

Hadoop & its Usage at Facebook

[Dhruba Borthakur](#)

Project Lead, Hadoop Distributed File System

dhruba@apache.org

Presented at the The Israeli Association of Grid Technologies

July 15, 2009



Outline

- Architecture of Hadoop Distributed File System
- Synergies between Hadoop and Condor
- Hadoop Usage at Facebook



Who Am I?

- **Hadoop FileSystem Project Lead**
 - Core contributor since Hadoop's infancy
- **Facebook** (Hadoop, Hive, Scribe)
- **Yahoo!** (Hadoop in Yahoo Search)
- **Veritas** (San Point Direct, Veritas File System)
- **IBM Transarc** (Andrew File System)
- **UW Computer Science Alumni** (Condor Project)



Hadoop, Why?

- **Need to process Multi Petabyte Datasets**
- **Expensive to build reliability in each application.**
- **Nodes fail every day**
 - Failure is expected, rather than exceptional.
 - The number of nodes in a cluster is not constant.
- **Need common infrastructure**
 - Efficient, reliable, Open Source Apache License
- **The above goals are same as Condor, but**
 - Workloads are IO bound and not CPU bound



Hadoop History

- **Dec 2004** — Google GFS paper published
- **July 2005** — Nutch uses MapReduce
- **Feb 2006** — Starts as a Lucene subproject
- **Apr 2007** — Yahoo! on 1000-node cluster
- **Jan 2008** — An Apache Top Level Project
- **Jul 2008** — A 4000 node test cluster
- **May 2009** — Hadoop sorts Petabyte in 17 hours



Who uses Hadoop?

- Amazon/A9
- Facebook
- Google
- IBM
- Joost
- Last.fm
- New York Times
- PowerSet
- Veoh
- Yahoo!



What is Hadoop used for?

- Search
 - Yahoo, Amazon, Zvents
- Log processing
 - Facebook, Yahoo, ContextWeb. Joost, Last.fm
- Recommendation Systems
 - Facebook
- Data Warehouse
 - Facebook, AOL
- Video and Image Analysis
 - New York Times, Eyealike

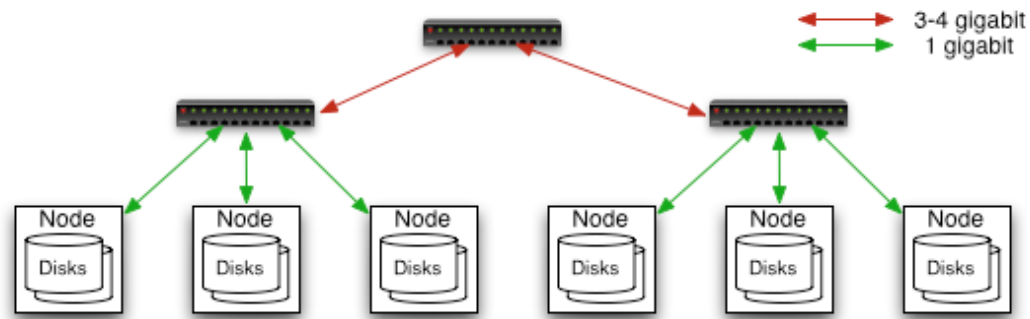


Public Hadoop Clouds

- Hadoop Map-reduce on Amazon EC2
 - <http://wiki.apache.org/hadoop/AmazonEC2>
- IBM Blue Cloud
 - Partnering with Google to offer web-scale infrastructure
- Global Cloud Computing Testbed
 - Joint effort by Yahoo, HP and Intel



Commodity Hardware



Typically in 2 level architecture

- Nodes are commodity PCs
- 30-40 nodes/rack
- Uplink from rack is 3-4 gigabit
- Rack-internal is 1 gigabit

