



Hadoop Security

Building a fence around your Hadoop cluster

Lars Francke

June 12, 2017

Berlin Buzzwords 2017

Introduction

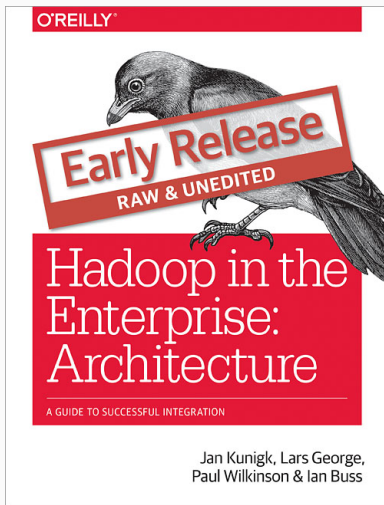
- Partner & Co-Founder at OpenCore
- Before that: EMEA Hadoop Consultant
- Hadoop since 2008/2009
- Apache Committer: Hive, ORC
- Contact:
 - lars.francke@opencore.com
 - [@lars_francke](https://twitter.com/lars_francke)

Overview





Source: oreilly.com



Source: oreilly.com

- Kerberos: The Definitive Guide¹
- Kerberos (German)²
- Active Directory, 5th Edition³
- Hadoop and Kerberos: The Madness beyond the Gate⁴
- HBase: The Definitive Guide, 2nd Edition⁵
- Bulletproof SSL and TLS⁶

¹<http://shop.oreilly.com/product/9780596004033.do>

²<http://www.kerberos-buch.de/>

³<http://shop.oreilly.com/product/0636920023913.do>

⁴https://www.gitbook.com/book/steveloughran/kerberos_and_hadoop/details

⁵<http://shop.oreilly.com/product/0636920033943.do>

⁶<https://www.feistyduck.com/books/bulletproof-ssl-and-tls/>

How is a typical project structured and where does *Hadoop security* come into play?

- Idea/Initiation

- Idea/Initiation
- Planning/Design

- Idea/Initiation
- Planning/Design
- Execution/Implementation

- Idea/Initiation
- Planning/Design
- Execution/Implementation
- “Production”

You can not begin thinking about Security *too early!*

Let's dig into the details of each step in a project

Project Planning

- What am I going to use the cluster for?

- What am I going to use the cluster for?
- Which tools am I going to need to achieve the use-cases?

- What am I going to use the cluster for?
- Which tools am I going to need to achieve the use-cases?
- What kind of data am I going to ingest, store or process?

- What am I going to use the cluster for?
- Which tools am I going to need to achieve the use-cases?
- What kind of data am I going to ingest, store or process?
- What kind of corporate guidelines exist that must be followed?

Do not *assume* companies/people know what a product is capable of just because they seem confident.

What is part of *Hadoop Security*?

- Authentication of users

- Authentication of users
- Authorization of users

- Authentication of users
- Authorization of users
- Auditing

- Authentication of users
- Authorization of users
- Auditing
- Data Protection: Encryption on-the-wire

- Authentication of users
- Authorization of users
- Auditing
- Data Protection: Encryption on-the-wire
- Data Protection: Encryption at-rest

Project Execution

So, does Hadoop Security mean I need to “enable” Kerberos and I’m done?

So, does Hadoop Security mean I need to “enable” Kerberos and I’m done?

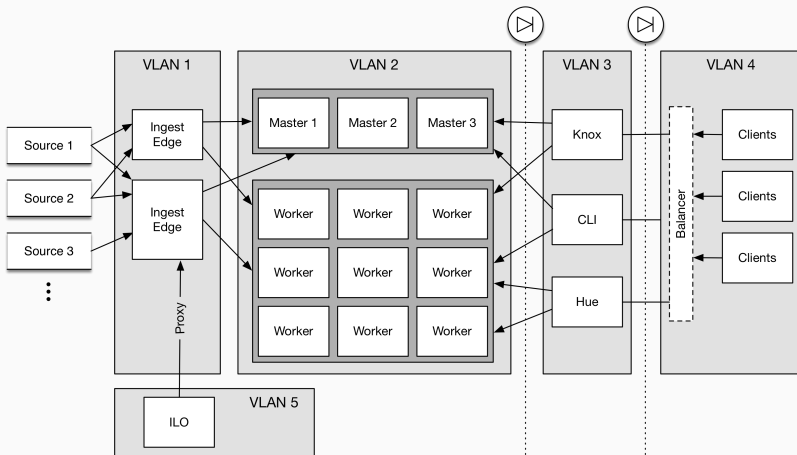
No! Far from it!

