# The Hadoop Distributed File System (HDFS):
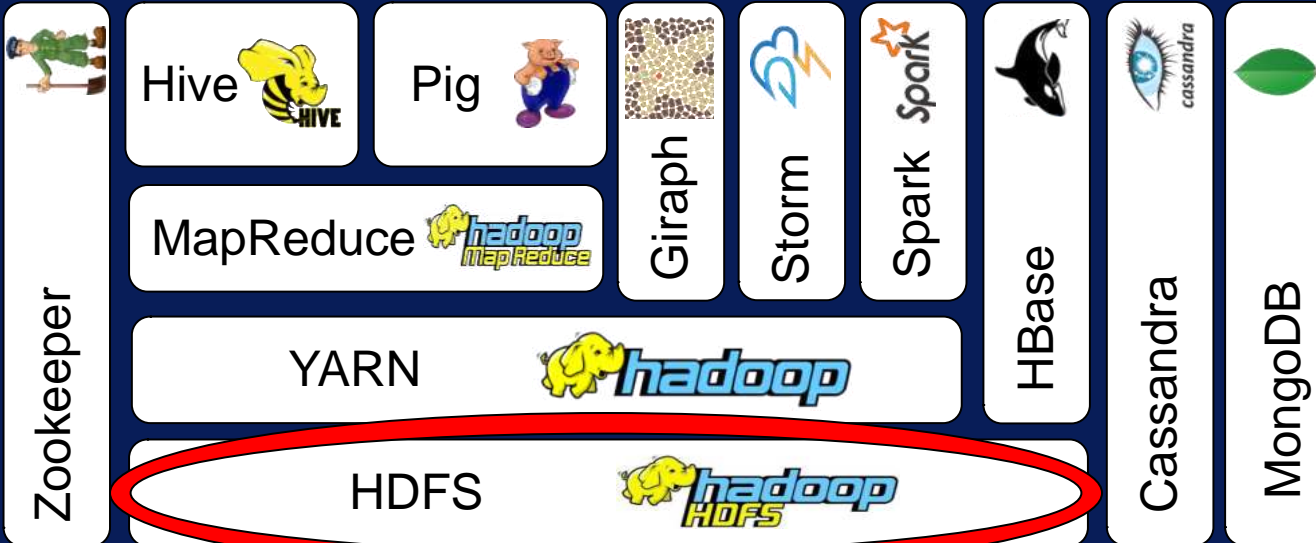
# A Storage System for Big Data

# After this video you will be able to..

- Describe how HDFS provides scalable and reliable storage