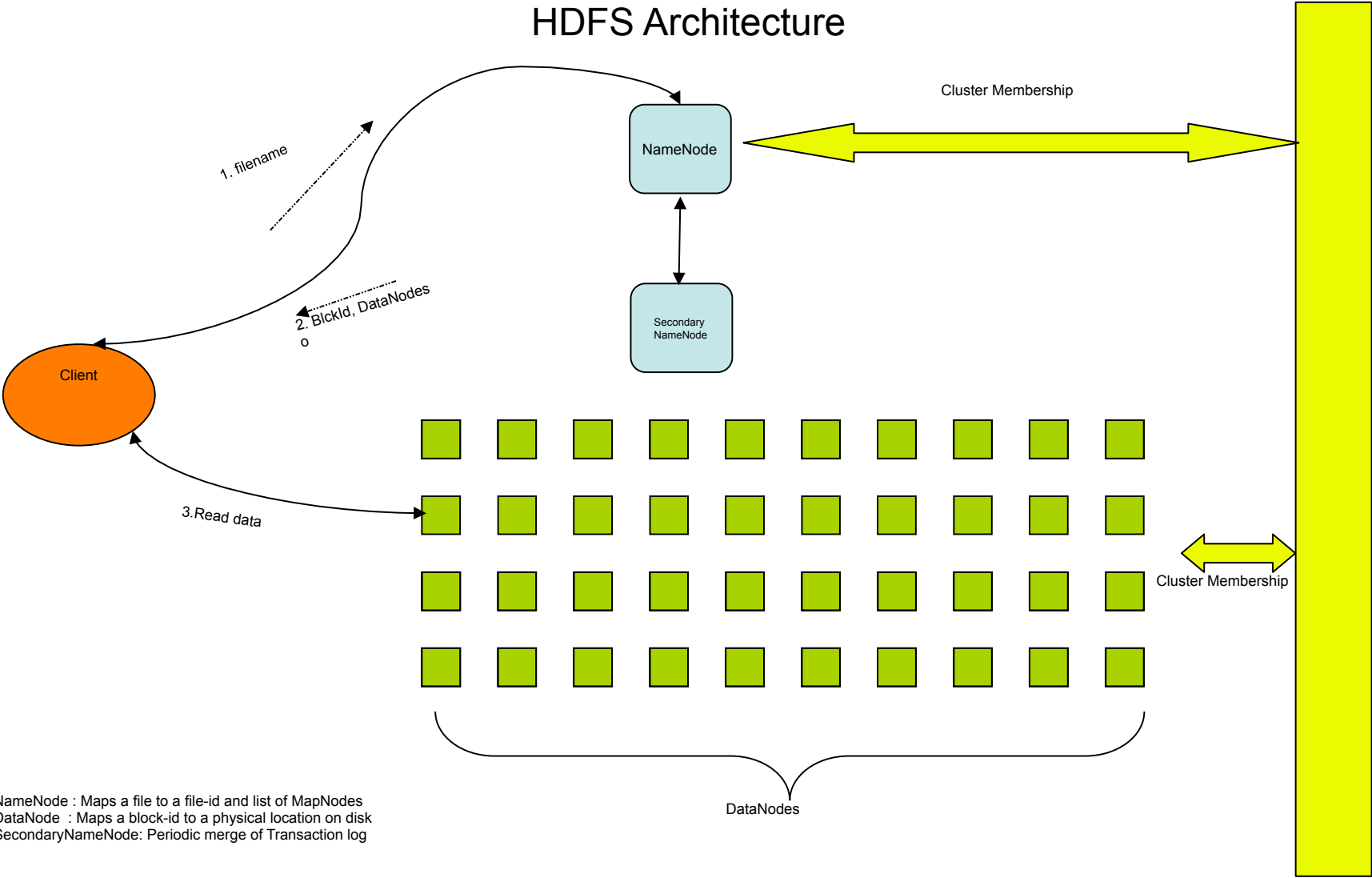


Goals of HDFS

- **Very Large Distributed File System**
 - 10K nodes, 100 million files, 10 PB
- **Assumes Commodity Hardware**
 - Files are replicated to handle hardware failure
 - Detect failures and recovers from them
- **Optimized for Batch Processing**
 - Data locations exposed so that computations can move to where data resides
 - Provides very high aggregate bandwidth
- **User Space, runs on heterogeneous OS**



HDFS Architecture



NameNode : Maps a file to a file-id and list of MapNodes
DataNode : Maps a block-id to a physical location on disk
SecondaryNameNode: Periodic merge of Transaction log

