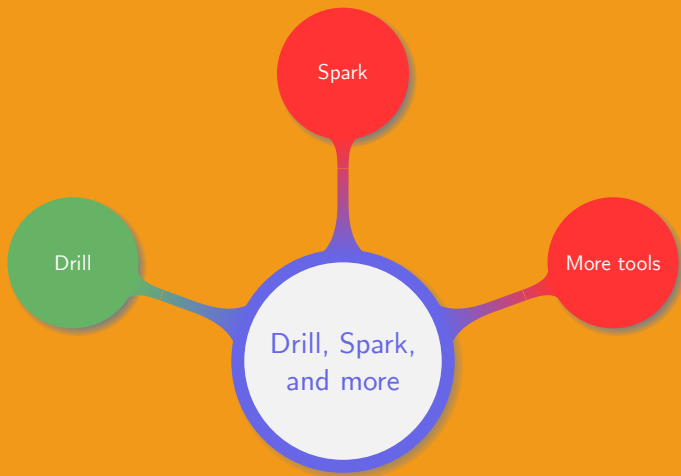


Methods and tools for big data

11. Drill, Spark, and more

Jing & Manuel – Summer 2020



Parties always involved in a Drill job:

- A client which initiates the job
- Zookeeper

Parties always involved in a Drill job:

- A client which initiates the job
- Zookeeper

Parties optionally involved in a Drill job:

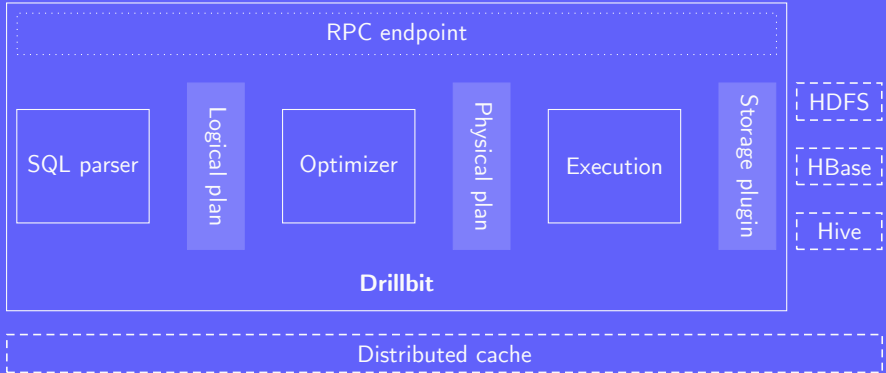
- YARN
- HDFS
- Hive
- HBase

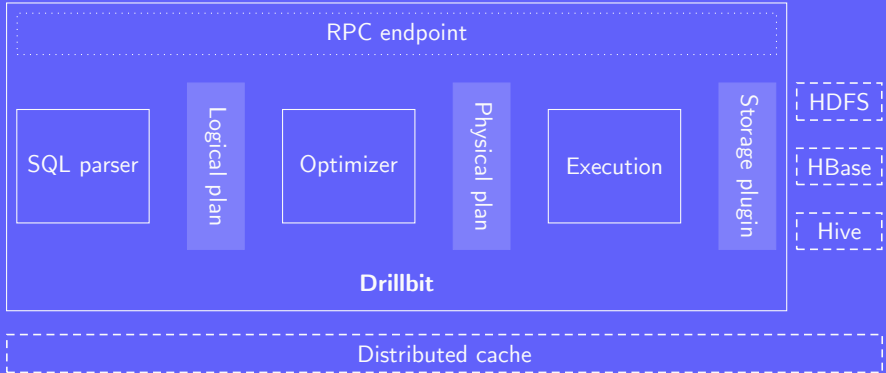
Running Drill as a YARN application:

- 1 Start Drill on the client machine
- 2 Upload resources to the FS and request resources for the application master
- 3 Ask a node manager to prepare and start a container for the application master
- 4 The application master contacts the resource manager to obtain more containers

Running Drill as a YARN application:

- ⑤ Request the start of Drill software on each assigned node
- ⑥ Start a “Drill process” called a *drillbit*
- ⑦ Each drillbit starts and registers with Zookeeper
- ⑧ The application master checks the health of each drillbit through Zookeeper
- ⑨ Use Zookeeper to retrieve information on the drillbits, run queries, etc.





