

# Hadoop & Kerberos



DevFest2016

# Sobre nosotros



# Sobre el taller

Preparar un servidor para instalar

Desplegar un cluster de un solo nodo [pseudo distribuido]

Kerberizar el cluster que hemos desplegado

Hablar sobre medidas de seguridad en Hadoop

# Entorno para Laboratorio

Maquina Virtual con Centos - Instalación con imagen mínima!

- 4GB RAM
- Networking en modo brigde
- 20GB HDD Thin provisioned



# Laboratorio Uno!



# HADOOP

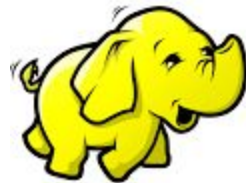
**Hadoop:** Creado por Doug Cutting, inspirado por Google File System y MapReduce.

Almacenamiento y procesamiento de grandes volúmenes de datos. Sobre “commodity hardware”.  
Tolerancia a errores incluida en el diseño.

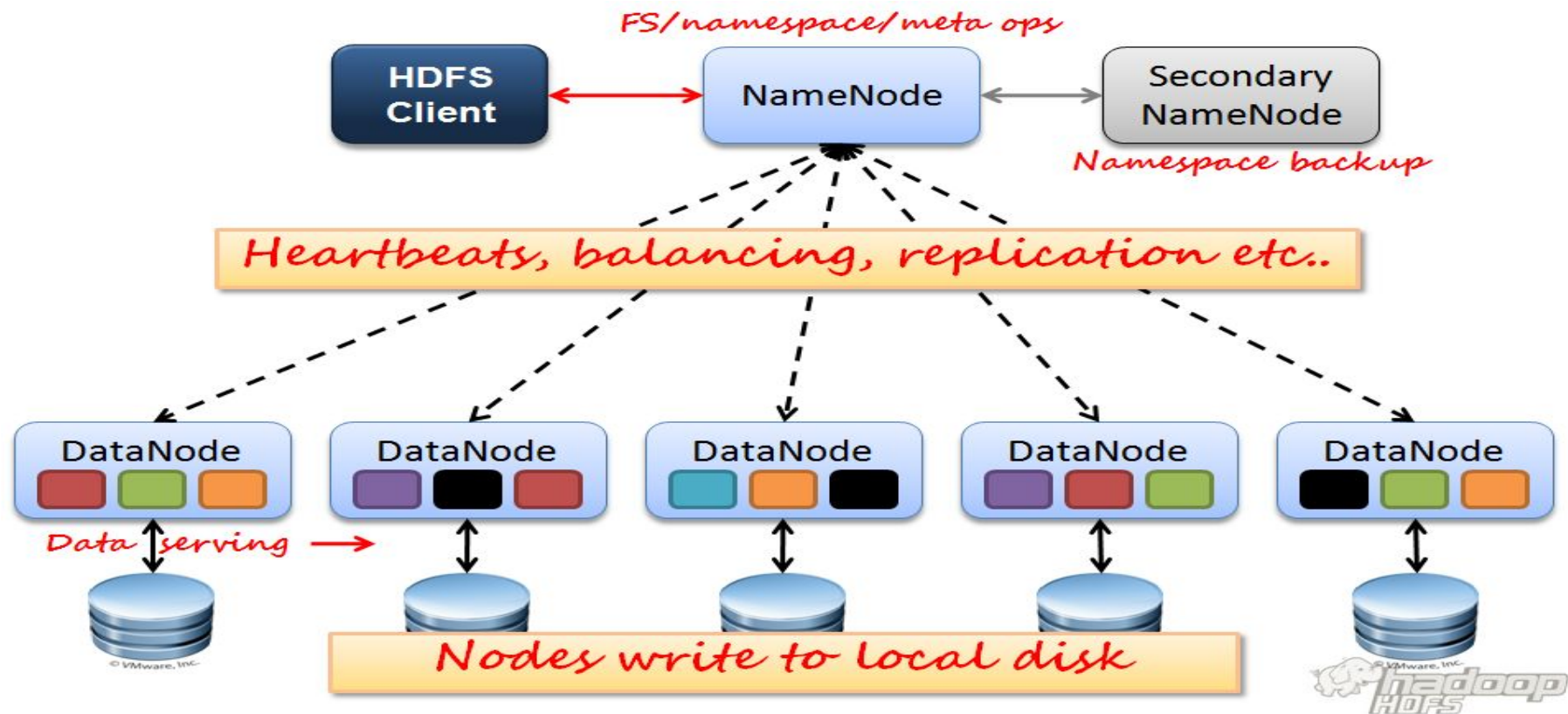
Compuesto por:

