

# Data-parallel processing with Hadoop



# Outline

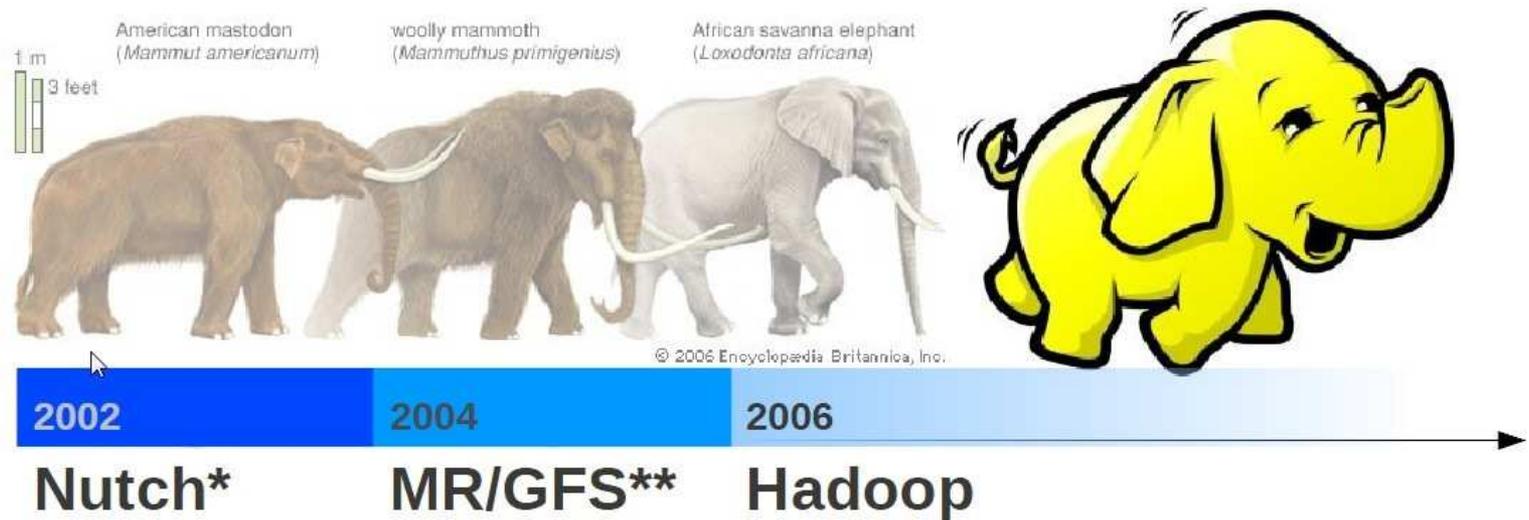
- ❑ **Big Data**
- ❑ **Distributed File System (HDFS)**
- ❑ **MapReduce**
- ❑ **Hadoop Ecosystem**