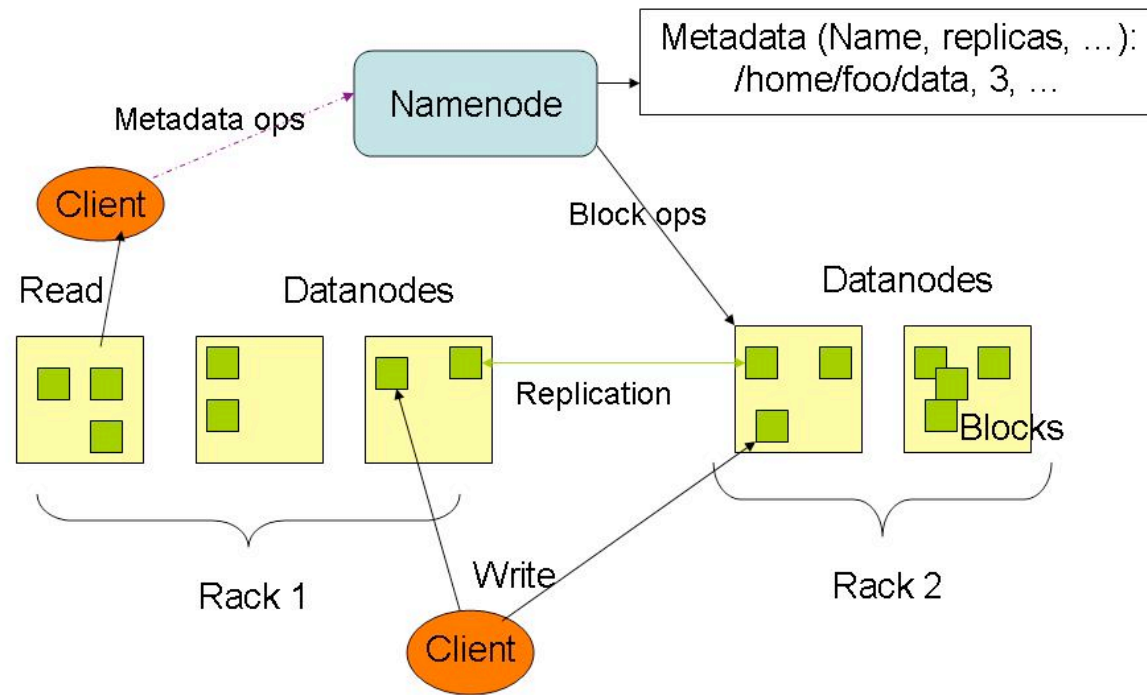


Distributed File System

- **Single Namespace for entire cluster**
- **Data Coherency**
 - Write-once-read-many access model
 - Client can only append to existing files
- **Files are broken up into blocks**
 - Typically 128 MB block size
 - Each block replicated on multiple DataNodes
- **Intelligent Client**
 - Client can find location of blocks
 - Client accesses data directly from DataNode



HDFS Architecture



NameNode Metadata

- **Meta-data in Memory**
 - The entire metadata is in main memory
 - No demand paging of meta-data
- **Types of Metadata**
 - List of files
 - List of Blocks for each file
 - List of DataNodes for each block
 - File attributes, e.g creation time, replication factor
- **A Transaction Log**
 - Records file creations, file deletions. etc



DataNode

- **A Block Server**
 - Stores data in the local file system (e.g. ext3)
 - Stores meta-data of a block (e.g. CRC)
 - Serves data and meta-data to Clients
- **Block Report**
 - Periodically sends a report of all existing blocks to the NameNode
- **Facilitates Pipelining of Data**
 - Forwards data to other specified DataNodes



Data Correctness

- **Use Checksums to validate data**
 - Use CRC32
- **File Creation**
 - Client computes checksum per 512 byte
 - DataNode stores the checksum
- **File access**
 - Client retrieves the data and checksum from DataNode
 - If Validation fails, Client tries other replicas



NameNode Failure

- **A single point of failure**
- **Transaction Log stored in multiple directories**
 - A directory on the local file system
 - A directory on a remote file system (NFS/CIFS)
- **Need to develop a real HA solution**



Rebalancer

- **Goal: % disk full on DataNodes should be similar**
 - Usually run when new DataNodes are added
 - Cluster is online when Rebalancer is active
 - Rebalancer is throttled to avoid network congestion
 - Command line tool



Hadoop Map/Reduce

- **The Map-Reduce programming model**
 - Framework for distributed processing of large data sets
 - Pluggable user code runs in generic framework
- **Common design pattern in data processing**
cat * | grep | sort | unique -c | cat > file
input | **map** | shuffle | **reduce** | output
- **Natural for:**
 - Log processing
 - Web search indexing
 - Ad-hoc queries



Hadoop and Condor



Condor Jobs on HDFS

- **Run Condor jobs on Hadoop File System**
 - Create HDFS using local disk on condor nodes
 - Use HDFS API to find data location
 - Place computation close to data location
- **Support map-reduce data abstraction model**



Job Scheduling

- **Current state of affairs with Hadoop scheduler**
 - FIFO and Fair Share scheduler
 - Checkpointing and parallelism tied together
- **Topics for Research**
 - Cycle scavenging scheduler
 - Separate checkpointing and parallelism
 - Use resource matchmaking to support heterogeneous Hadoop compute clusters
 - Scheduler and API for MPI workload