Foreman drillbit:

- Drillbit that receives the query
- It drives the entire query

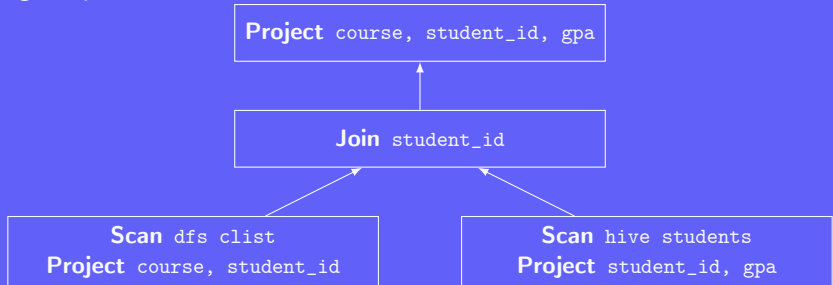
Initial SQL query:

```
1 SELECT course, student_id, gpa
2 FROM clist.json l, hive.students s
3 WHERE l.student_id = s.student_id
```

Initial SQL query:

```
1 SELECT course, student_id, gpa
2 FROM clist.json l, hive.students s
3 WHERE l.student_id = s.student_id
```

Logical plan:



Major fragment:

- Concept representing a phase of the query execution
- Composed of minor fragments

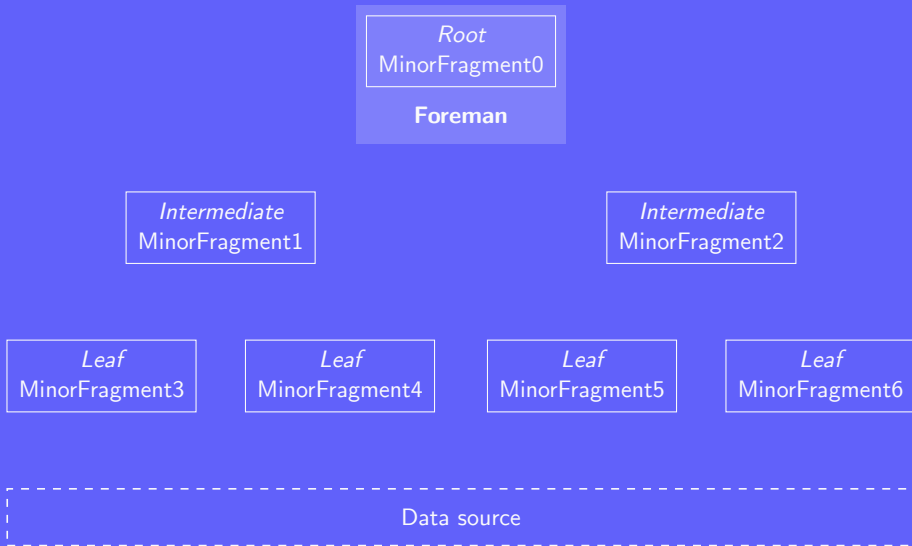
Major fragment:

- Concept representing a phase of the query execution
- Composed of minor fragments

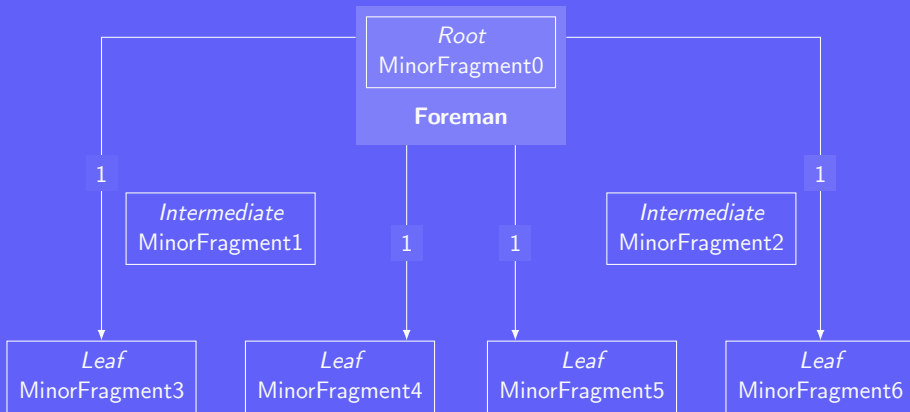
Minor fragment:

- Logical unit of work running in a thread
- Contain one or more relational operators
- Usually as numerous as the number of available drillbits
- Scheduled based on data locality when possible and round-robin otherwise

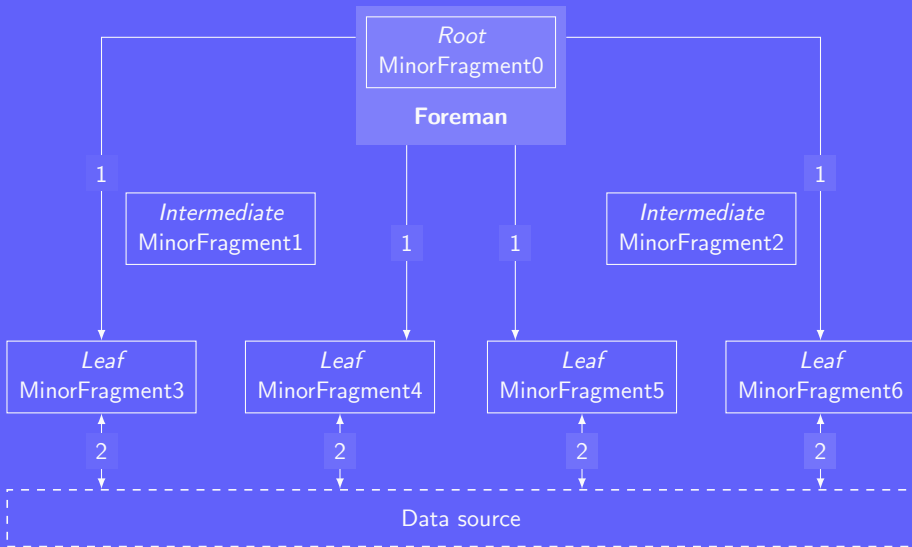
Minor fragments execution tree



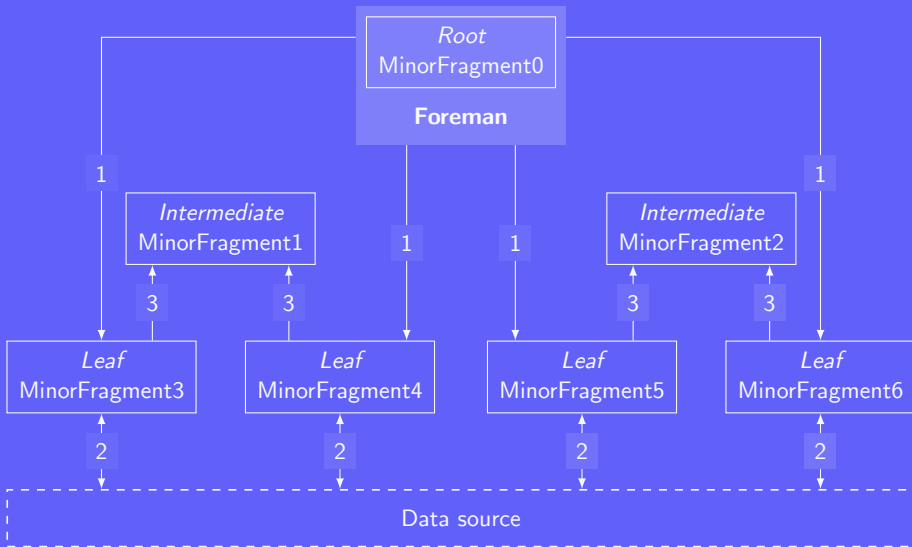
Minor fragments execution tree



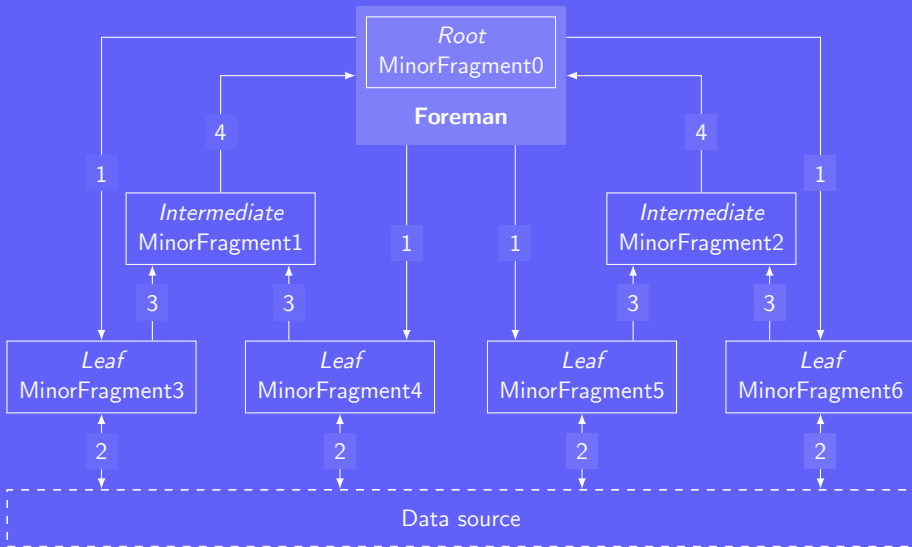
Minor fragments execution tree



Minor fragments execution tree



Minor fragments execution tree



Architecture:

- No central server, no master-slave concept
- Each drillbit contains all the services and capabilities of Drill
- Nodes can be added or removed at no cost

Columnar execution:

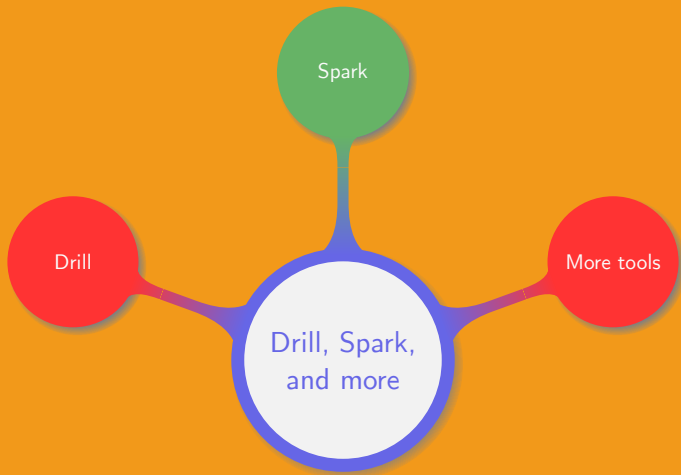
- Avoid access for columns not involved in the query
- Directly performs SQL processing on columns

Optimistic query execution:

- Assume no failure will occur during query execution
- Rerun the query in case of failure
- Only write on disk when memory overflows

Vectorization: allow the CPU to operate on vectors

Runtime compilation: generate efficient code for each query



Spark organisation:

- Application: user program built on Spark and composed of:
 - A driver program: process running the main function
 - Executors: processes launched for an application on worker nodes
- The driver program connects to a cluster manager
 - Standalone: cluster manager provided with Spark
 - YARN
 - Mesos
 - Kubernetes

Two modes available:

- Client mode:
 - Driver runs in the client
 - Required in the case of interactive programs
 - Useful when building a Spark program
- Cluster mode:
 - The entire application runs in the cluster
 - Appropriate from production jobs
 - YARN application master failure strategy (slide 10.31) is applied

Spark job workflow:

- 1 Start the driver program on a client node
- 2 The driver requests a container to the resource manager
- 3 A container starts and runs an Executor Launcher application master
- 4 The Executor Launcher requests more resources to start Executor backends processes in new containers
- 5 Each Executor Backend registers with the driver

Workflow similar to client mode but:

- The driver program runs in a YARN application master process
- The client submit a job but does not run any user code
- The application master starts the driver program
- The driver program “replaces” the Executor Launcher

Workflow similar to client mode but:

- The driver program runs in a YARN application master process
- The client submit a job but does not run any user code
- The application master starts the driver program
- The driver program “replaces” the Executor Launcher

Remark.

Data locality:

- Executors are launched before data locality information becomes available
- The driver can optionally specify preferred locations

Resilient Distributed Dataset (RDD):

- Core abstraction in Spark
- Collection of objects distributed across a cluster
 - Read-only: do not alter a dataset, transform it into a new one
 - Resilient: no disk write, reconstruct the RDD in case of partition loss

Resilient Distributed Dataset (RDD):

- Core abstraction in Spark
- Collection of objects distributed across a cluster
 - Read-only: do not alter a dataset, transform it into a new one
 - Resilient: no disk write, reconstruct the RDD in case of partition loss
- Loaded as input
 - Created from an external dataset
 - From an existing RDD
 - Parallelising an existing collection

Two types of operations on an RDD:

- Transformation:
 - Create a new dataset from an existing one
 - Only compute the result when an action is run
 - Do not return any result to the driver program
- Action:
 - Run a computation on a dataset
 - Return the value to the driver program

Two types of operations on an RDD:

- Transformation:
 - Create a new dataset from an existing one
 - Only compute the result when an action is run
 - Do not return any result to the driver program
- Action:
 - Run a computation on a dataset
 - Return the value to the driver program

Benefits of this approach:

- Transformed RDD is in memory when performing an action
- No large dataset to send back to the driver program

Datasets are cached in memory across operations:

- An RDD is stored on the node where it was computed
- An old RDD is dropped following the LRU algorithm
- A lost RDD is automatically recomputed if needed

Datasets are cached in memory across operations:

- An RDD is stored on the node where it was computed
- An old RDD is dropped following the LRU algorithm
- A lost RDD is automatically recomputed if needed

Caching levels:

- Memory only: no compression, lost partitions are recomputed
- Memory and disk: partitions that do not fit in the memory are spilled on disk
- Memory only serialized: compression enabled
- Replication: all the above but also replicate on another node

Serialization of data and functions:

- Used to share information among the executors
- Transparent to the user

Serialization of data and functions:

- Used to share information among the executors
- Transparent to the user

Task closure:

- Cannot share variables among executors
- Determine what variables and methods an executor needs
- Serialize this closure and send it to the executor
- Each executor receives a copy of the original variable
- Variables are not updated on the driver

Broadcast variables:

- Read-only variables broadcasted to each executor
- Data sent in an efficient way to minimize traffic
- Useful for data needed over several stages of the computation

Broadcast variables:

- Read-only variables broadcasted to each executor
- Data sent in an efficient way to minimize traffic
- Useful for data needed over several stages of the computation

Accumulators:

- Variables that can be added to using associative and commutative operations
- The driver can retrieve their value
- They are only updated on action tasks
- Update only occurs once, even if an action is rerun

Job submission and execution:

- A job is submitted when an action is performed on an RDD
- The transformations on the RDD are organised into a logical execution plan
- Spark DAG scheduler transforms the logical plan into a execution physical plan
- The physical plan defines stages, split into tasks
- Spark task scheduler constructs a mapping of tasks to executors
- The executor runs the task
- Executors send status updates to the driver when a task is completed or has failed

Advantages of DAG:

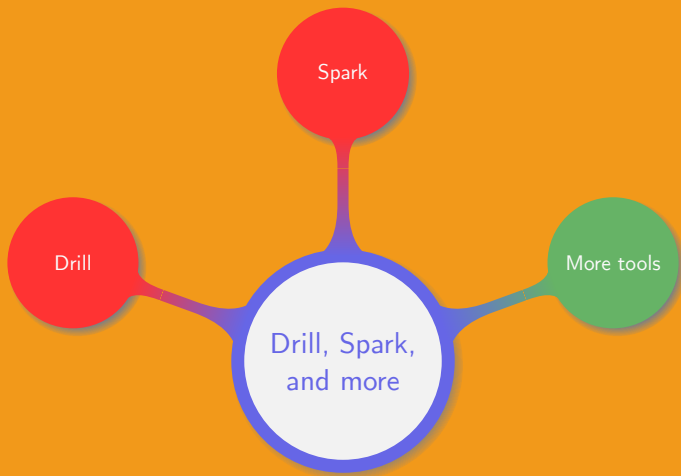
- Any lost RDD can easily be recovered
- Offers more possibilities than a simple Map and Reduce approach
- Transformation on RDDs are not directly applied:
 - Allows better optimizations
 - Decreases disk writes and data transfer

Advantages of DAG:

- Any lost RDD can easily be recovered
- Offers more possibilities than a simple Map and Reduce approach
- Transformation on RDDs are not directly applied:
 - Allows better optimizations
 - Decreases disk writes and data transfer

Remark.

Spark generalises the MapReduce approach, is much faster, and features many more high-level operators



Most common ways to integrate R with Hadoop:

- MapReduce:
 - Rhipe: Hadoop 2 supported, not recent update
 - RHadoop: Hadoop 2 supported, no recent update
 - Hadoop streaming: use MapReduce streaming engine
- Drill:
 - DrillR: R driver for Apache Drill
 - Sergeant: tools to query Drill data from R
- Spark:
 - SparkR: light-weight R frontend for Spark
 - SparklyR: R interface for Apache Spark

Copying data into Hadoop using CLI:

```
1  hadoop fs -put local-file1.txt hdfs-file.txt
2  hadoop fs -put local-file1.txt local-file2.txt \
3    hdfs-directory/
4  hadoop fs -touch hdfs-file.txt
```

Copying data into Hadoop using CLI:

```
1  hadoop fs -put local-file1.txt hdfs-file.txt
2  hadoop fs -put local-file1.txt local-file2.txt \
3    hdfs-directory/
4  hadoop fs -touch hdfs-file.txt
```

Other strategies:

- Use WebHDFS and its REST API
- Access HDFS as an NFS mount
- Use DistCp to copy data from a cluster into another one

Simple strategy:

- Write a Map in R to drive the Map part
- Write a Reduce in R to drive the reduce part

Map outputs each record on a separate line to feed Reduce

Simple strategy:

- Write a Map in R to drive the Map part
- Write a Reduce in R to drive the reduce part

Map outputs each record on a separate line to feed Reduce

```
1  hadoop jar hadoop-streaming.jar \  
2    -input InputDirs \  
3    -output OutputDir \  
4    -mapper Map.R \  
5    -reduce Reduce.R
```

Basics on Sergeant:

- Implements a Database Interface using Drill REST API
- Allows connection through JDBC (RJDBC library needed)
- Implements a `dplyr` interface
- Allows passing direct SQL requests
- Perfect if no complex data analysis is to be performed on Drill

sergeant.R

```
1 library(sergeant)
2 ds <- src_drill("localhost")
3 libgen_d<- tbl(ds,
4   "dfs.`/home/courses/ve572/datasets/libgen.csvh`")
5 start_time <- Sys.time()
6 a<-libgen_d %>% filter(str_detect("hadoop",tolower(title))) %>%
7   select(title,author,year,publisher,edition)
8 b<-libgen_d %>% select(title, author) %>%
9   filter(str_detect("hadoop",tolower(title))) %>%
10  group_by(author) %>% count(author) %>% arrange(author) %>%
11  print(n=200) %>% collect()
12
13 end_time <- Sys.time()
14 print(end_time - start_time)
15
16 # c<-libgen_d %>% select(title, author, pages) %>%
17 #   filter(str_detect("hadoop",tolower(title))) %>% group_by(author) %>%
18 #   summarize(pp=mean(pages,na.rm=TRUE)) %>% arrange(author) %>% print(n=200)
```