# Big Data

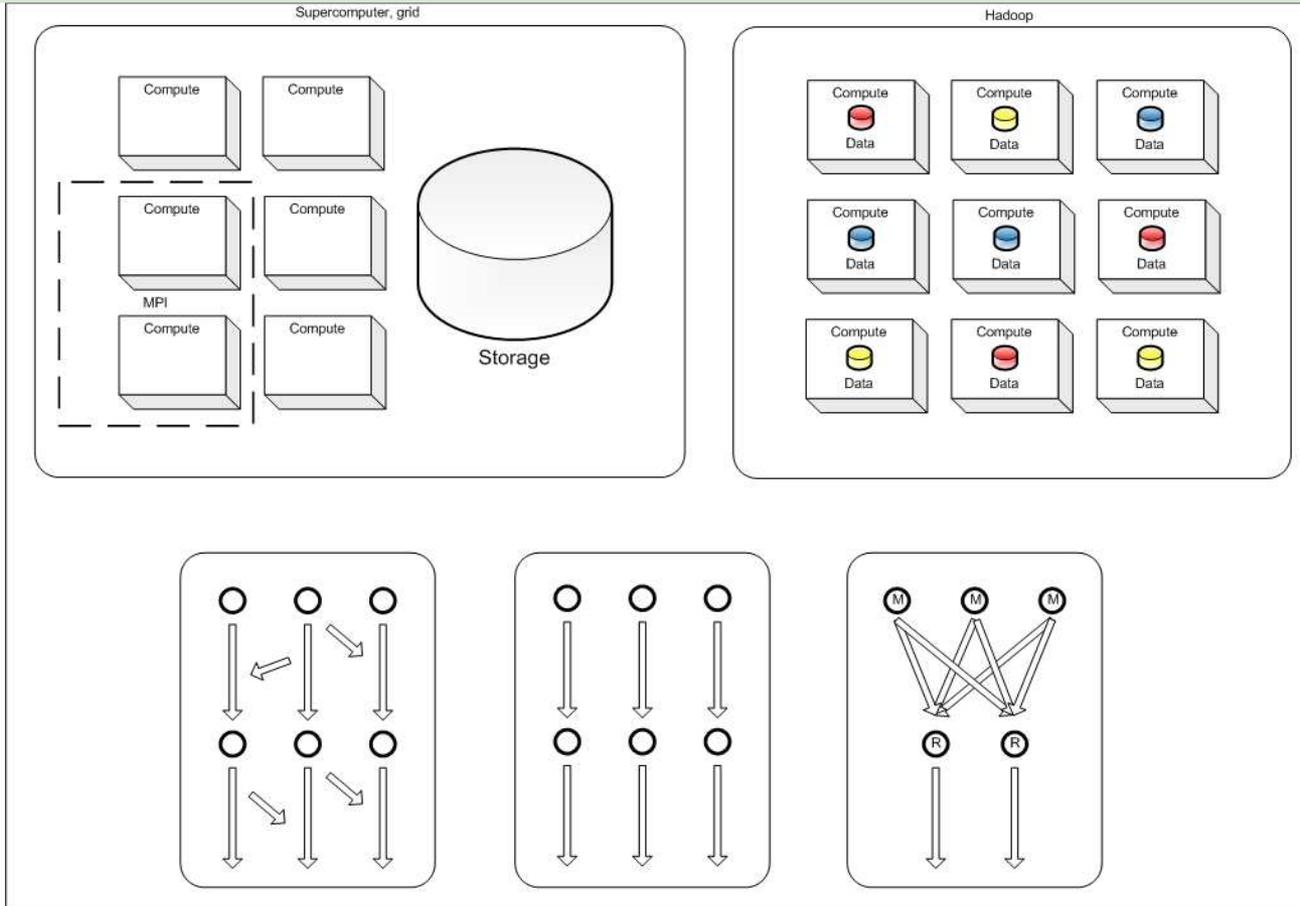
- ❑ **Growing data volumes**
  - How to store this cost-effectively
  - How to scan all data in a reasonable time
  
- ❑ **Not (just) about Volume**
  - Variety: multiple sources, multiple formats
  - Velocity: speed of data in / out
  - Veracity: uncertainty of data

# Hadoop's Google roots

- Hadoop is an open source Java implementation of Google's DFS & MapReduce
  - *The Google File System (2003)*  
Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung
  - *MapReduce: Simplified Data Processing on Large Clusters (2004)*  
Jeffrey Dean and Sanjay Ghemawat