Before you're able to secure anything...

Before you're able to secure anything...
...something must be running that *can* be secured.

Overview - contd.

6. Clients		<i>e.g. Qlik, SPSS, BI</i>		Business Administrator
5a. Applications	<i>e.g. Spark, MR, Giraph</i>	5b. Access Engines	<i>e.g. Big SQL, Hive</i>	
4. Core Hadoop		<i>e.g. Kerberos, TLS, ACLs, RPC, GRANTs</i>		Hadoop Administrator
3. Host		<i>e.g. firewall, ACLs, SELinux, SSSD, PAM (authz)</i>		
2. Network		<i>e.g. routing, VLANs, firewalls</i>		System Administrator
1. Perimeter		<i>e.g. datacenter, cage, servers</i>		



- SELinux
- NTP/Chrony
- Firewalls
- Antivirus
- Proxies

- You usually install a cluster manager first

- You usually install a cluster manager first
- Then you install Agents

- You usually install a cluster manager first
- Then you install Agents
- Followed by a distribution with lots of components...

- Core Hadoop requires Kerberos for strong authentication
- Multiple choices on how to implement that:
 - Cluster-local standalone KDC
 - Cluster-local standalone KDC with one-way cross-realm trust to a central KDC
 - Direct integration into a central KDC

Authentication: Cloudera Manager/Ambari Wizards



Source: unsplash.com by Sirotorn Sumpunkulpak

What you need:

- Names/IPs for your KDC
- Supported Encryption Types (see my blog post⁷ for more details)
- Firewall must allow all cluster machines to access the KDC
- Potentially a bunch of information to configure `krb5.conf` properly

⁷<http://www.opencore.com/blog/2017/3/kerberos-encryption-types>

Ideally you want the automatic option, but...

- You need an account in your AD/KDC that is allowed to create other accounts!
- (AD only) You need to talk via LDAPS and make sure that you have all the truststores set up properly
- (AD only) You cannot have multiple SPNs per User
- For SPNEGO you need a `HTTP/<host>` principal, sometimes doesn't match corporate guidelines

There's a manual option...

Kerberos is only part of the authentication story.

Knox, Cloudera Manager, Ambari, Ranger, Sentry, Hive, Impala etc. all have their own authentication layers which also support LDAP(S) or other mechanisms.

Your users need to *exist* on the workers so that YARN can start jobs using `setuid()`.

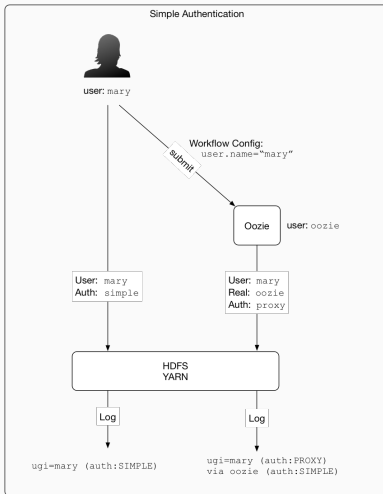
- This is usually done using a tool like SSSD (free) or Centrify (commercial)

- This is usually done using a tool like SSSD (free) or Centrify (commercial)
- When using Centrify, configure it to *not* create the default HTTP principal for each server

- This is usually done using a tool like SSSD (free) or Centrify (commercial)
- When using Centrify, configure it to *not* create the default HTTP principal for each server
- This does *not* mean that your users need to be able to SSH into the workers, quite the opposite

- You need to have the details needed to fetch user & group information from LDAP(S)
- Some solutions require a domain join

Impersonation/Proxy Users



- All tools have some form of authorization built-in
- Core Hadoop components have a first line defense: Service Level Authorizations
- Ranger & Sentry promise cross-cutting RBAC functionality

HDFS Transparent Encryption

All I need is a KMS, right?

- KMS should be “close” to your clients and NameNode
- KMS should be separately administered
- AFAIK only Hadoop 3 will allow the / (root) directory to be its own zone and have sub-zones
- Secure access to the backing store & host where KMS runs
- You can use a HSM
- What about Impersonation and things like Knox or HttpFS?

