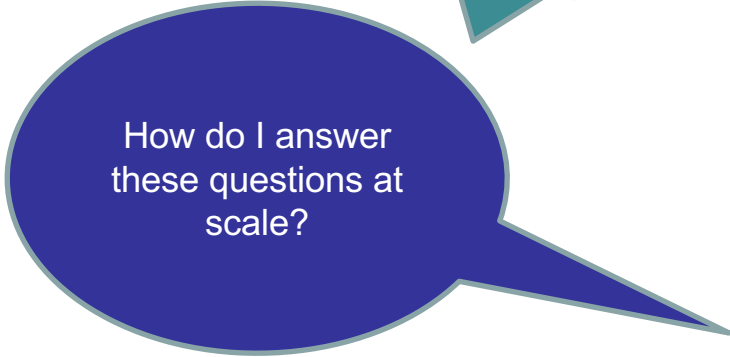
How did the data get here?



Who used the data?



How has the data been used?



How do I answer these questions at scale?

What makes big data governance different?

Governing big data
requires governing
petabytes of diverse types
of data

No one application will
solve every big data
governance problem

Applications are shifting to
the cloud, and data
governance must still be
applied consistently

Self-service data
discovery is mandatory for
big data

Compliance + Productivity = Adoption



Compliance/Governance

- Am I prepared for an audit?
- Who's accessing sensitive data?
- What are they doing with the data?
- Is sensitive data governed and protected?

End User Productivity

- How can I find explore data sets on my own?
- Can I trust what I find?
- How do I use what I find?
- How do I find and use related data sets?

What makes governance so difficult?

Hadoop governance challenges

- Variety, Volume, Velocity
- Multiple compute types: Spark, Hive, Pig, MR, MR2, Sqoop, etc.
- Multiple third-party tools

Cloud governance challenges

- Multiple storage types: HDFS, S3, ADLS, etc.
- Transient clusters
- Long-running clusters
- Shared Hive Metastores

Yet the business still needs one set of trusted governance artifacts

Governance: the Foundation of Data Management

Compliance

Track, understand and protect access to data

Am I prepared for an audit?

Who's accessing sensitive data?

What are they doing with the data?

Is sensitive data governed and protected?

Stewardship

Manage and organize data assets at scale

How can I efficiently manage data lifecycle, from ingest to purge?

How can I efficiently organize and classify all my data?

How can I efficiently make data available to my end users?

End User Productivity

Effortlessly find and trust data sets

How can I find explore data sets on my own?

Can I trust what I find?

How do I use what I find?

How do I find and use related data sets?

Administration

Boost user productivity and cluster performance

Is my data optimized to support current access patterns?

How can I optimize for future workloads?

How can I migrate workloads to Hadoop risk-free?

