

# Universidade da Beira Interior



Performance evaluation of a  
MongoDB and Hadoop platform for  
scientific data analysis

Alfredo Fernandes / José Chantre

# Investigação Científica

- Aumento de forma exponencial no volume de dados
- Aumento da variedade de dados
- Aumento da taxa a que os dados são analisados
- Dificuldade com as ferramentas tradicionais
- Necessidade de novas ferramentas para tratar enormes quantidades de dados semi-estruturados

# Ferramentas

**Hadoop** - Implementação OpenSource de modelo de programação  
MapReduce

**MongoDB** – Base de dados orientada a documentos – (NoSql) (Dados  
semi-estruturados)

# Apache Hadoop

- Plataforma Open-Source de [software](#) em [Java](#) para computação distribuída
- Orientado a clusters
- Adequado para o processamento de grandes quantidades de dados.

- Hadoop Distributed File System (HDFS)  
Sistema de ficheiros distribuído  
Permite acesso com elevado débito a dados .

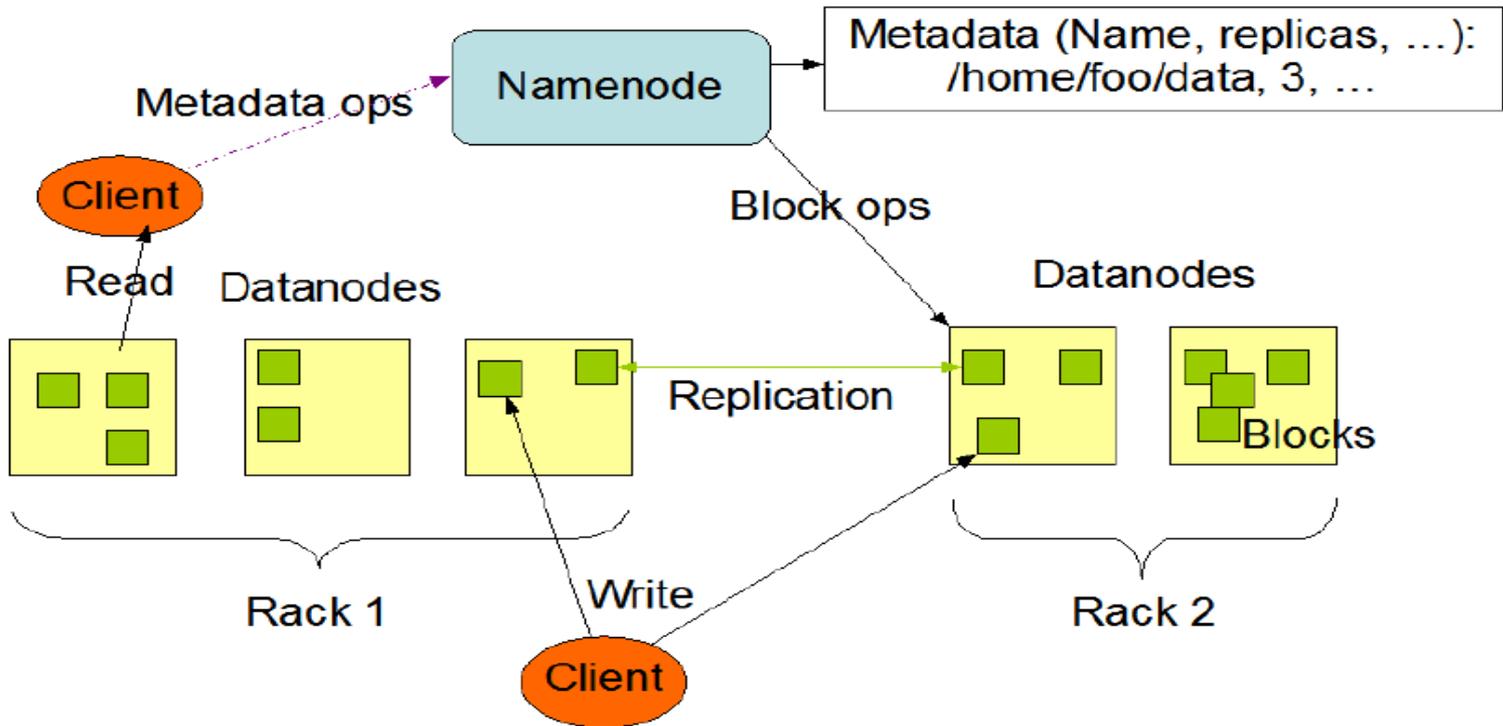
# Apache Hadoop

- Detecção de falhas e recuperação automática das mesmas
- Desenhado para processamento em “batch”
- Desenhado para trabalhar com grandes Conjuntos de dados
- Modelo simples de coerência . Modelo Map-Reduce adapta-se perfeitamente. Write once / read many
- Portabilidade entre plataformas
- Replicação de dados
- Robustez

# Apache Hadoop

Arquitetura do tipo Mestre / Escravo

HDFS Architecture



# MongoDB

- Base de dados orientada a documentos – (NoSql) (Dados semi-estruturados)
- Embora não seja um Modelo Relacional disponibiliza muitas das funções das bases de dados relacionais
- Documentos “Serializados” como JSON e guardados como BSON.
- Usa “sharding” - técnica de dividir dados através do cluster de modo a disponibilizar acesso em paralelo
- Existe MongoDB-Hadoop Connector