# Cluster Architecture



# Distributed File System

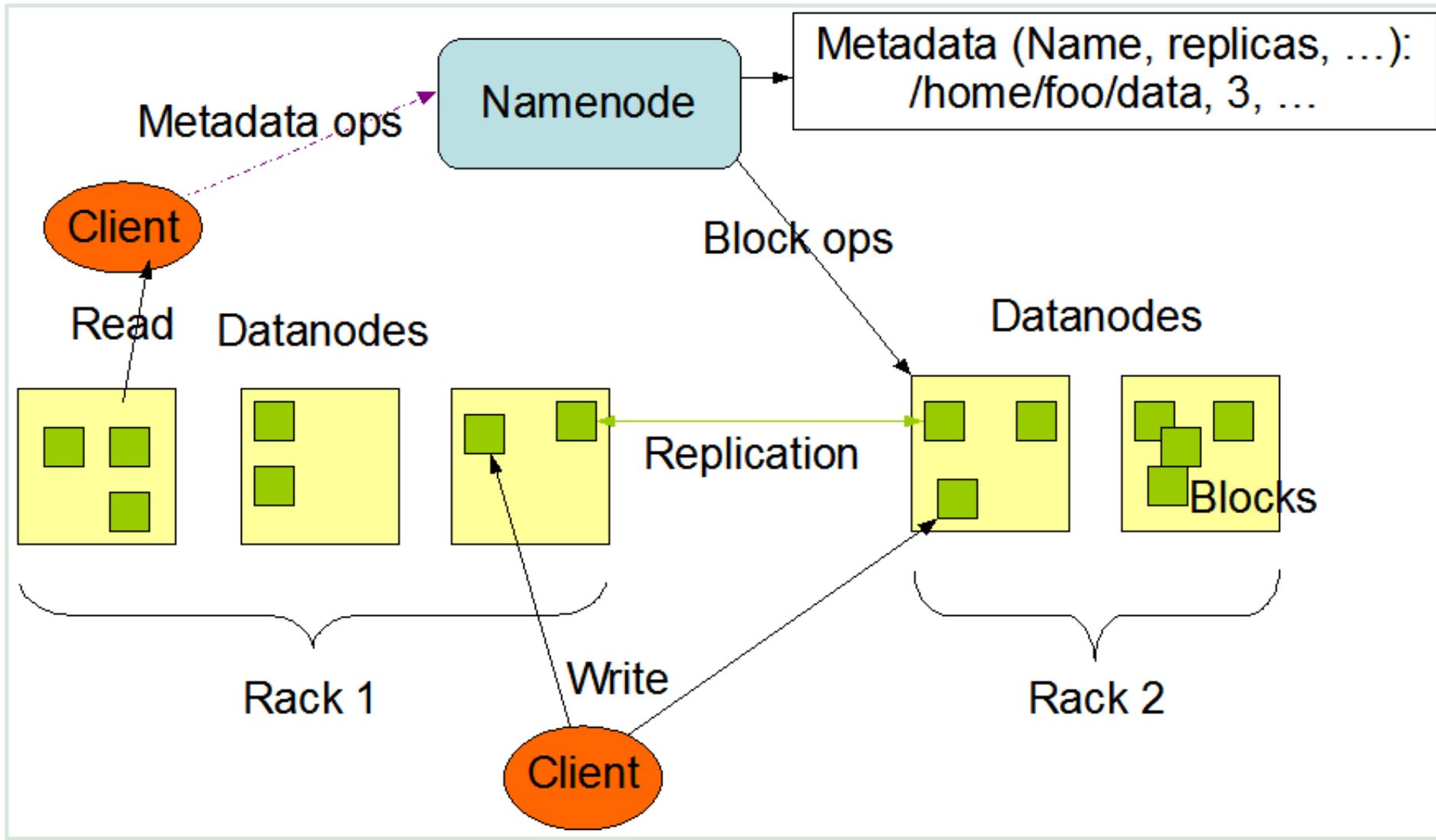
## ❑ Conventional cluster

- Data stored on dedicated storage elements
- I/O performance limited by SE bandwidth
- Low-level (such as RAID) to guarantee availability

## ❑ HDFS file system

- Compute elements are also storage elements
- Data is distributed / can be accessed in parallel
- Multiple copies (3) of every data block to protect against single node failure

