

Hadoop and

Erika and Nathan



History of Hadoop



- Started with Google
- Google utilized a DFS and MapReduce
- The Hadoop people basically stole it

Hadoop Distributed File System

- Hadoop is used for **BIG** DATA
- Data is stored across many nodes (cluster)
- Hadoop makes 2 back ups of your data
- Namenode calibrates the data across the cluster

HDFS - A Farmer's Story

His "file"



Basket ID	Contents
0	apple, apple, pear, pear, orange, orange, apple, pear
1	orange, pear, pear, apple, orange, apple, apple, orange
2	apple, pear, orange, apple, apple
3	pear, pear, apple, apple, pear, orange, apple



64 MB



64 MB

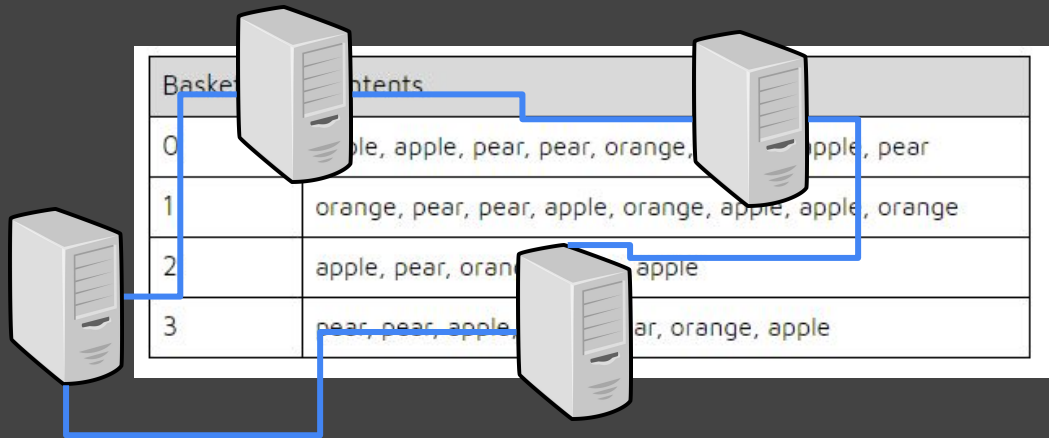


64 MB

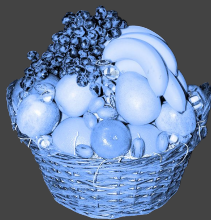


64 MB

Hadoop Distributed File System