# MongoDB-Hadoop Connector

- Permite a ligação de MongoDB ao Hadoop ao invés de HDFS
- MongoDB como fonte de dados
- Permite ao utilizador efectuar consultas, dividindo o resultado em “input splits”.
- Para Servidores com suporte de “sharding” o split é de 64 MB

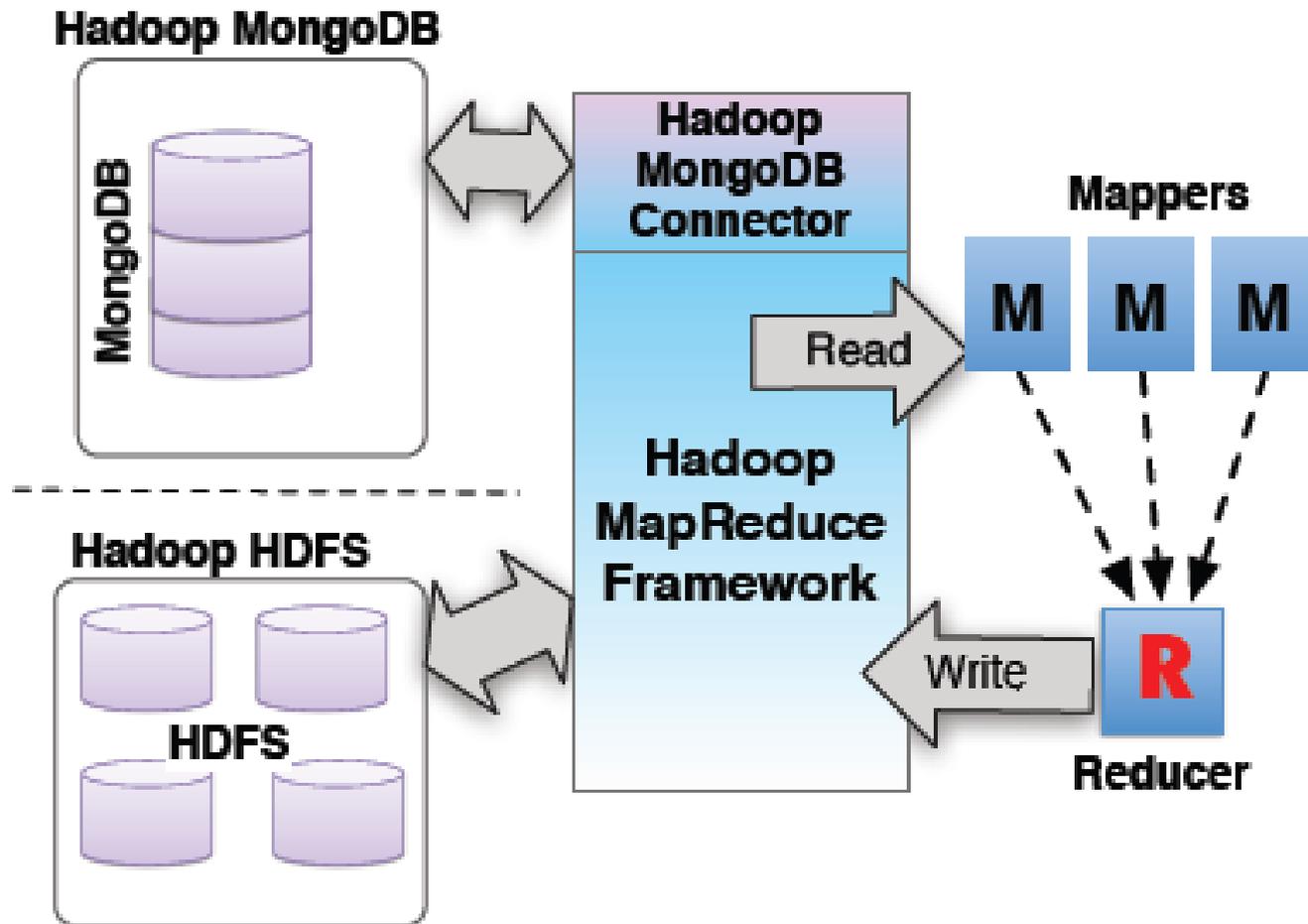
# HDFS vs MongoDB

Característica	HDFS	MongoDB
Armazenamento	Sistema de ficheiros distribuído	BD distribuída, em memória em cada nó
Leitura	Sequencial, por blocos	Sequencial e aleatório Btree Índices
Escrita	Em cache local, enviada depois para o DataNode	Escrita para índices e blocos de dados Bloqueio de escrita global por servidor
Confiança/Fiabilidade	Replicação	Replicação

HDFS – otimizado para operações de leitura e escrita sequenciais de dados em blocos relativamente grandes

MongoDB - otimizada para acesso aleatório e paralelo a dados, ou seja “queries” a dados

# Estrutura do Sistema para Avaliação



# DATASET UTILIZADO

- Tamanho aproximado - 300GB
- Dataset de census (U.S.A.), onde cada registo tem 111 campos separados por virgula.
- Uso de sub-datasets nesta experiência.