Hadoop Governance Foundation

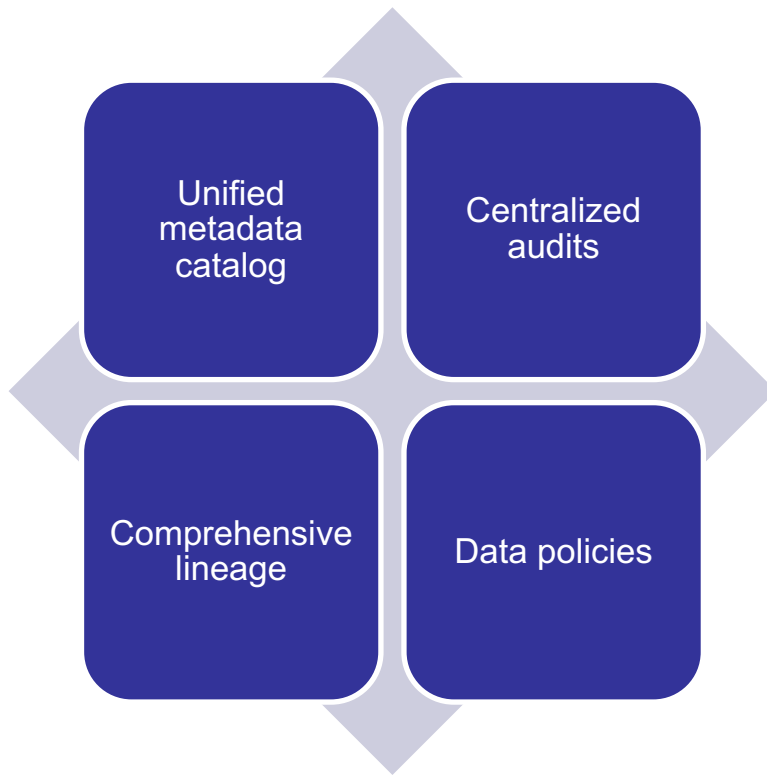
Centralized audits

Unified metadata catalog

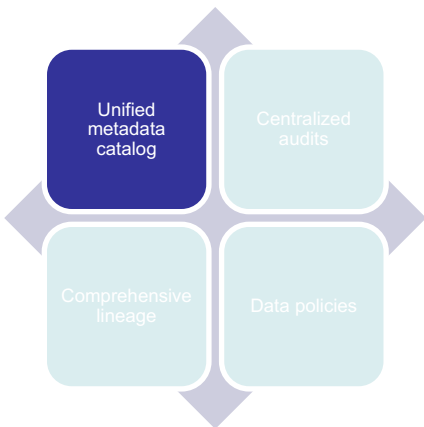
Comprehensive lineage

Data policies

Hadoop Governance Requirements



Unified Metadata Catalog



Technical Metadata

All files in directory /sales

All files with permissions 777

Anything older than 7 years

Any not accessed in the past 6 months

Managed Metadata

Sales data from last quarter for the Northeast region

Protected health information

Business glossary definitions

Data sets associated with clinical trial X

Custom Metadata

Tables that I want to share with my colleagues

Data sets that I want to retrieve later

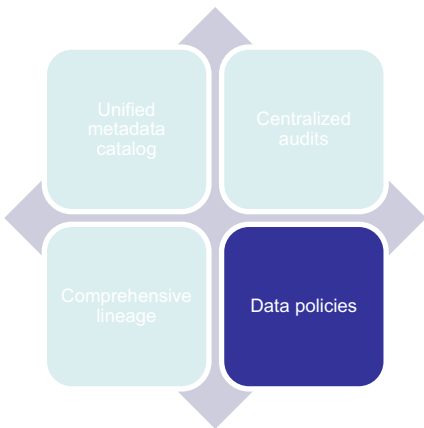
Data sets that are organized by my personal classification scheme (e.g., "quality = high")

Challenges

- Technical metadata in Hadoop is component-specific
- Curated/custom attributes: Hive meta store has comments, and HDFS has extended attributes, but:
 - Not searchable
 - No validation
- Aggregated analytics are not possible
 - How many files are older than two years?

Data Policies

- **Goal:** Manage and automate the information lifecycle from ingest to purge/cradle to grave, based on the unified metadata catalog
- Once you find data sets, you'll likely need to do something with them
 - Tag every new file that lands in /sales as sales data
 - Send an alert whenever a sensitive data set has permissions 777
 - Purge all files that are older than seven years

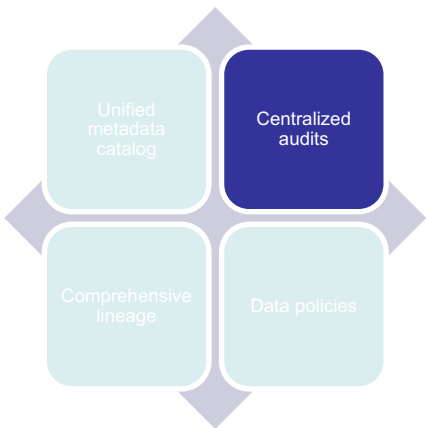


Challenges

- Oozie workflows can be difficult to configure
- Event-triggered oozie workflows are limited to very few technical metadata attributes, such as directory path
- Data stewards prefer to define, view, and manage data policies in a metadata-centric fashion

