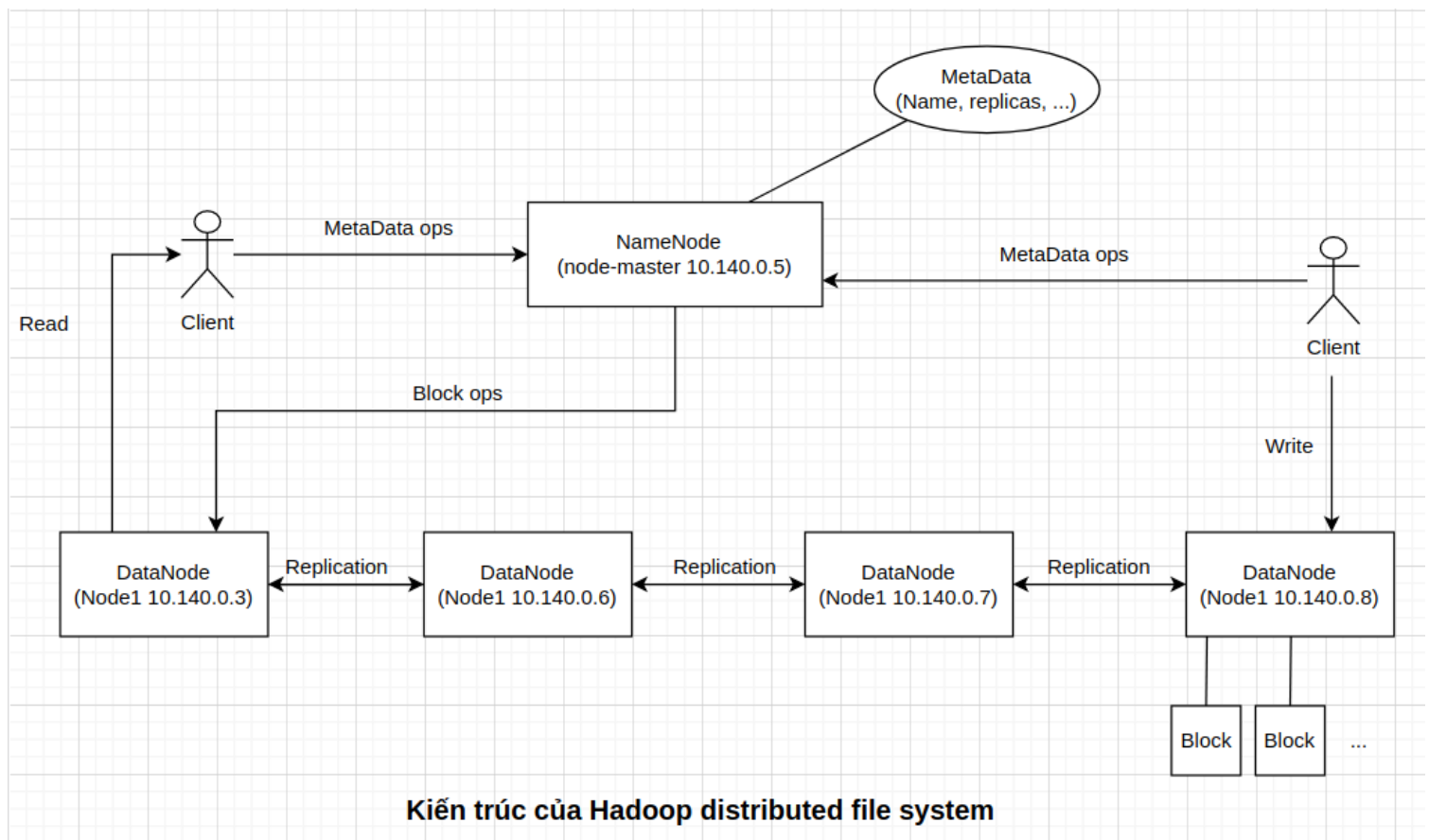


Bài tập tuần 4 - Hadoop - Chương trình masterDev mùa 3

Họ và tên: Phạm Văn Ngọc - ngocpv22

1. Kiến trúc Hadoop distributed file system (HDFS Architecture)



Mô tả các thành phần trong kiến trúc HDFS:

1.NameNode:

NameNode là nơi lưu trữ và cập nhật các meta data như địa chỉ của block đặt tại datanode, quyền truy cập của client tới các block, ... => Nó kiểm soát quyền truy cập từ phía client. NameNode có nhiệm vụ chính là quản lí và duy trì các dataNode, kiểm soát thông tin tình trạng của các dataNode, luôn lắng nghe để đảm bảo các dataNode còn sống. Nó điều khiển sự cân bằng về lưu lượng lưu trữ và lượng truy cập tới các dataNode. Nếu NameNode chết thì cả cụm hadoop sẽ chết theo.

2.DataNode

DataNode có vai trò:

Là nơi lưu trữ các dữ liệu khi dữ liệu được đưa vào HDFS

Là nơi chạy các tiến trình xử lý dữ liệu (Khi các job được chia cho các dataNode)

Gửi tình trạng health về cho nameNode thường xuyên (khoảng 3s một lần để NameNode biết tình trạng của nó).

