

the Hadoop distributed file system

wbg231

January 2023

1 Introduction

- hadoop is designed to store very large datasets reliably and to stream those data sets at high bandwidth to user applications
- the idea is to distribute storage and computation across many servers allowing resources to grow with demand and remain economical ie cheap lol
- was made at yahoo

Introduction and related work

- hadoop provides a distributed file system and a framework for computation using the map reduce paradigm
- hadoop works by partitioning data and computation across many hosts and keeping in parallel close to where the data is kept
- hadoop scales simply by just adding more servers
- components of hadoop
 1. hdfs the file system
 2. map reduce the framework for distributed computation
 3. hbase column oriented table service
 4. pig data flow language with parallel execution
 5. hive data warehouse infrastructure
 6. zoo keeper distributed coordination service
 7. avro data serialization service
- hadoop is an apache project
- most of the hdfs interface is modeled after unix but some changes were made to improve performance

- HDFS stores meta data on a dedicated server called the [name node](#)
- all servers are fully connected and able to communicate with another via TCP based protocols
- file is made durable by being replicated on multiple file nodes
- this also allows improvements to data transfer bandwidth as it is more likely that data will be located near the computer that needs it

architecture

name node

- the HDFS name space is a hierarchy of files, files and directories are represented on the name node by inodes which record attributes like permissions, modifications and access times
- the file content is split into large blocks and each block of the file is replicated at multiple usually three data nodes.
- the name node maintains the name space tree and mapping of file blocks to the data node (ie where the data is physically stored)
- an hdfs client wanting to read a file first contact the name node to locate the data blocks comprising the file and then reads the blocks contents from the nearest data node/ '
- when something is written to a file the name node nominates three data nodes to store the updated data.
- HDFS keeps the entire name space in ram.
- all the meta data is called the image
- the back up of the image is called a checkpoint

data nodes

- each data block replica on a data node is represented by two files on the hosts native file system.
- the first file has the data the second has the blocks meta data
- host computers can share multiple blocks so if one block only takes up as much space as it can fill
- when starting a program the name node checks that each data node has the right name space id and software versions

- the namespace id is assigned to the file system instance when it is formatted, the namespace id is present on all nodes of the cluster, so nodes with different namespace id can not join a cluster keeping integrity (ie no weird computers on your stuff)
- the name node stores the storage id of each data node which registers which is a unique identifier.
- datanode identifies block replicas to the name node by sending a block report
- the data node sends ping the name node regularly with heart beats if these do not come through for a period of time the data node is thought to be dedicated

hdfs client

- users work with hdfs using the hdfs client
- user can write or read files. gives instructions to name node writes and reads from data nodes

image and journal

- the namespace image is the file system meta data that describes the organization of data files and directories
- the journal is a commit log for changes to the file system
- the name node in addition to serving client requests can also work as either a checkpoint or backup node
- the checkpoint node combines the existing checkpoint and journal to make a new checkpoint (the checkpoint node typically runs on a different host than the name node)
- having checkpoints is a good way to ensure system data will be ok in the event of a crash

backup node

- the backup node is a read only name node that works kind of like a checkpoint node but stores the most recent journal in its local memory so if the name node goes down the information is quickly accessible

file and io operation

file read and write

- apps can add data to hdfs by creating a new file and writing data to it. after the file is closed the it can not be written to except by append's that is adding stuff to the end
- when a client tries to write to a file they ping the name node and are the only system allowed to write to that file until there transaction is done
- when file is to large for current block, the name node makes new blocks and tells certain data nodes to hold replications
- if a user tries to read a file, it is connected to the nearest data node, if that read does not work it is connected to the data node with the nearest replica and so on
- a user can read the file from a data node that is being written to, but it gets the last line (ie where that file could be modified) from another replication of the file

block placement

- nodes are organized in racks.
- nodes of a rack share a switch
- racks are connected by core switches
- communication between two nodes in different racks had to go through multiple switches
- the network bandwidth between nodes in the same rack is greater than that between racks
- network bandwidth between nodes is estimated by their distance
- when a new file is created hdfs places the first replica on the node where the writer is located. the second and third on two nodes in a different rack and the rest are placed randomly to ensure that if any single node or rack goes down that file is not lost

replication management

- the name node tries to always ensure that each block has the correct number of replicas

- it prioritizes adding replicas to blocks that have too few over deleting replicas from blocks that have too many
- always makes sure that not all replicas of a block are on a single node

balancer

- there is a balancer tool that tries to spread disk usage evenly across the hdfs cluster
- there is also a block scanner that periodically makes sure that the meta data of each block is in line with the block id
- the rest of this paper is about practices at yahoo