Centralized Audits

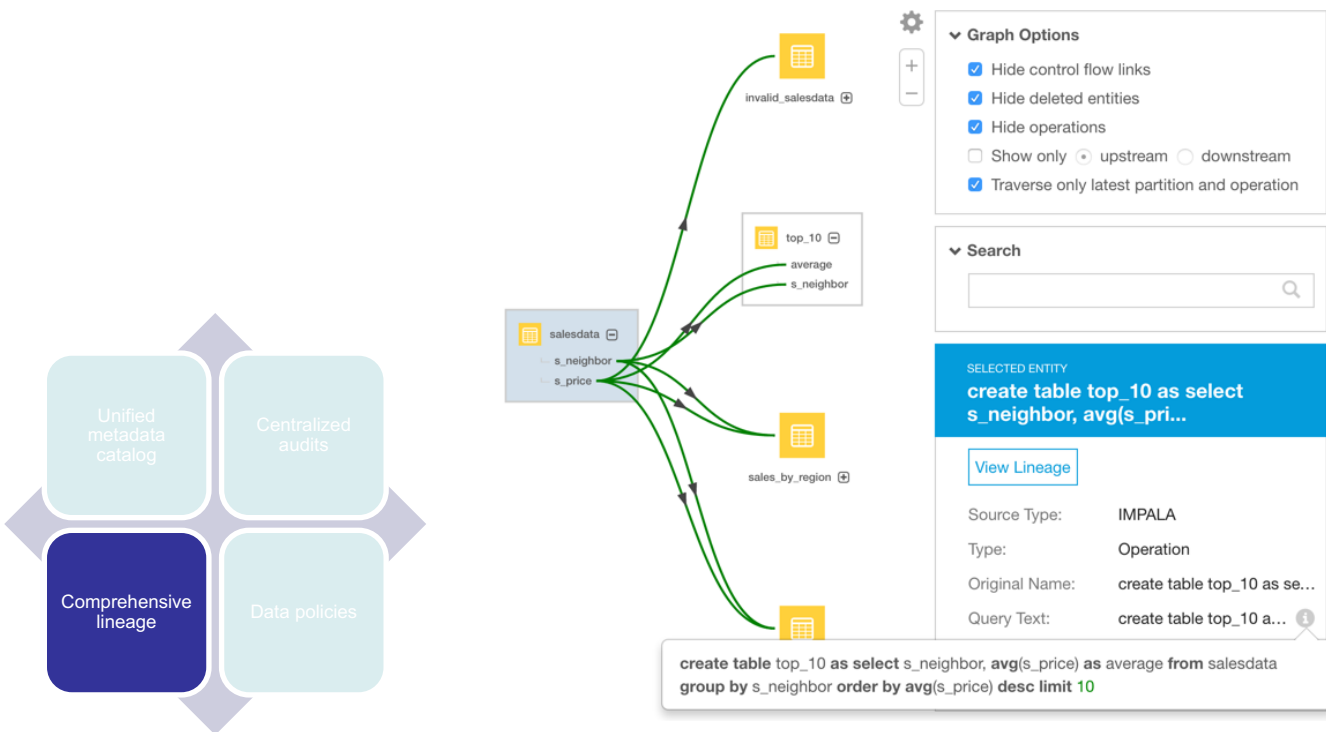
- **Goal:** Collect all audit activity in a single location
 - Redact sensitive data from the audit logs to simplify compliance with regulation
 - Perform holistic searches to identify data breaches quickly
 - Publish securely to enterprise tools



Challenges

- Each component has its own audit log, but:
- Sensitive data may exist in the audit log
 - `Select * from transactions where cc_no = "1234 5678 9012 3456"`
- It's difficult to do holistic searches
 - What did user *a* do yesterday?
 - Who accessed file *f*?
- Integration with enterprise SIEM and audit can be complex

Comprehensive Lineage



Challenges

- Most uses of lineage require column-level lineage
- Hadoop does not capture lineage in an easily-consumable format
- Lineage must be collected automatically and cover all compute engines
- Third-party tools and custom-built applications need to augment lineage

