

# Big Data – Hadoop cluster



Jens Baetens

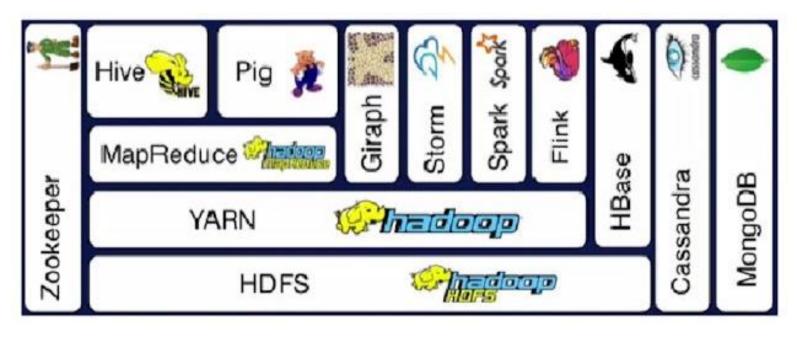
# Hadoop

- Gebaseerd op Google File System (2003)
- Ontwikkeld door Apache
- Open source
- Uitgegroeid tot omgeving met veel verschillende applicaties
  - Elke applicatie runt typisch in aparte containers/servers

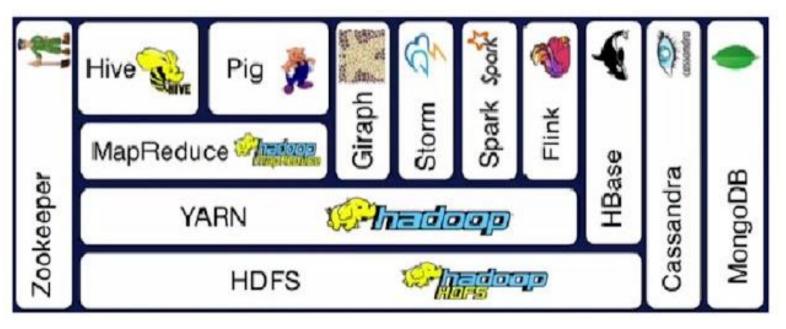




- HDFS core functionality
- Distributed File System
- Op HDD



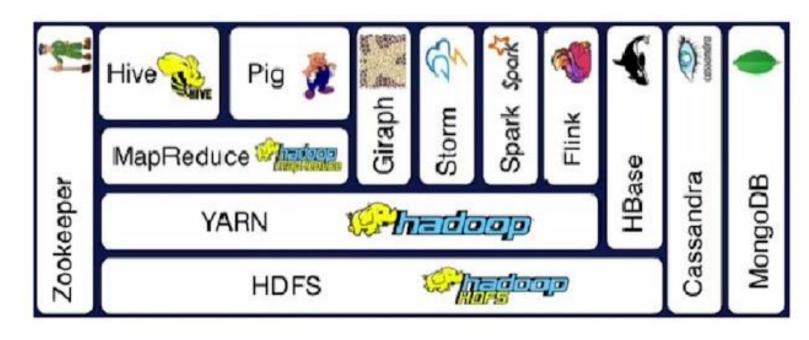
- YARN Yet Another Resource Negotiator
- Beheer van computing power
- Welke code op welke node



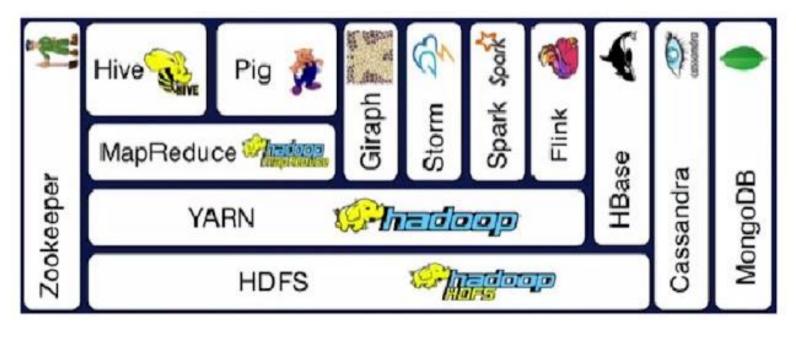
- MapReduce
- Distributed Computing
- Ontwikkeld door Google
- 2 fases
  - Mapping (Divide)
  - Reduce (Conquer)



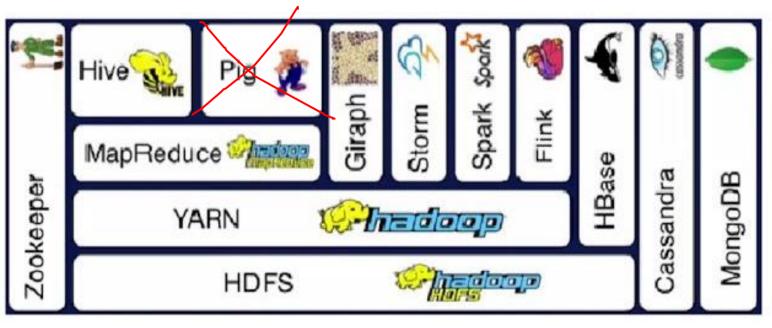
- Zookeeper
- Beheren van alle applicaties die lopen op de verschillende nodes



