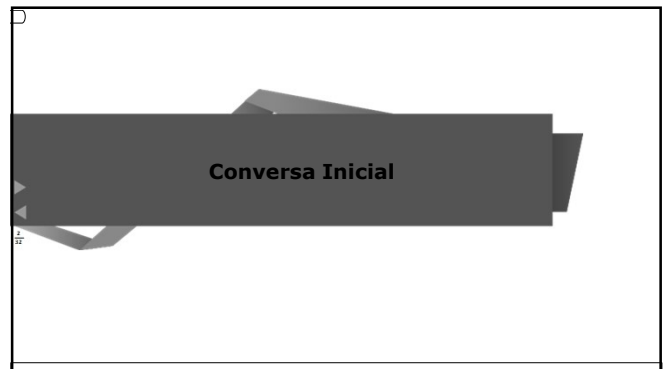
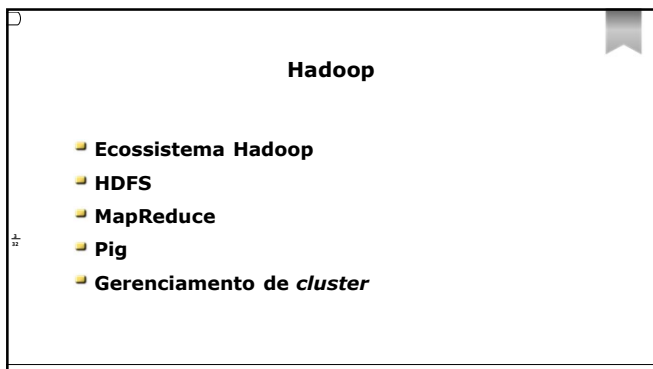


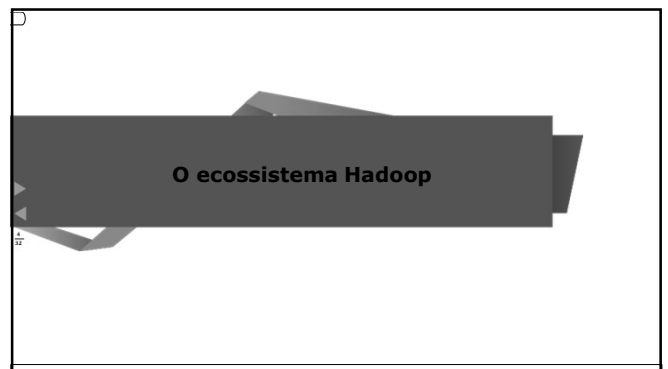
1



2



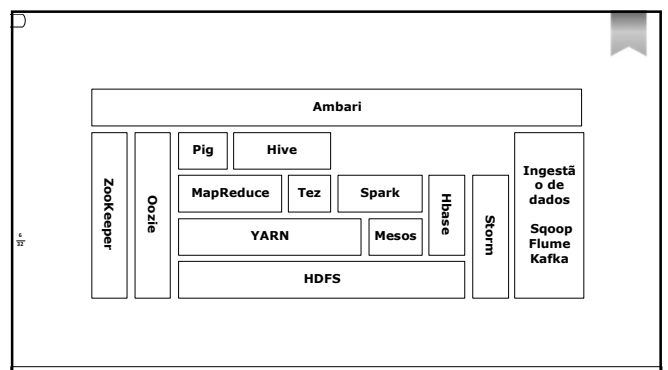
3



4



5



6

## Hadoop Distributed File System

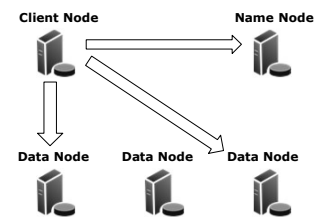
## HDFS

- Sistema distribuído para *cluster*
- Tolerante a falhas (partição)
- Alta disponibilidade
- Escalabilidade horizontal
- Acesso paralelo aos dados

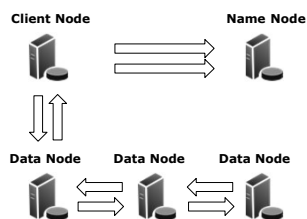
## Arquitetura HDFS

- Os dados são distribuídos em blocos de 128 MB
- Os blocos são replicados nos nós do *cluster*
- Name Nodes
- Data Nodes

## HDFS – leitura de arquivos



## HDFS – escrita de arquivos

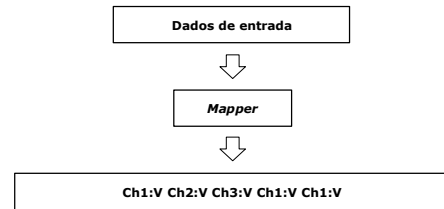


## MapReduce

## MapReduce

- Os *mappers* transformam os dados
- Os *reducers* agregam os resultados
- Tolerante a falhas