strataconf.com
#StrataData

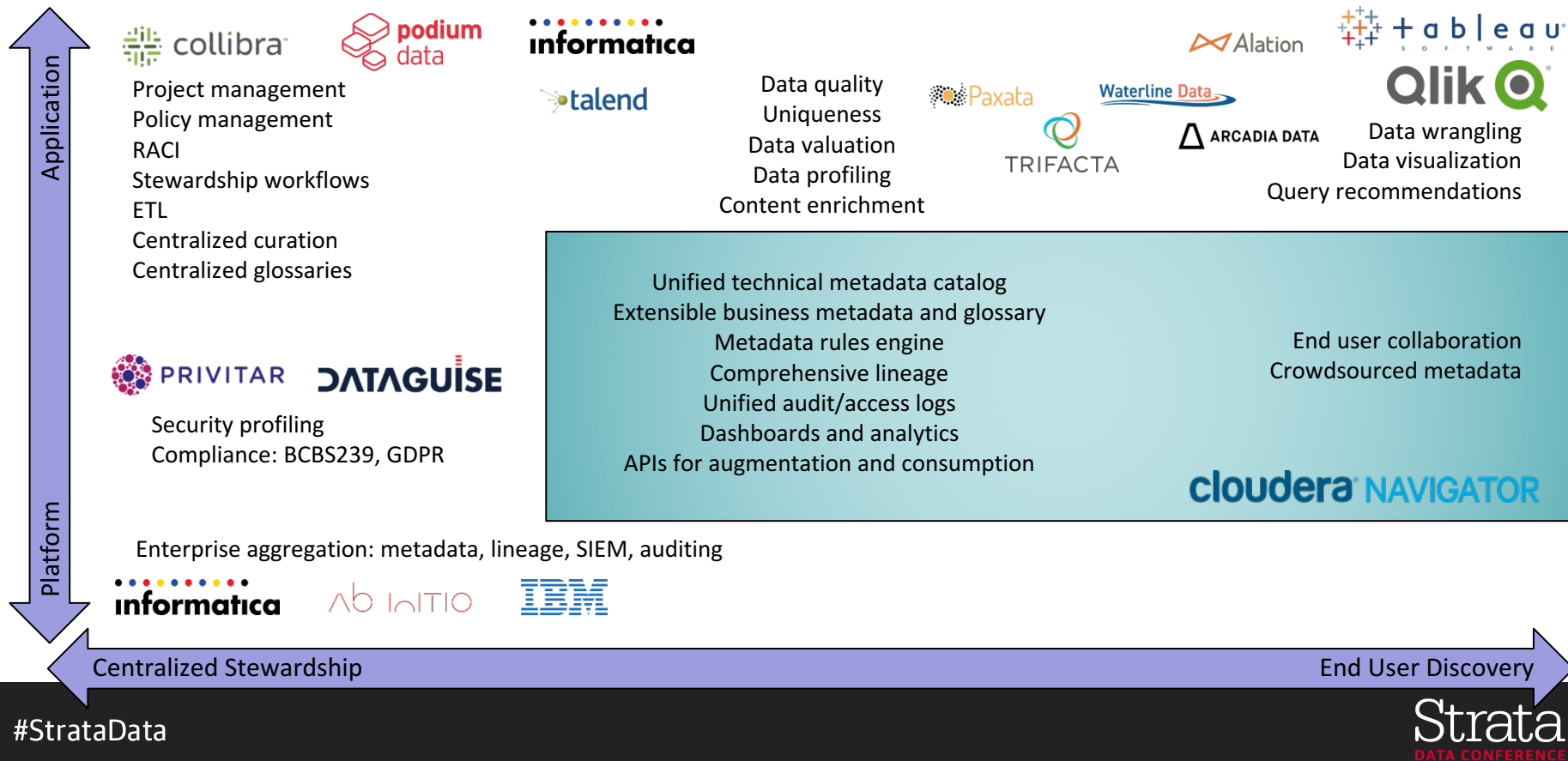
PRESENTED BY

O'REILLY

cloudera

Use Cases

Data Stewardship and Governance Activities



Use Cases: Compliance

Compliance

Track, understand and protect access to data

Am I prepared for an audit?

Who's accessing sensitive data?

What are they doing with the data?

Is sensitive data governed and protected?

ENTERPRISE METADATA REPOSITORY

 informatica

 Data Advantage Group

 *adaptive*

ENTERPRISE AUDITING & SECURITY

 splunk>

 IMPERVA

 IBM

 RSA
SECURITY™

HADOOP DATA GOVERNANCE & MANAGEMENT

Centralized audits

Unified metadata catalog

Comprehensive lineage

Data policies

Common use cases:

- Security breach detection
- Data access tracking for PCI compliance
- Audit defense

Use Cases: Administration

Administration

Boost user productivity
and cluster performance

Is my data optimized
to support current
access patterns?

How can I optimize
for future
workloads?

How can I migrate
workloads to
Hadoop risk-free?