Dynamic-size HDFS clusters

- **Hadoop Dynamic Clouds**
 - Use Condor to manage HDFS configuration files
 - Use Condor to start HDFS DataNodes
 - Based on workloads, Condor can add additional DataNodes to a HDFS cluster
 - Condor can move DataNodes from one HDFS cluster to another



Condor and Data Replicas

- **Hadoop Data Replicas and Rebalancing**
 - Based on access patterns, Condor can increase number of replicas of a HDFS block
 - If a condor job accesses data remotely, should it instruct HDFS to create a local copy of data?
 - Replicas across data centers (Condor Flocking?)



Condor as HDFS Watcher

- **Typical Hadoop periodic jobs**
 - Concatenate small HDFS files into larger ones
 - Periodic checksum validations of HDFS files
 - Periodic validations of HDFS transaction logs
 - Convert data from lzo to gzip compression
- **Condor can intelligently schedule above jobs**
 - Schedule during times of low load



HDFS High Availability

- **Use Condor High Availability**
 - Failover HDFS NameNode
 - Condor can move HDFS transaction log from old NameNode to new NameNode



Power Management

- **Power Management**
 - Major operating expense
- **Condor Green**
 - Analyze data-center heat map and shutdown DataNodes if possible
 - Power down CPU's when idle
 - Block placement based on access pattern
 - Move cold data to disks that need less power



Hadoop Cloud at Facebook



Who generates this data?

- Lots of data is generated on Facebook
 - 200+ million active users
 - 30 million users update their statuses at least once each day
 - More than 900 million photos uploaded to the site each month
 - More than 10 million videos uploaded each month
 - More than 1 billion pieces of content (web links, news stories, blog posts, notes, photos, etc.) shared each week



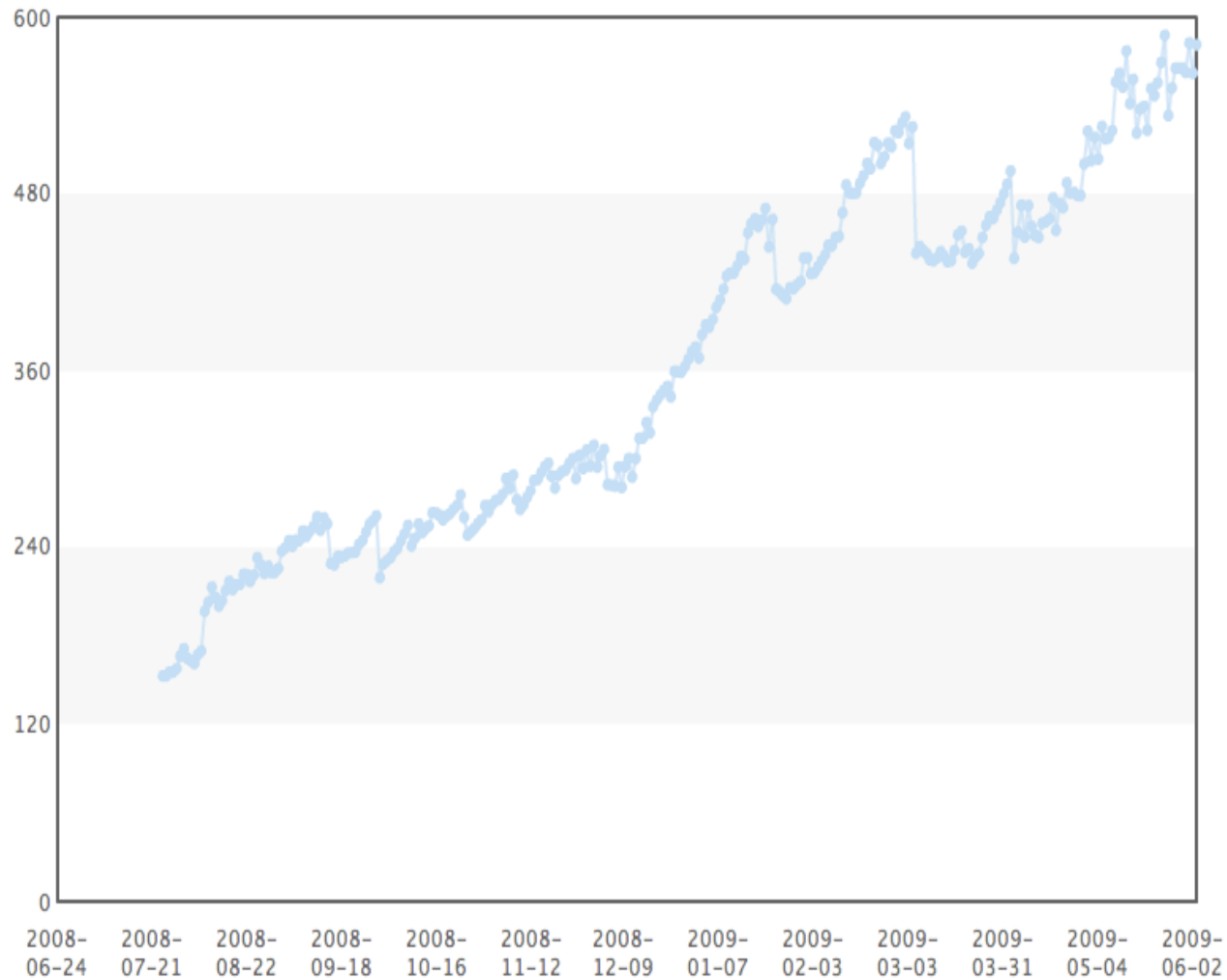
Where do we store this data?