## Mapper



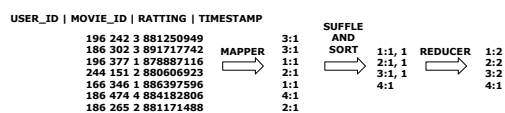
## Shuffle and sort

- Embaralhar e ordenar
- Processo automático
- Agrega dados por chave

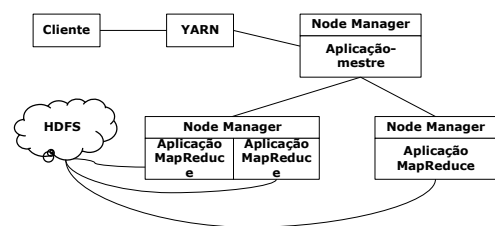
## Reducer

- Extraí as informações dos conjuntos de dados

## Exemplo



## Fluxo de processamento MapReduce



## Além do MapReduce

## Alternativas

- ▀ Mais simples?
- ▀ Mais eficientes?
- ▀ Mais versáteis?

## Pig

- ▀ Estende a capacidade do MapReduce
- ▀ Oferece uma linguagem semelhante à do SQL
- ▀ Executa *mappers* e *reducers* otimizados

## Pig Latin

- ▀ Facilidade de programação
- ▀ Oportunidade de otimização
- ▀ Extensibilidade

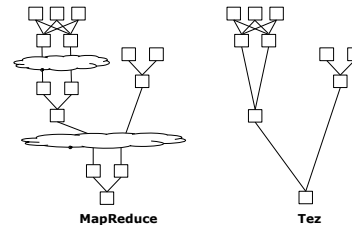
- ▀ Correspondência com bancos de dados tradicionais
  - Relação → tabela
  - Tupla → registro (linha)
  - Campo → coluna

- ▀ Permite operações JOIN
- ▀ Pode ser executado em conjunto com substitutos do MapReduce
  - Tez, Spark
- ▀ Interfaces
  - CLI Grunt
  - Ambari
  - Hue

## Tez

- Substitui o MapReduce
- Grafos acíclicos direcionados
- Pode otimizar a execução de outros componentes

## Processamento distribuído baseado em grafos acíclicos direcionados

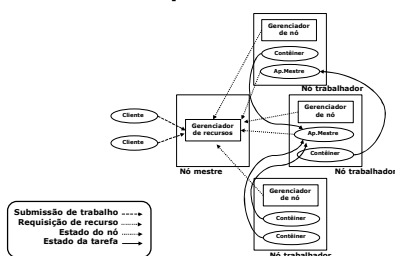


## Gerenciamento de cluster

## YARN

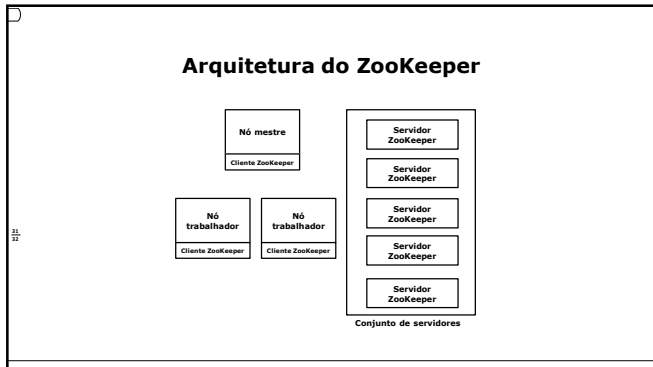
- "Yet Another Resource Navigator"
- Gerenciador de recursos global
- Distribui e gerencia o processamento pelo cluster
- Desmembramento do MapReduce original

## Arquitetura do YARN



## ZooKeeper

- Coordenador de aplicações distribuídas
- Alta disponibilidade
- Estrutura de árvores de diretório
- Permite manter informações dos nós e tarefas de uma aplicação

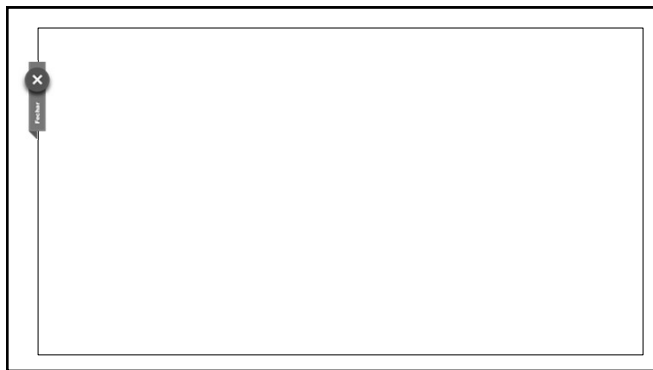


31

### Oozie

- Orquestração de trabalhos
- Fluxo de dados por vários componentes Hadoop
- Agendamento e execução de trabalhos
- Grafos acíclicos direcionados

32



33