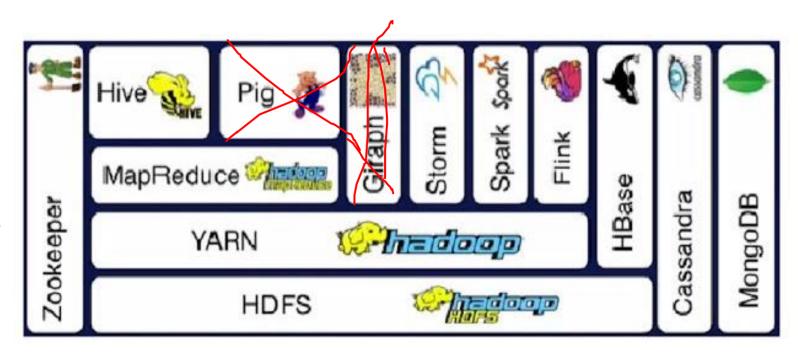
- Hive
- Distributed Datawarehouse
- Sql-like
- Queries via MapReduce



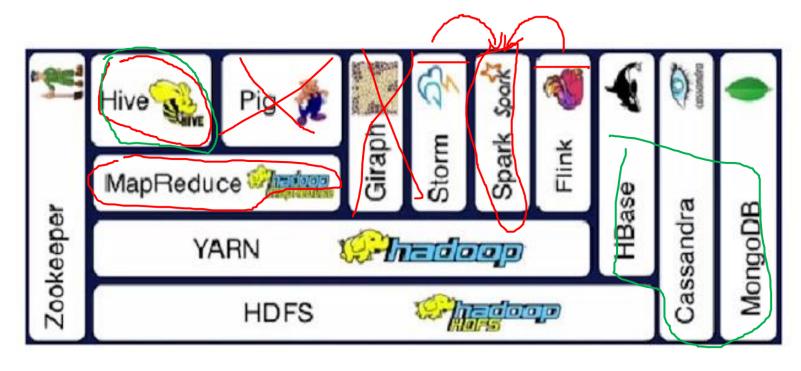
- Pig
- Data analysis
- Using MapReduce/Spark/...
- Taal: Pig Latin



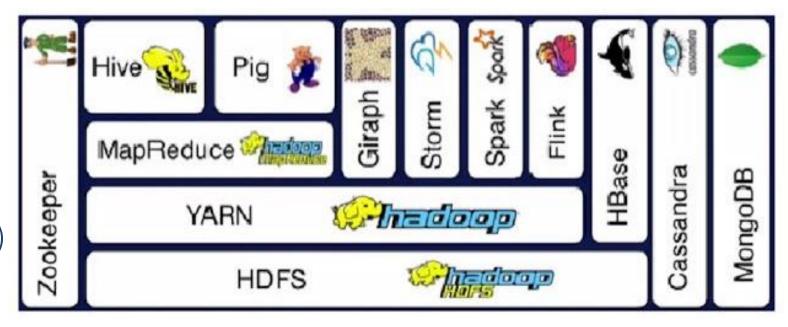
- Giraph
- Bestuderen van een graaf
- Social graph
  - Facebook
  - Twitter
- Gebruikt geen mapreduce



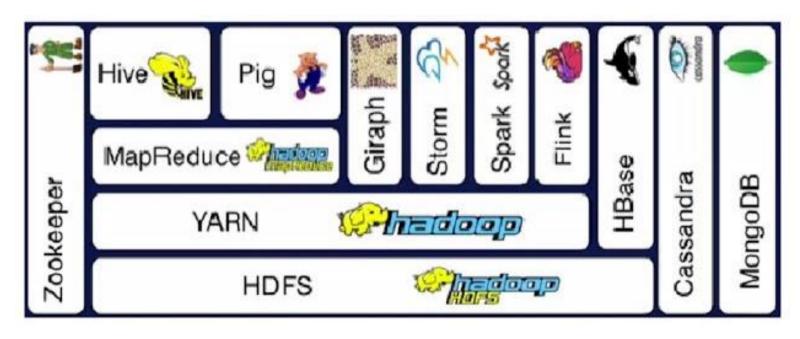
- Storm / Flink
- Verwerken van data streams continue inkomende datastromen
  - Classificeren
  - Opslaan



- Spark
- Alternatief voor MapReduce
- Computing in Ram
- Op Hadoop/Cloud/...
- Gebruikt voor
  - SQL (Spark SQL)
  - Streaming (Spark Streaming)
  - Machine Learning (MLlib)
  - Graph analysis (GraphX)



- **■** HBase
- Distributed NoSQL Database
- Geen SQL maar in JAVA



- Cassandra / Mongo DB
- Maken geen gebruik van HDFS
- NoSql databases
- Stand-alone solutions