- Hadoop/Hive Warehouse
 - 4800 cores, 2 PetaBytes (July 2009)
 - 4800 cores, 12 PetaBytes (Sept 2009)
- Hadoop Archival Store
 - 200 TB

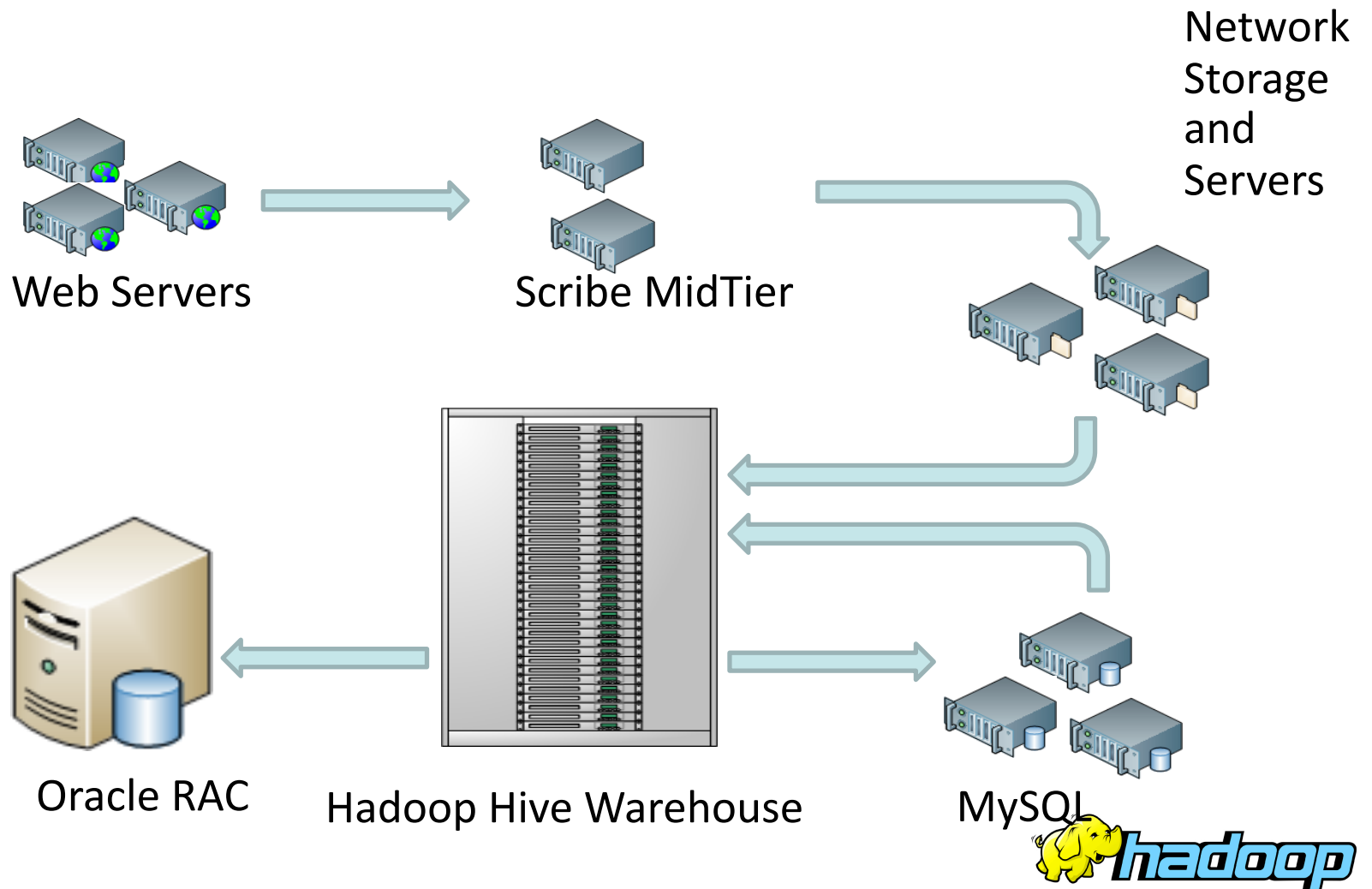


Rate of Data Growth

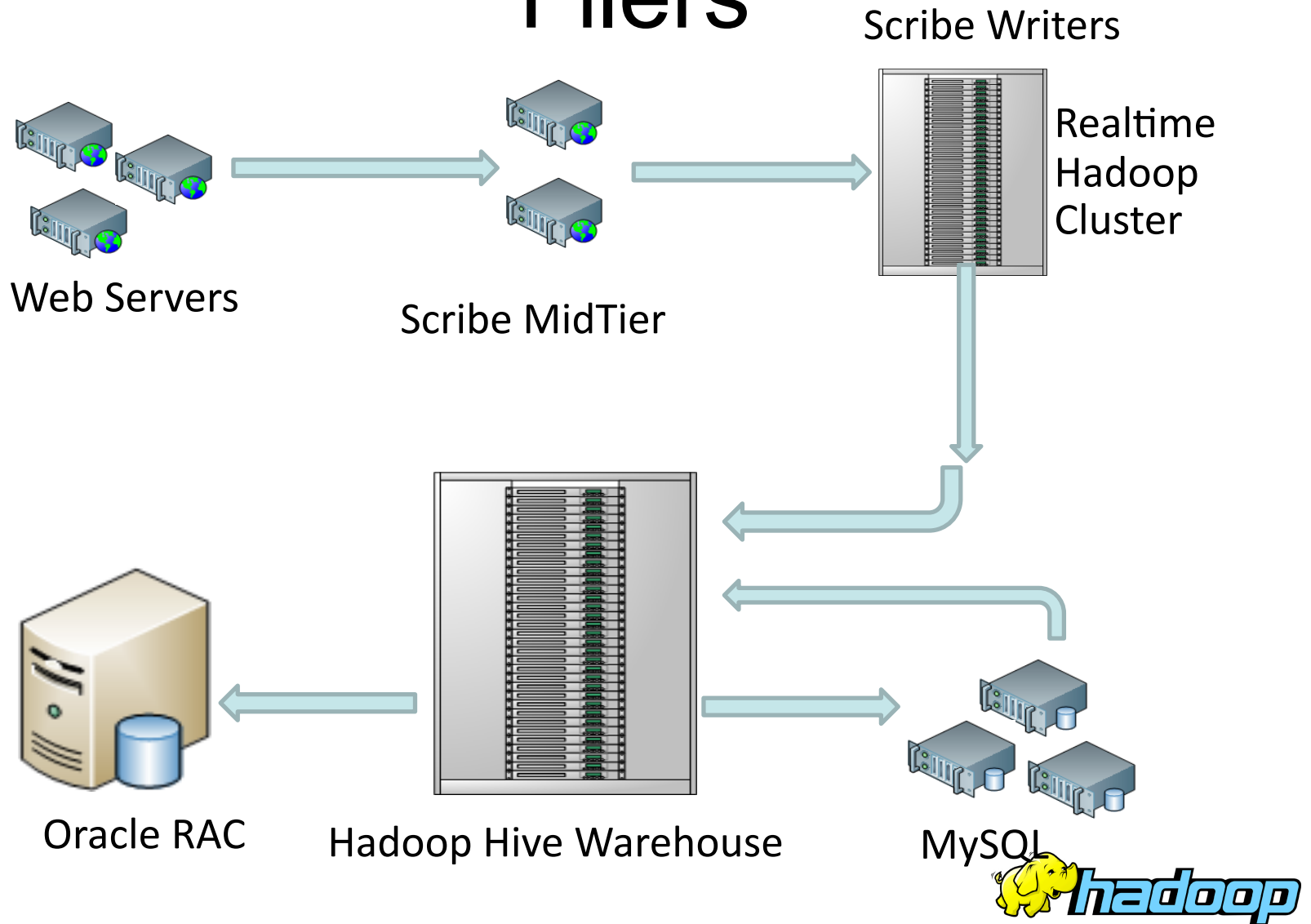
Hadoop File System Size (Terabytes) by Date