Basics on SparklyR:

- Allows advanced data analysis, including Machine Learning
- Provides a complete `dplyr` backend
- Enables SQL queries through SparkSQL
- Can be extended to use H2O and `sparkling` libraries
- Fully compatible with `tidyverse`

sparklyr.R

```
1 library(sparklyr)
2 sc <- sparklyr::spark_connect(master = "local")
3 libgen_s <- sparklyr::spark_read_csv(sc, "ms",
4   "/home/courses/ve572/datasets/libgen.csvh")
5 start_time <- Sys.time()
6 a<-libgen_s %>% filter(str_detect("hadoop",tolower(title))) %>%
7   select(title,author,year,publisher,edition) %>% collect()
8 end_time <- Sys.time()
9 print(end_time - start_time)
10
11 b<-libgen_s %>% select(title, author, pages) %>%
12   filter(str_detect("hadoop",tolower(title))) %>%
13   group_by(author) %>% summarize(pp=mean(pages,na.rm=TRUE)) %>%
14   arrange(author) %>% print(n=200)
15 spark_disconnect(sc)
```


Linux Containers (LXC):

- Operating-system-level virtualization method
- Relies on the kernel's *cgroup* and *namespace* isolation functionalities
- Concurrently run multiple isolated Linux OS on a machine
- Each container must be individually maintained
- Containers can be either privileged or unprivileged

Linux Containers (LXC):

- Operating-system-level virtualization method
- Relies on the kernel's *cgroup* and *namespace* isolation functionalities
- Concurrently run multiple isolated Linux OS on a machine
- Each container must be individually maintained
- Containers can be either privileged or unprivileged

Limitation: necessitate the setup of a whole OS for each container

Docker uses a different approach than LXC:

- Initially based on LXC but *now* completely independent
- A daemon manages the docker containers
- A container is an encapsulated environment that runs applications
- An image containing an application and its dependencies is built based on a “configuration file”
- To update simply replace the old image with a new one

Docker uses a different approach than LXC:

- Initially based on LXC but *now* completely independent
- A daemon manages the docker containers
- A container is an encapsulated environment that runs applications
- An image containing an application and its dependencies is built based on a “configuration file”
- To update simply replace the old image with a new one

Limitation: need to “manually” deploy, manage, scale containers

Kubernetes is a container orchestration tool:

- Mostly used with, but not limited to Docker
- Initially developed as an internal Google project
- Kubernetes is a cluster application that handles a cluster:
 - *Master*: controls all other machines in the cluster
 - *Nodes*: the machines onto which applications are running
 - *Pods*: single instance of an application or running process

Main tasks of Kubernetes:

- Monitor the health of the running applications
- Balances the load
- Manages hardware resources allocation
- Eases the deployment of preconfigured applications
- Allows access to storage in the same way as any other resources

During the lifespan of an application:

- Containers can live, die, be resurrected
- Kubernetes handles everything without any human interaction

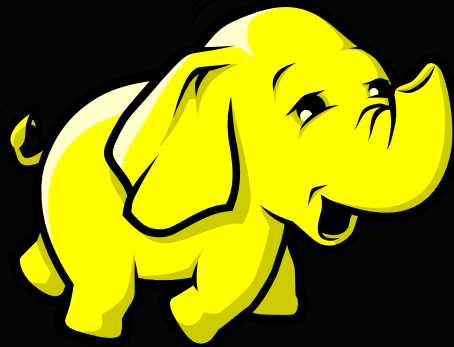
Main tasks of Kubernetes:

- Monitor the health of the running applications
- Balances the load
- Manages hardware resources allocation
- Eases the deployment of preconfigured applications
- Allows access to storage in the same way as any other resources

During the lifespan of an application:

- Containers can live, die, be resurrected
- Kubernetes handles everything without any human interaction

Kubernetes can be coupled to Hadoop or work independently



Thank you!