- Differentiate two key HDFS components: the NameNode and the DataNode
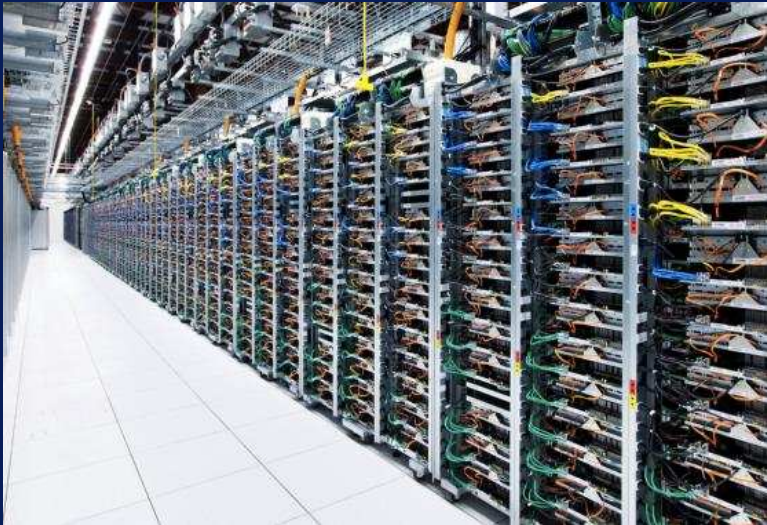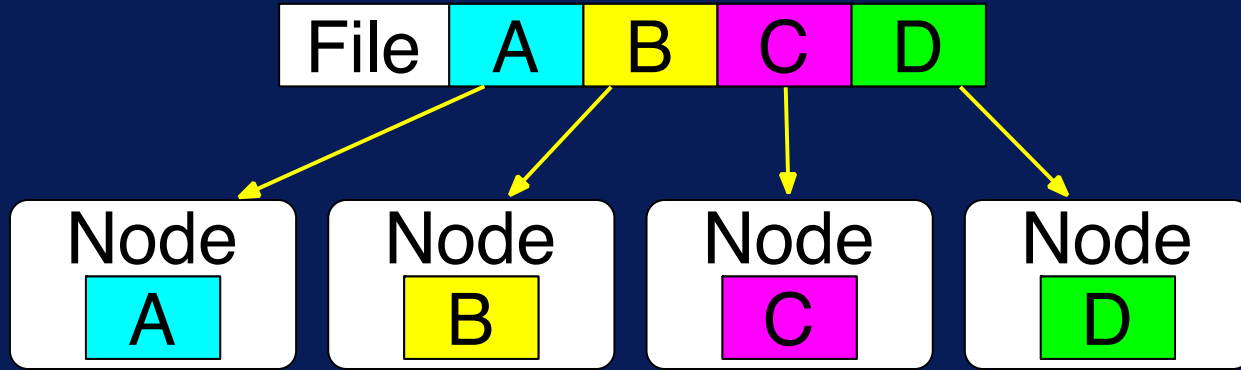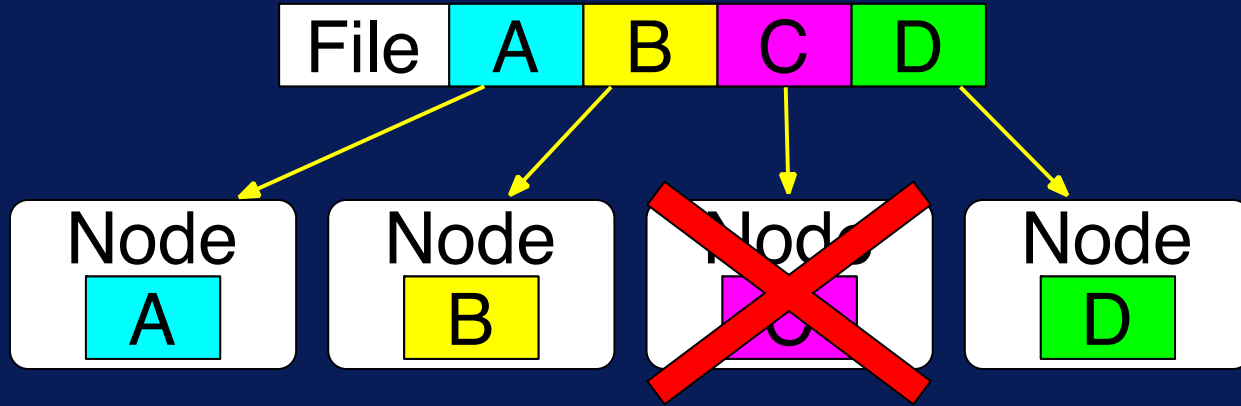
# Store massively large data sets

up to 200 Petabytes,
4500 servers,
1 billion files and blocks!