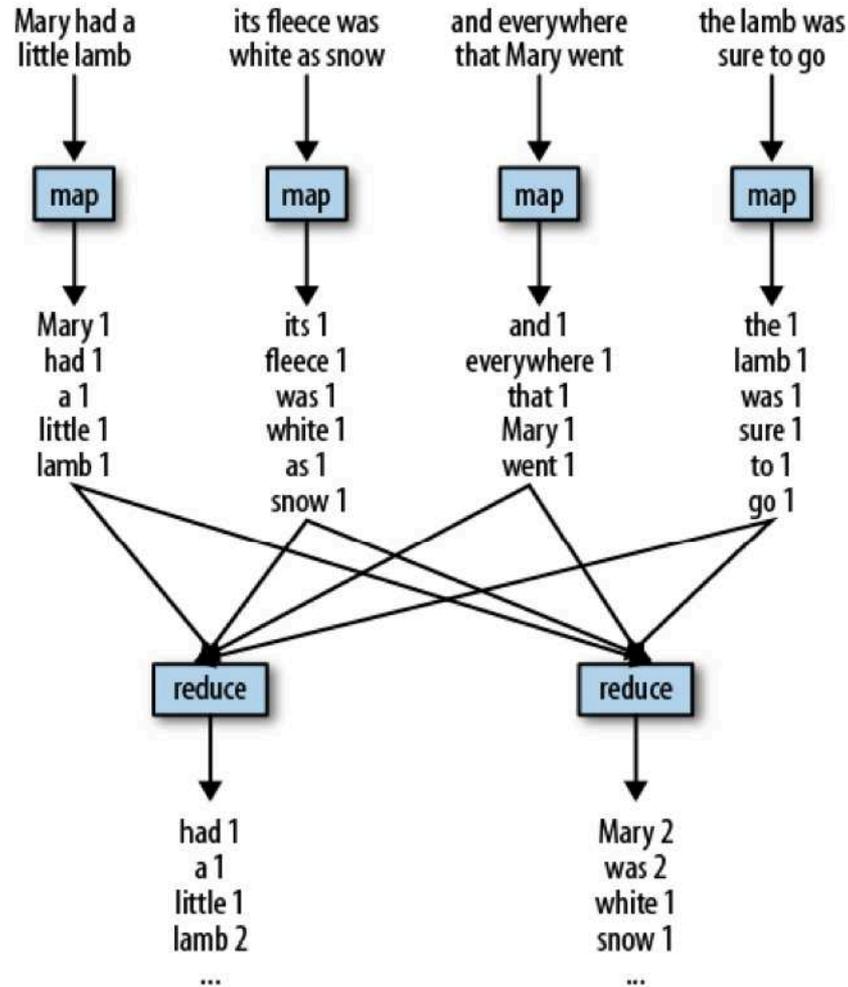
# HDFS



# MapReduce

- ❑ **Natural approach to independent data processing**
  - Read data sequentially: streaming reads faster than random reads
  - Treat the data as records
    - ✓ For every record, apply some transformation => **MAP**
    - ✓ Group and sort intermediate results
    - ✓ For every grouping, aggregate the results => **REDUCE**
  - Bring computation to the data
  - Detect failure and restart only the failed processes
  - The MapReduce framework handles this for the developer