Trong quá trình write từ phía client, DataStream sẽ chuyển packet vào dataNode phù hợp nhất, sau đó dataNode này sẽ chuyển tiếp đến các dataNode khác (Quá trình replication)

3.Block

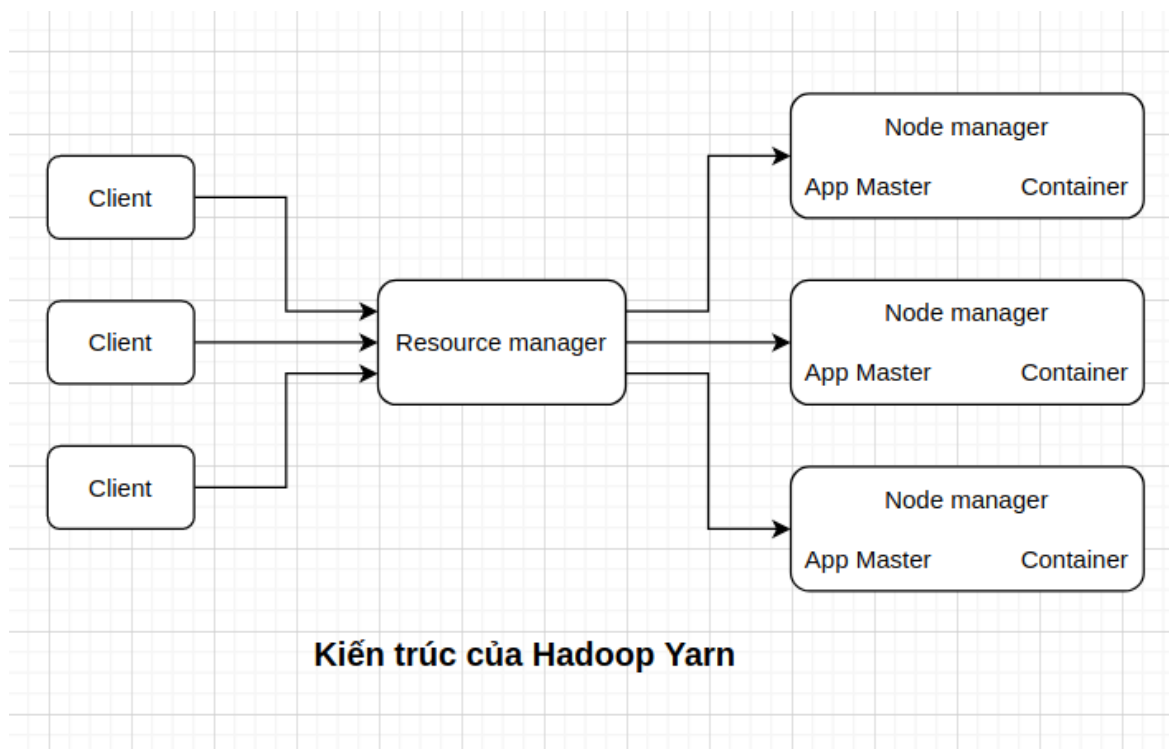
Dữ liệu trong HDFS sẽ được chia nhỏ thành nhiều mảnh, mỗi mảnh như vậy được gọi là block, các block này nằm trong các dataNode. Thông tin về dữ liệu được chia thành bao nhiêu block, block nằm ở dataNode nào sẽ được lưu trong metaData của nameNode.

Các block được nhân bản (quá trình replication) và lưu tại các dataNode, để đảm bảo không mất mát dữ liệu khi một dataNode nào đó gặp sự cố hoặc chết, đây là một trong những đặc trưng của HDFS về độ tin cậy và khả năng phục hồi sau lỗi tốt.

2. Kiến trúc Yarn

Yarn (“Yet-Another-Resource-Negotiator”) là framework hỗ trợ phát triển ứng dụng phân tán, được sử dụng làm hệ thống quản lý resource trong Hadoop.

Kiến trúc của Yarn để quản lý tài nguyên và quản lý job thành chia thành 2 components riêng biệt: Resource manager và Application Master.



Mô tả các thành phần trong kiến trúc của Yarn:

Resource manager: Quản lý toàn bộ tài nguyên tính toán của toàn bộ cụm hadoop.

Node manager: Mỗi Node manager quản lý các job trên mỗi node. Các node tính toán trong cluster bây giờ sẽ chạy NodeManager để quản lý các tiến trình chạy trên máy đó.

Application Master: Quản lý vòng đời của job trên các DataNode. Nó hoạt động cùng với Node Manager và giám sát việc thực thi các task.

Container: Quản lý tài nguyên vật lý cho các node như RAM, CPU, ...

Client: Đóng vai trò đưa các job map-reduce cho hệ thống xử lý.

3. Tóm tắt các lệnh cơ bản thao tác với HDFS

*Thao tác với hệ thống file HDFS:

Xem các mục có trong path:	<code>hdfs dfs -ls +path :</code>
Tạo thư mục:	<code>hdfs dfs -mkdir /user/ngocpv22</code>
Tạo file:	<code>hdfs dfs -touch /user/ngocpv22/myFile</code>
Tạo file có nội dung:	<code>hdfs dfs -echo "abc" >myFile</code>
Xóa thư mục:	<code>hdfs dfs -rm -r -f /user/ngocpv22</code>
Trao quyền:	<code>hdfs dfs -chown -R ngocpv22 /user/ngocpv22</code>
Đẩy 1 file từ ngoài vào hdfs:	<code>hdfs dfs -put myFile2 /user/ngocpv22</code>

Các lệnh sử dụng với quyền admin (super user)

Format lại hệ thống file:	<code>hadoop namenode -format</code>
Khởi động/Dừng hệ thống hdfs:	<code>sbin/start-all.sh sbin/stop-all.sh</code>

Kiểm tra tình trạng hệ thống hdfs:	<code>hadoop fsck /</code>
------------------------------------	----------------------------

Báo cáo và thống kê chi tiết tình trạng HDFS: `hadoop dfsadmin -report`