Data Flow into Hadoop Cloud



Hadoop Scribe: Avoid Costly Filers



Data Usage

- Statistics per day:
 - 4 TB of compressed new data added per day
 - 55TB of compressed data scanned per day
 - 3200+ Hive jobs on production cluster per day
 - 80M compute minutes per day
- Barrier to entry is significantly reduced:
 - New engineers go through a Hive training session
 - 140+ people run jobs on Hadoop/Hive jobs
 - Analysts (non-engineers) use Hadoop through Hive

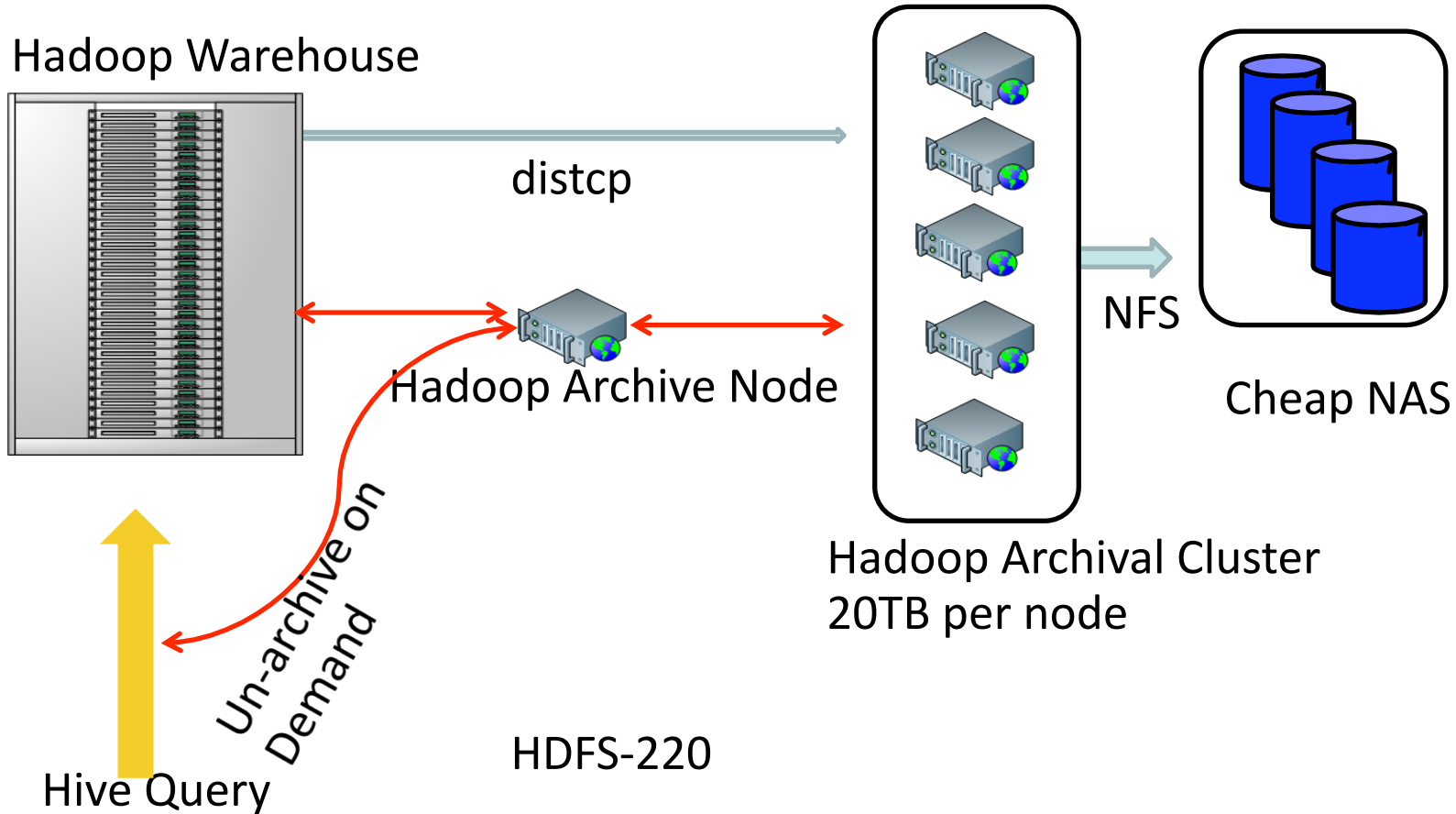


Hive Query Language

- SQL type query language on Hadoop
- Analytics SQL queries translate well to map-reduce
- Files are insufficient data management abstractions
 - Need Tables, schemas, partitions, indices