- HDFS - Hadoop Distributed Filesystem
- YARN - Yet Another Resource Negotiator
- MapReduce - Sistema de procesamiento en paralelo basado en YARN
- Librerías comunes

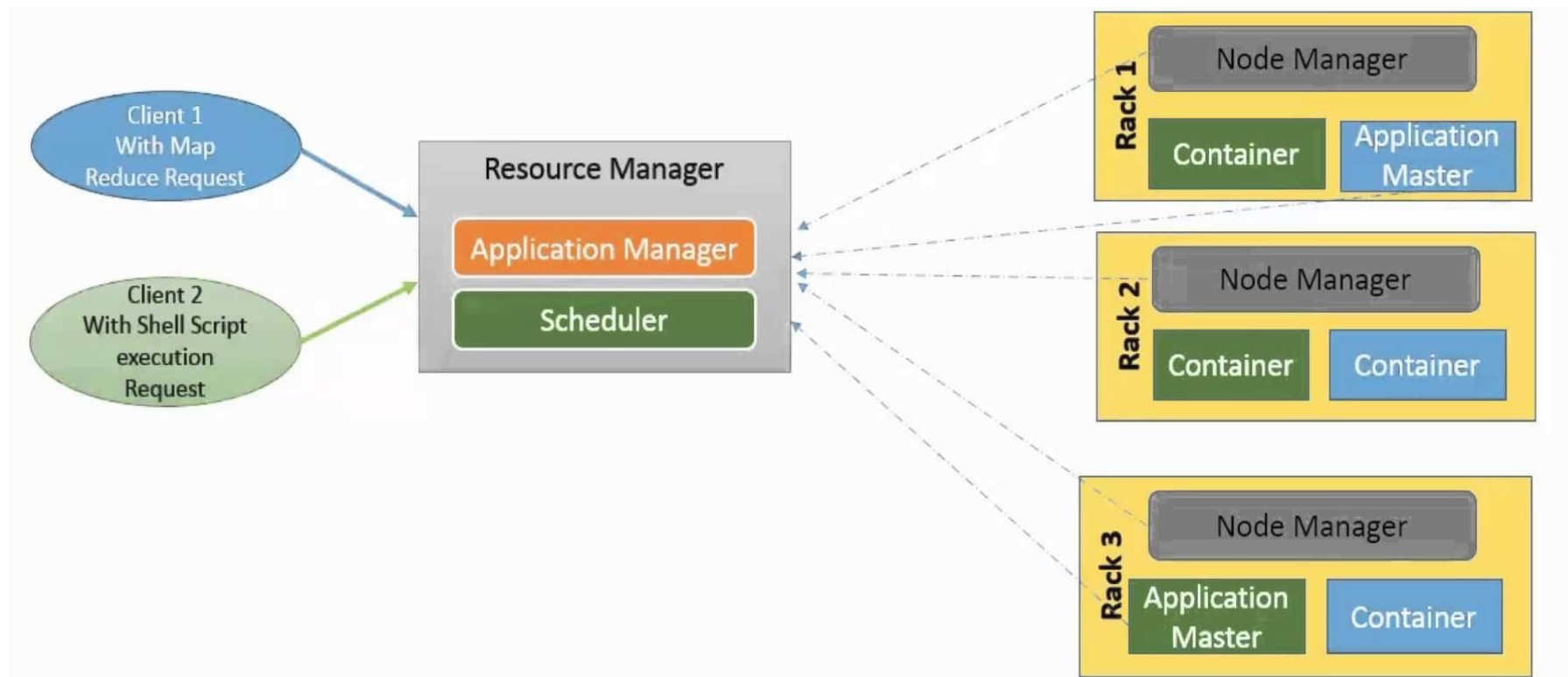
**Ecosistemas Hadoop:** conjunto de software relacionado con el proyecto Hadoop.