# MapReduce



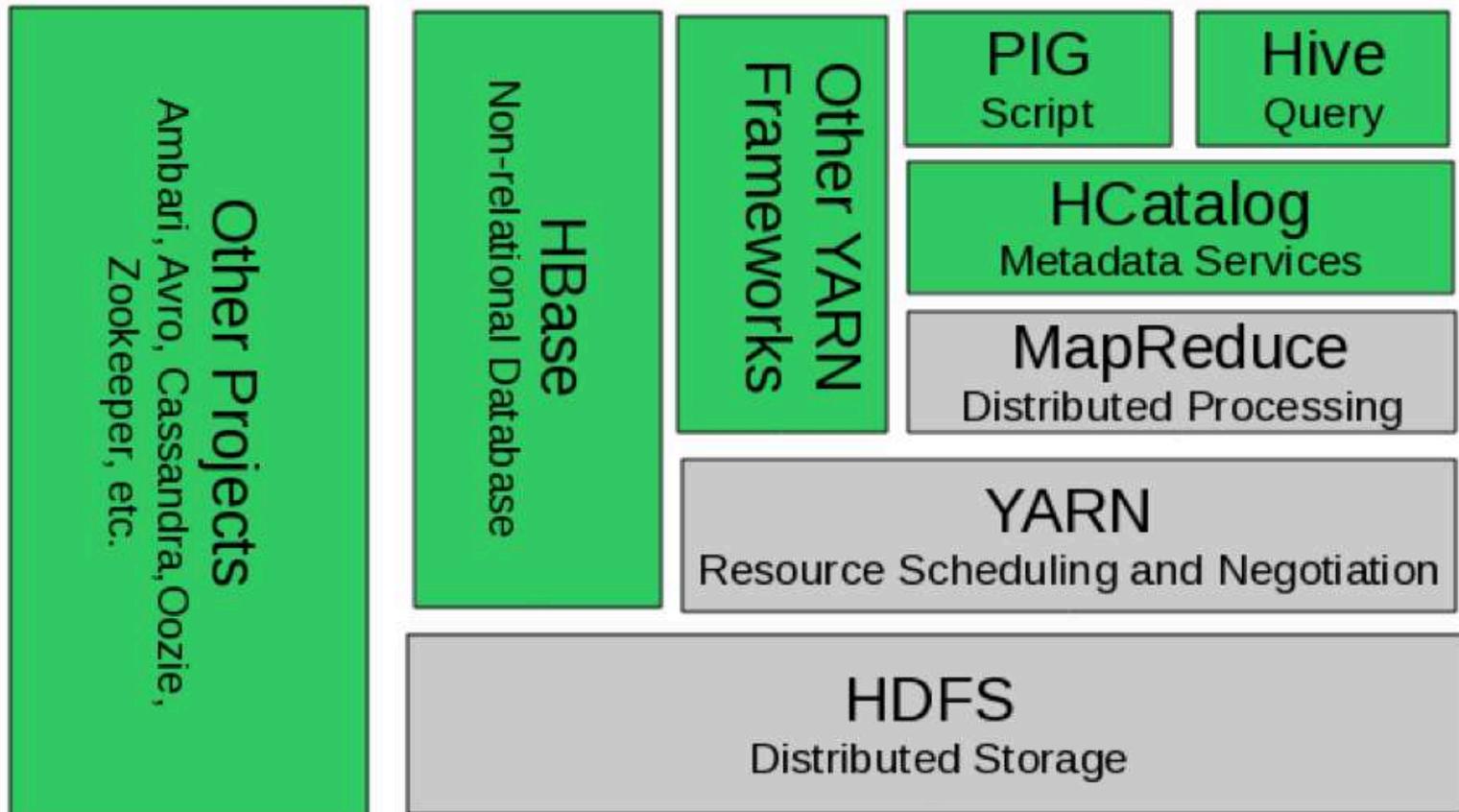
# Example application

- ❑ Generate word statistics for all Wikipedia articles
- ❑ Create an inverted index of all words



**WIKIPEDIA**  
The Free Encyclopedia

# Hadoop Ecosystem



# Summary

- ❑ **HDFS offers a scalable solution to data storage**
- ❑ **MapReduce good fit when:**
  - Data can be process independent
  - Small number of iterations
- ❑ **Hadoop ecosystem offers other tools:**
  - SQL: Hive
  - Big Table: HBase
  - Real-time processing: Storm

# Getting Started

- ❑ **Hortonworks Sandbox**
  - <http://hortonworks.com/products/hortonworks-sandbox/>
- ❑ **Amazon EMR (Elastic MapReduce)**
  - <http://aws.amazon.com/elasticmapreduce/>
- ❑ **Dutch researchers: SURFsara Hadoop cluster**
  - <https://www.surfsara.nl/project/hadoop>
- ❑ **MapReduce will return later in the course**