Archival: Move old data to cheap storage

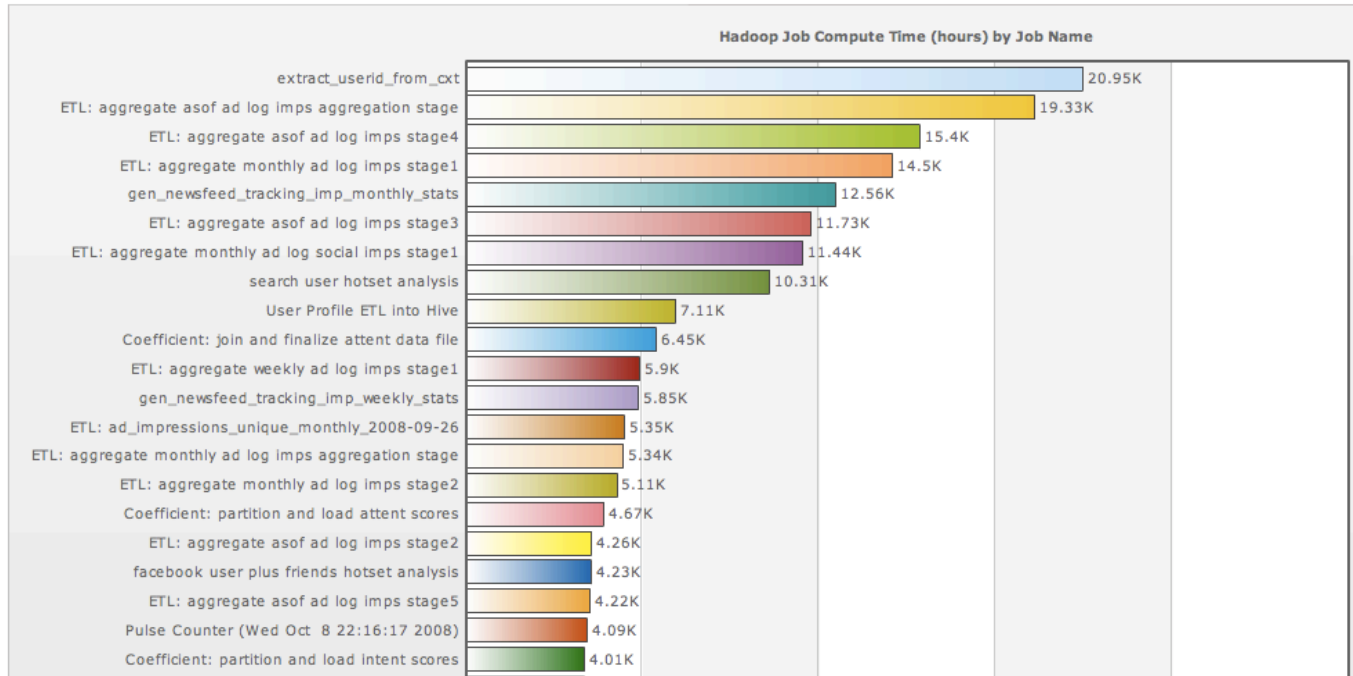


Cluster Usage Dashboard

Jobs | Compute Time | Map Time | Reduce Time | Job Durations | Map Durations | Reduce Durations | Jobs by Date | Compute Time by Date

Task Time by Date | I/O by Date | HDFS Size | HDFS Metadata | Largest Hive Tables | Largest Home Directories | Largest Facebook Project Dirs

Days back: Break down by:



Confidential Materials — For Internal Use Only
Hadoop Job Compute Time (hours) by Job Name



Add Comment



Summary

- Hadoop is the platform of choice for Storage Cloud
- Facebook a big contributor to Open Source Software
- Lots of synergy between Hadoop and Condor



Useful Links

- **HDFS Design:**
 - http://hadoop.apache.org/core/docs/current/hdfs_design.html
- **Hadoop API:**
 - <http://hadoop.apache.org/core/docs/current/api/>