HDFS splits files across nodes for parallel access
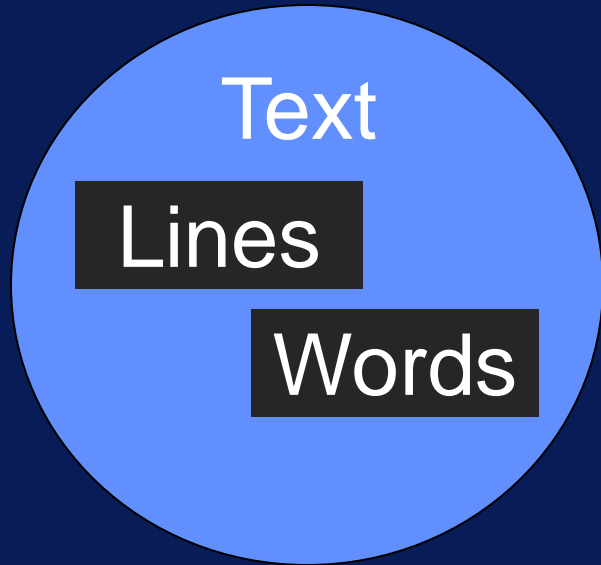
File | A | B | C | D

Node A

Node B

Node C

Node D

# What happens if node fails?

File | A | B | C | D

Node A

Node B

Node ~~C~~

Node D

# Replication for fault tolerance

# Customized reading to handle *variety* of file types

Text

Lines

Words

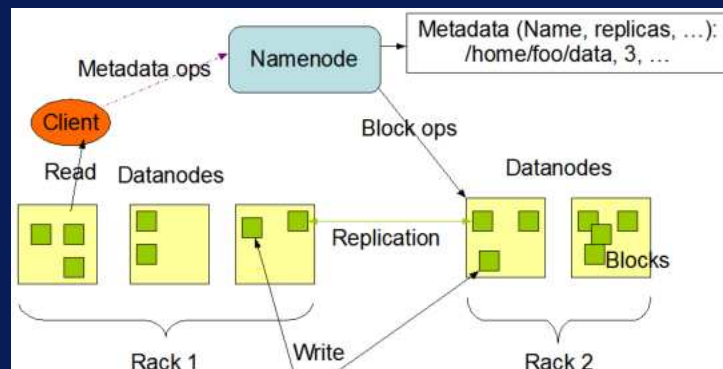Customized reading to handle *variety* of file types

Text

Lines

GIS

Vectors

Rasters

# Customized reading to handle *variety* of file types

Text

GIS

Bio

Line

Vec

FASTA

FASTQ

# Two key components of HDFS



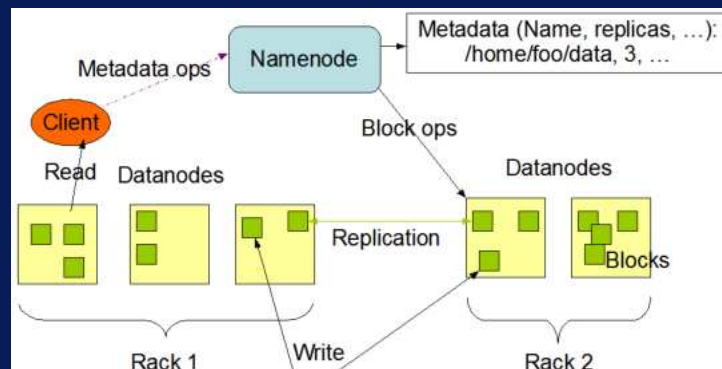1. NameNode for metadata

2. DataNode for block storage

# Two key components of HDFS



1. NameNode for metadata

   *Usually one per cluster*

2. DataNode for block storage

   *Usually one per machine*

The NameNode coordinates operations

Keeps track of file name, location in directory, etc.

Mapping of contents on DataNode.

# DataNode stores file blocks

Listens to NameNode for block creation, deletion, replication
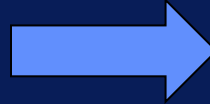
DataNode stores file blocks

Listens to NameNode for block creation, deletion, replication
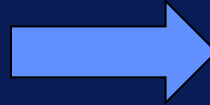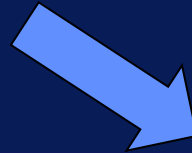
Fault Tolerance

Data locality