64 MB



64 MB

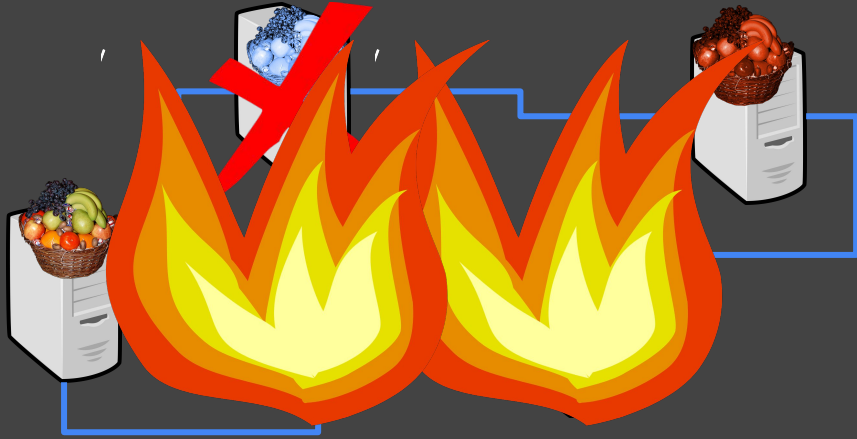


64 MB

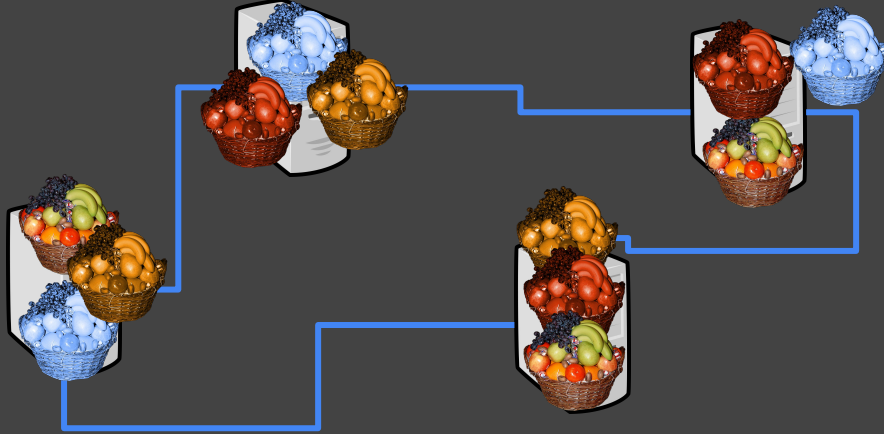


64 MB

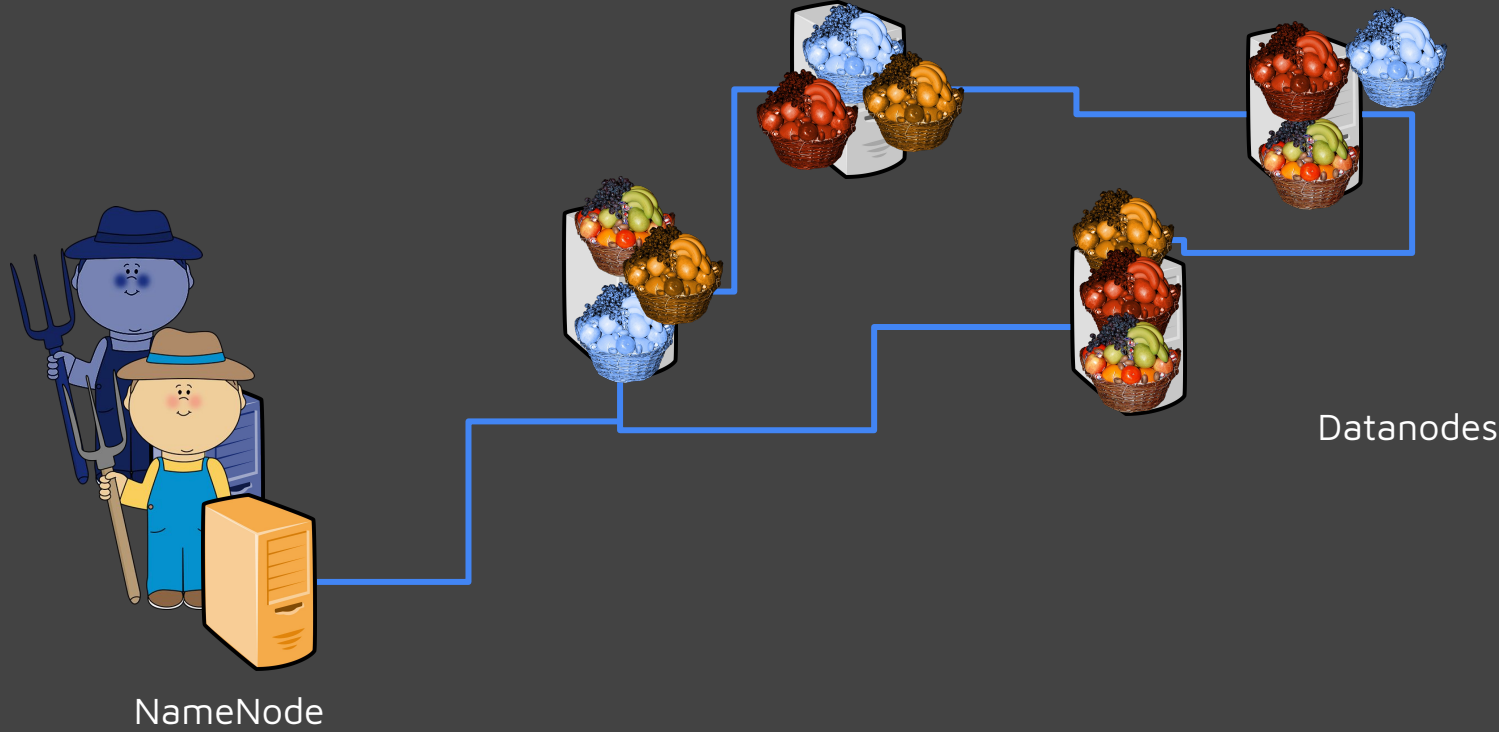
Hadoop Distributed File System



Hadoop Distributed File System



Hadoop Distributed File System



MapReduce

Map

Sort/Shuffle

Reduce

- Deals with big data more efficiently
- Utilizes multiple nodes for **PARALLEL PROCESSES**
- Divides/conquers (backs up)

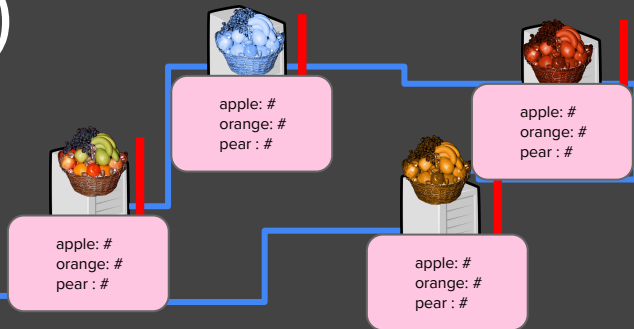
MapReduce

Map

Sort/Shuffle

Reduce

MAP()



Basket ID	Contents
0	apple, apple, pear, pear, orange, orange, apple, pear
1	orange, pear, pear, apple, orange, apple, apple, orange
2	apple, pear, orange, apple, apple
3	pear, pear, apple, apple, pear, orange, apple

On each node:



MAP()



apple: #
orange: #
pear : #

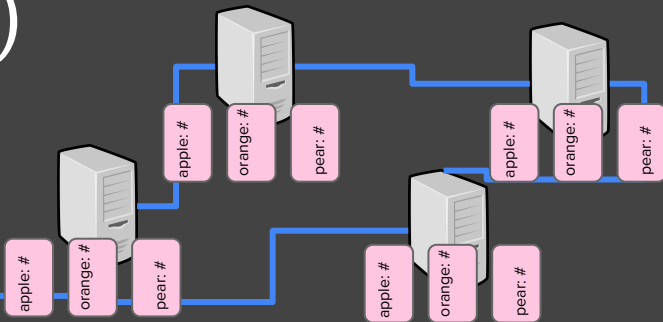
MapReduce

Map

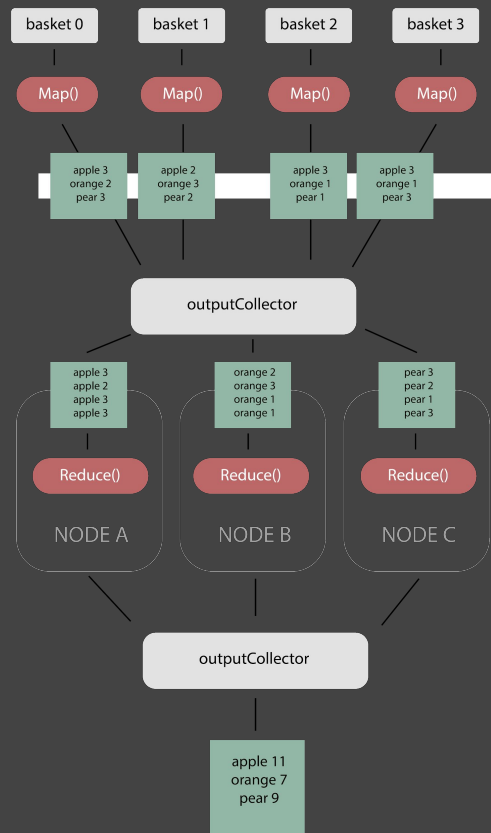
Sort/Shuffle

Reduce

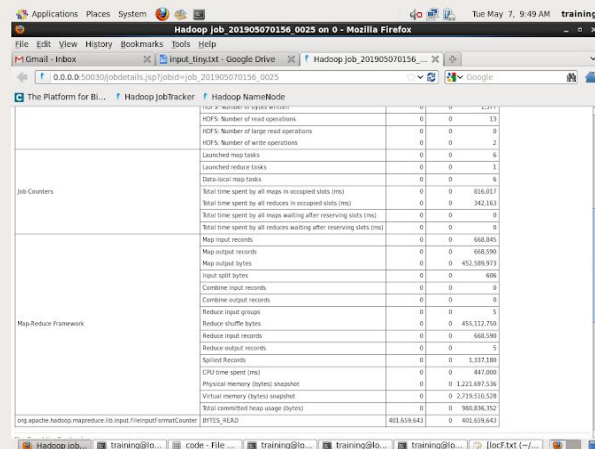
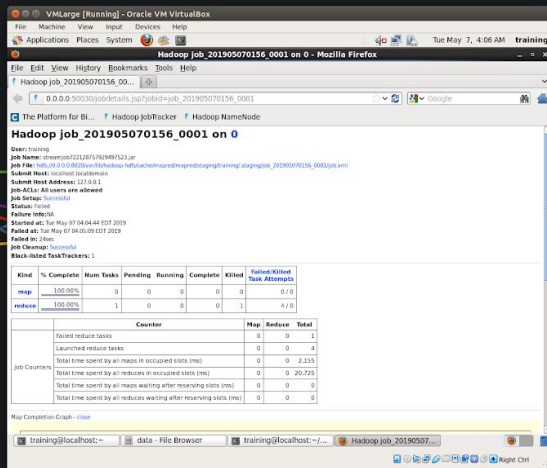
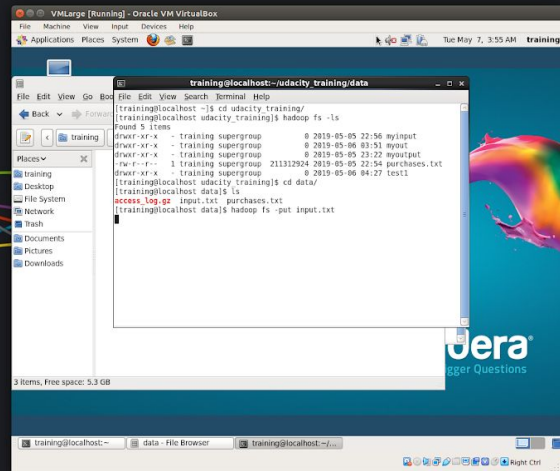
MAP()



Basket ID	Contents
0	apple, apple, pear, pear, orange, orange, apple, pear
1	orange, pear, pear, apple, orange, apple, apple, orange
2	apple, pear, orange, apple, apple
3	pear, pear, apple, apple, pear, orange, apple



Demo:



Class Interaction Yay

Guess: 1 Star 3 Stars 5 Stars



Class Interaction Yay

Guess: 1 Star 3 Stars 5 Stars



QUESTIONS?