# HADOOP - HDFS

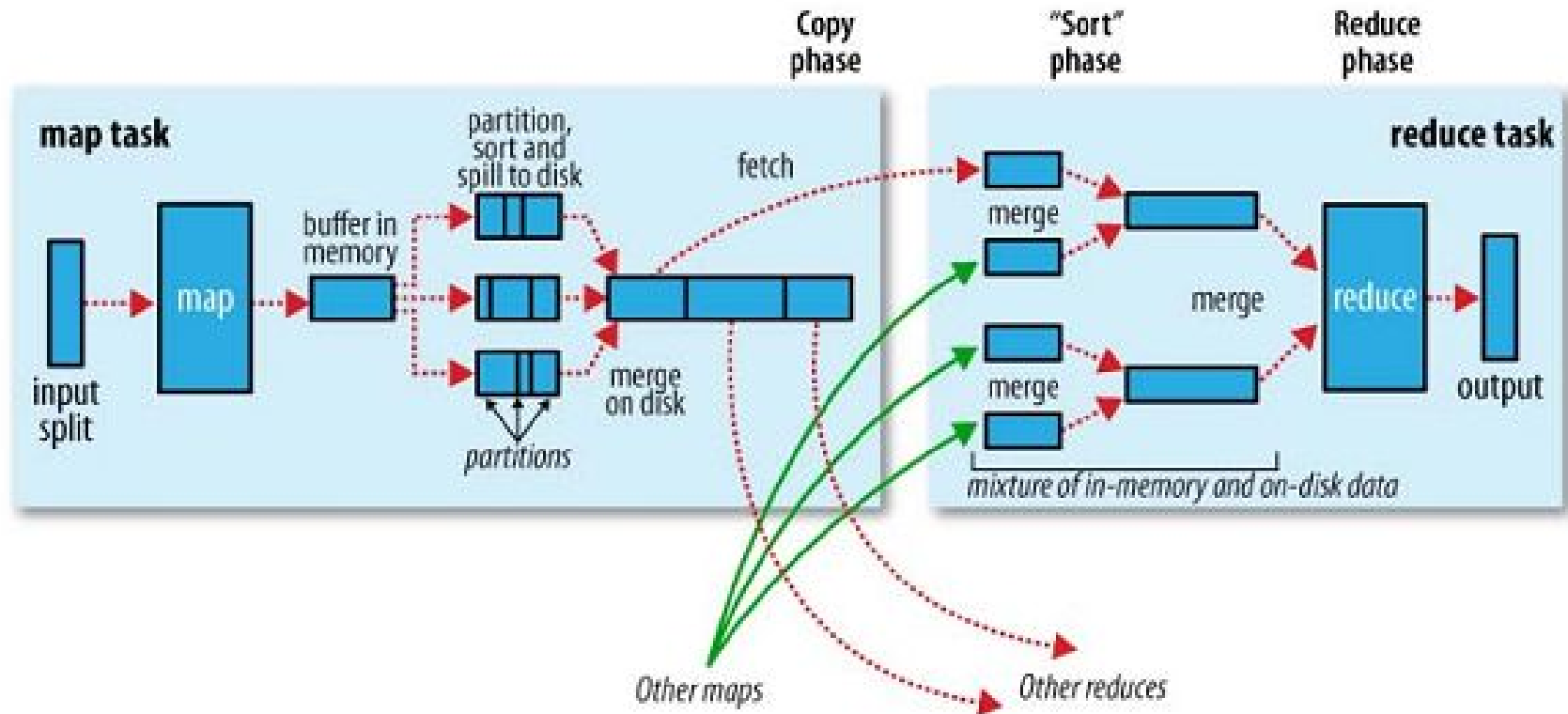


# HADOOP - YARN





# HADOOP - MapReduce



# Laboratorio Dos!



# Kerberos



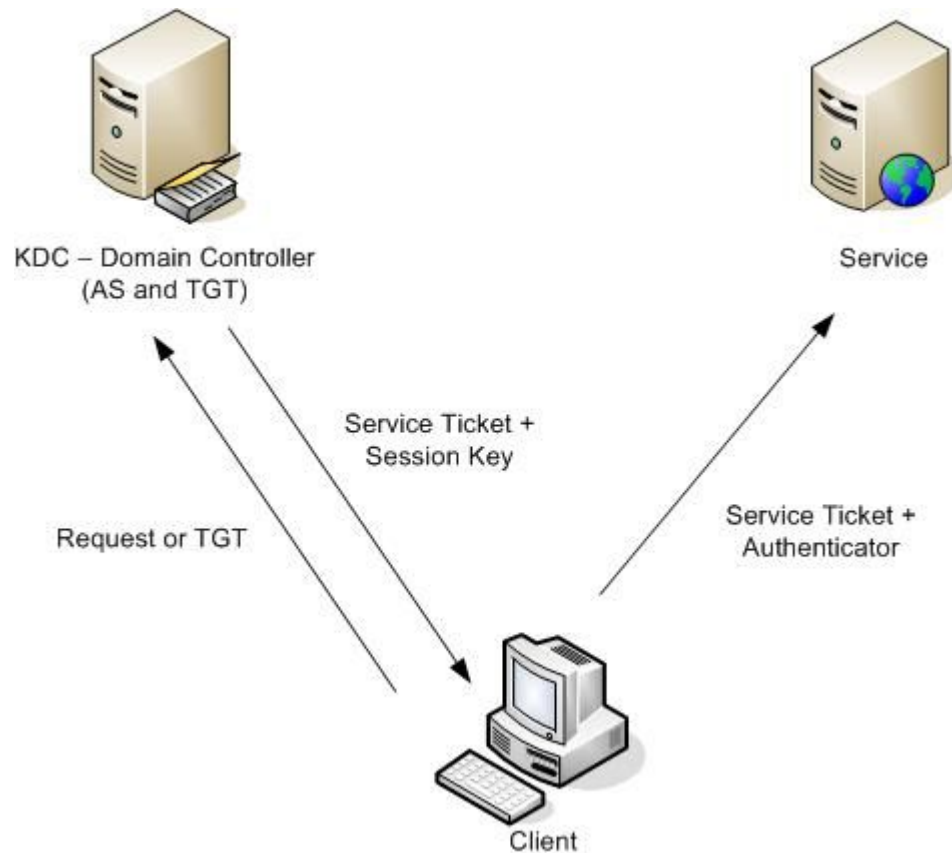
Autenticación en sistemas distribuidos

# Kerberos

kinit <identidad>

klist -ef

kdestroy



# Laboratorio Tres!



# Seguridad en Hadoop

## Perimeter

Guarding access to the cluster itself

### Technical Concepts:

Authentication  
Network isolation

## Data

Protecting data in the cluster from unauthorized visibility

### Technical Concepts:

Encryption  
Data masking

## Access

Defining what users and applications can do with data

### Technical Concepts:

Permissions  
Authorization

## Visibility

Reporting on where data came from and how it's being used

### Technical Concepts:

Auditing  
Lineage

# Perímetro - Autenticación

Proveedor de Autenticación:

- MIT Kerberos o compatible

Proveedor de usuarios:

- Sistema operativo / Integración mediante NSS [LDAP, NIS, etc ...]
- LDAP

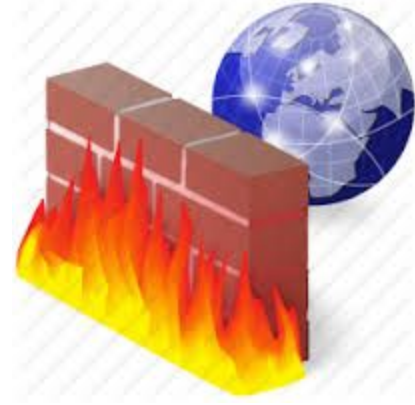
Gestores de Identidad:

- FreeIPA
- Active Directory

# Perímetro - Seguridad en Red

## Defensa perimetral

- filtrado a nivel de red (access list / firewalls)
- filtrado a nivel de host (iptables)
- uso de gateway o maquina de acceso





# Datos - Cifrado

## Cifrado de datos en tránsito

- Habilitar soporte TLS para HTTP
- Habilitar cifrado
  - RPC
  - Bloques



## Cifrado de datos en reposo

- Hadoop Key Management Server

# Datos - Protección de datos

## Cloudera - Recordservice

- Nuevo servicio y protocolo - Requiere cambios
- Enmascaramiento y filtrado - columnas y filas

## Hortonworks - Ranger

- Enmascaramiento y filtrado - columnas y filas
- Solo en Hive

# Acceso - Permisos

## Permisos

- Permisos Unix
- ACLs



# Acceso - Autorización

Autorización de uso

- ACL a nivel de servicio

Cloudera

- Sentry



Hortonworks

- Ranger

Apache  
Ranger

# Visibilidad - Auditoría

Cloudera

- Cloudera Navigator

Hortonworks

- Ranger

Apache  
Ranger

# Acceso - Linaje

## Cloudera

- Cloudera Navigator

## Hortonworks

- Atlas

Apache **Atlas**