

## PROGRAMMING MODELS FOR BIG DATA :

1. Support Big Data Operations
  - Split volumes of data
  - Access data fast
  - Distribute computations to nodes
2. Handle fault Tolerance
  - Replicate Data Partitions
  - Recover files when needed
3. Enable adding more Racks/Resources
4. Optimized for specific data types
  - Document, Table, key-value, Graphs, Multimedia, stream.

MAP REDUCE



A PROGRAMMING MODEL FOR BIG DATA



MANY IMPLEMENTATIONS

## THE HADOOP ECOSYSTEM

HDFS : Hadoop Distributed File System

→ Scalable Storage

→ fault tolerance

→ Yahoo

Yarn : Flexible scheduling and resource management

MapReduce : Programming model that simplifies parallel computing

→ Google

map → apply()

Reduce → summarize()

→ Yahoo  
Pig : Dataflow scripting

Hive : SQL-like queries

→ Facebook

Giraph : built for processing large scale graphs efficiently

→ Facebook uses it to analyse the social graphs of its users.

Storm

Spark

Flink

Real time & in-memory processing

→ used in Facebook's messaging platform

HBase

Cassandra

MongoDB

- NoSQL for non-files  
- Key-values  
- Sparse tables

→ created at Facebook

Zookeeper : for management

- synchronization
- configuration
- high availability

created by  
Yahoo to wrangle  
services named  
after animals