Visibility

- Distribution of data objects
- Workloads by engine

Patterns

- Data churn over time
- Table clusters
- Frequent users

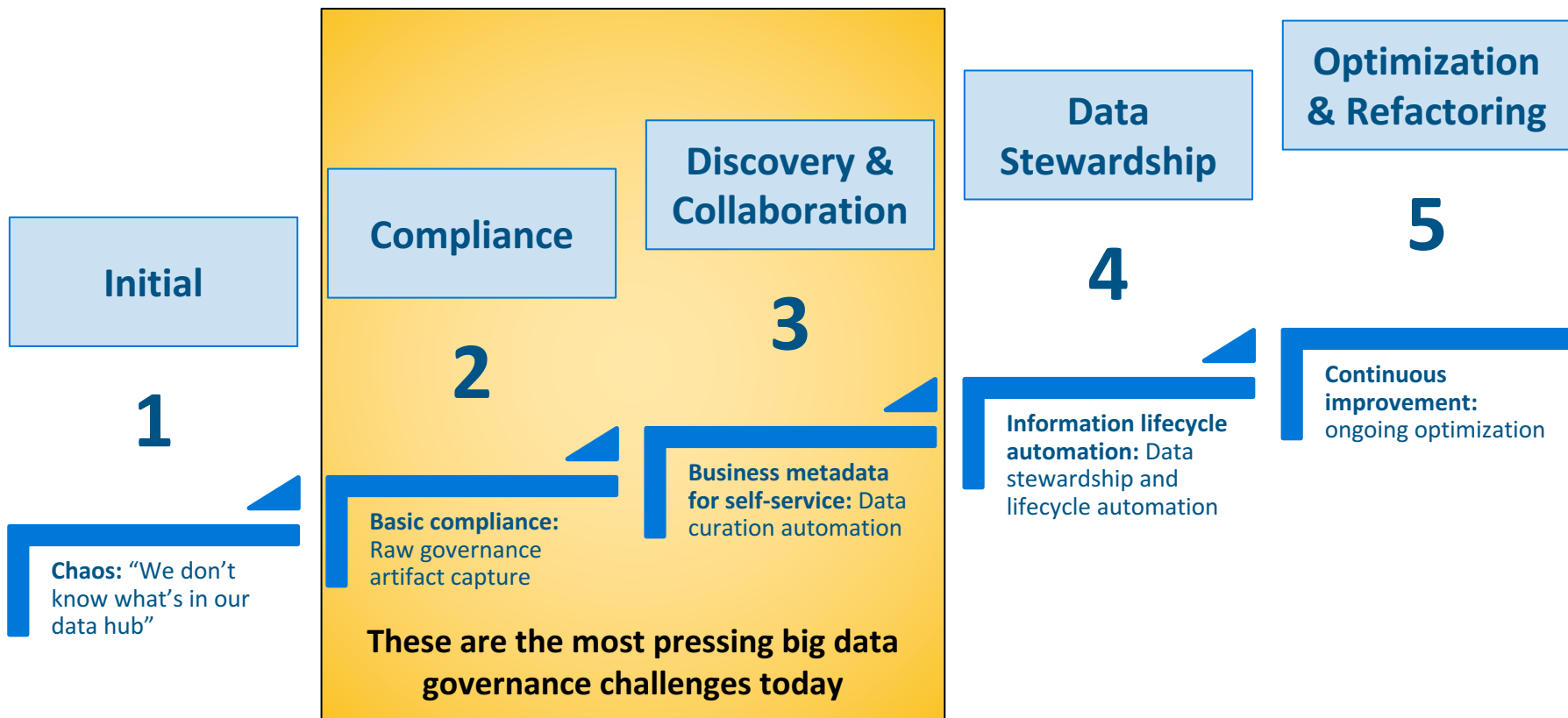
Optimization

- Sub-optimal query patterns
- “Rogue” users
- Capacity planning

Unexpected Behavior

- Hive tables suddenly missing
`rm -rf /user/hive/warehouse`

The current state of big data governance



Strata

DATA CONFERENCE

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

Demo



strataconf.com
#StrataData

PRESENTED BY



Questions

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

Final Thoughts

Compliance

- We have shown how an EDH environment can be secured end-to-end
- Is this enough to be compliant?
 - PCI DSS, HIPAA, GDPR
 - Internal compliance – PII data handling
- All of the security features discussed (and others not covered because of time) are enough to cover technical requirements for compliance
- However, compliance also requires additional **people** and **process** requirements
- Cloudera has worked with customers to achieve PCI DSS compliance as well as others – **you can do it too!**

Public Cloud Security

- Many Hadoop deployments occur in the public cloud
- Security considerations presented today all still apply
- Complementary to native cloud security controls

- **Cloudera blog post - How-to: Deploy a secure enterprise data hub on AWS**
- <http://blog.cloudera.com/blog/2016/05/how-to-deploy-a-secure-enterprise-data-hub-on-aws/>

Looking Ahead

- The Hadoop ecosystem is vast, and it can be a daunting task to secure everything
- Understand that **no system is completely secure**
- However, the proper security controls coupled with regular reviews can **mitigate** your exposure to threats and vulnerabilities
- Pay attention to new components in the stack, as these components often **do not** have the same security features in place
 - Kafka only recently added wire encryption and Kerberos authentication
 - Spark only recently added wire encryption
 - Many enterprises were using both of these in production before those features were available!

strataconf.com
#StrataData

PRESENTED BY

O'REILLY

cloudera

Final Questions?

Thank you!