Now that my data in HDFS is secure I'm good, right?

- Metadata in databases
- Log & Audit files
- YARN localized stuff
- Temporary files
- (Spark) spill files & cached data

Only Cloudera has a (paid) solution for this: Cloudera NavEncrypt

- Static vs. dynamic masking
- Ranger & BigSQL support masking, Sentry does not
- RecordService might be a potential integration point
- 3rd party tools

- Web Interfaces
- REST/Thrift Interfaces
- MapReduce Shuffle
- Spark Shuffle (only in 2.1+)
- Server to Agent communication
- Ingest/Egress tools (Flume, Informatica, Kafka, ...)
- RPC
- HDFS Data traffic
- Traffic from your YARN apps

- 2-way TLS possible?
- Which cipher suites are supported?
- Some tools automatically disable HTTP others don't!
- Missing documentation all around

- No documentation (on auditing events and lots of other things)
- Ranger aggregates audit logs, async by default
- Cloudera has Navigator
- No integrity protection
- Need to protect against admins

- Not a single product has good documentation
- How do they access the cluster?
- Do they themselves authenticate and authorize users? How?
- Auditing/Logging?

Production

What now?

So you've finished the project and handed it over to the operations team...what now?

- Plan regular updates to your OS
- Use a Configuration Management tool (e.g. Ansible)
- Plan *regular* updates to your distribution!

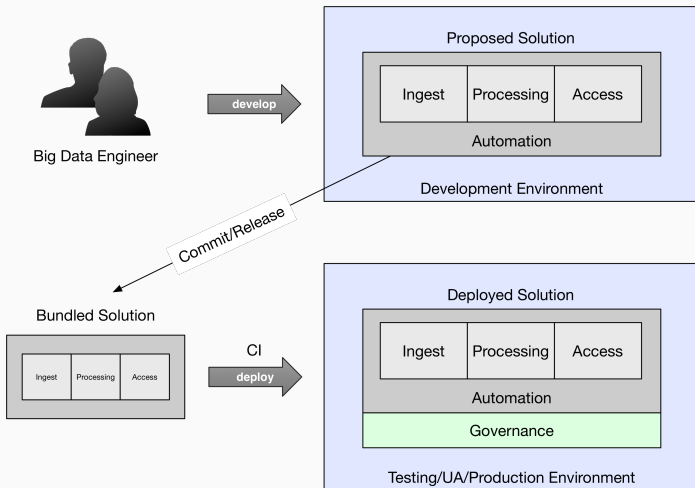
All our logging & auditing is irrelevant if you're not monitoring & alerting on those logs.

You need a place where you can test upgrades.

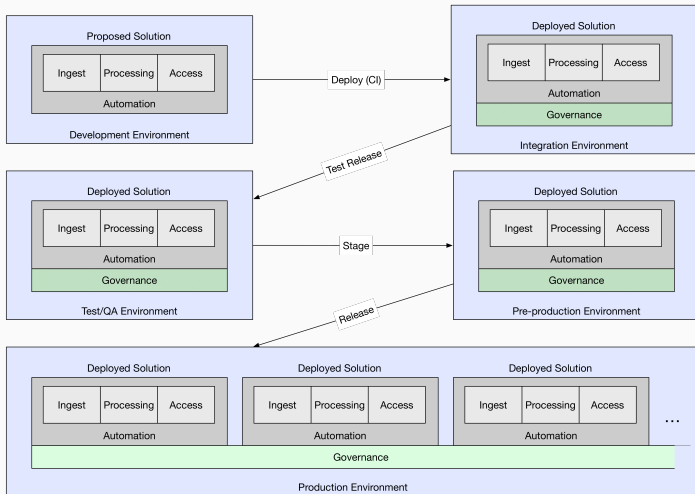
You might also need a place for backups.

Devs would like a cluster as well...

Environments



Environments



- What happens when a user leaves your organization?
- Caches
- Ranger Usersync
- Existing sessions
- Shared accounts

Misc

What about the cloud providers?

- Cloudera

- Cloudera
- Hortonworks

- Cloudera
- Hortonworks
- IBM

- Cloudera
- Hortonworks
- IBM
- Microsoft HDInsight

Thank you for listening!

Questions?

Contact me at:

- lars.francke@opencore.com
- [@lars_francke](#)

Visit us at opencore.com

And if this stuff interests you: We're looking to expand our team!