# 2016 International Conference on Platform Technology and Service (PlatCon)

## Proceedings

15-17 February 2016
Jeju, Korea

# Twister Platform for MapReduce Applications on a Docker Container

Yunhee Kang
Division of Information and Communication
Baekseok University
115 Anseo Dong, Cheonan, 31065, Korea
yunh.kang@gmail.com

R. Young Chul Kim
SELab., Dept. of Computer Information Communications
Hongik University
2639, Sejong Campus, 30016, Korea
bob@hongik.ac.kr

*Abstract*—Docker is one of ways to provide more light-weight for agile computing resource based on container technique to handle this problem. For this work we have chosen this specific tool due to the increasing popularity of MapReduce and cloud container technologies such as Docker. This paper aims at automatically configuring Twister workloads for container-driven clouds. Basically this is the first attempt towards automatic configuration of Twister jobs on container-based cloud platform for many workloads.

*Keywords—Docker; agile computing; MapReduce; Twister*

## I. INTRODUCTION

When prototyping a distributed application like MapReduce, a developer needs both to ensure the application execution corresponds to the specification and that its performance is not impacted by the number of nodes or by some failure scenarios [1-3]. Indeed, MapReduce relies on successive computing-commination steps that, if not coordinated with care, lead to performance bottlenecks and a poor scalability. However, this technology is both grounded on extremely complex platforms that are often difficult to understand, configure and optimize details.

As in rare situations both objectives can be reached at once, it is usual to start prototyping in a small set of nodes and then, when the execution was proved, perform additional scalability and fault tolerance tests.

Because all these steps require the execution of well-defined scenarios, we decided to rely on container-based virtualization, which allow the researchers to control both system images and network interconnections. The configuration can have a profound impact on application performance, security as well as availability, and the values they provide to end-users. Clearly, making these configurations manually is cumbersome or even impractical.

For this work we have chosen this specific tool due to the increasing popularity of MapReduce and cloud container technologies such as Docker. Basically this is the first attempt towards automatic configuration of Twister jobs on container-based cloud platform VM for many workloads. This paper aims at automatically configuring Twister workloads for container-driven clouds. Basically this is the first attempt towards automatic configuration of Twister jobs on container-based cloud platform VM for many workloads.

## II. RELATED WORKS

### A. Twister

Twister is one of MapReduce implementations, which is an enhanced MapReduce runtime with an extended programming model that supports an iterative MapReduce computing efficiently [4-5]. In addition it provides programming extensions to MapReduce with broadcast and scatter type for transferring data. These improvements allow Twister to support iterative MapReduce computations highly efficiently compared to other MapReduce runtimes. The demanding requirements have led to the development of a new programming model like Twister based MapReduce. It reads data from local disks of the worker nodes and handles the intermediate data in distributed memory of the worker nodes.
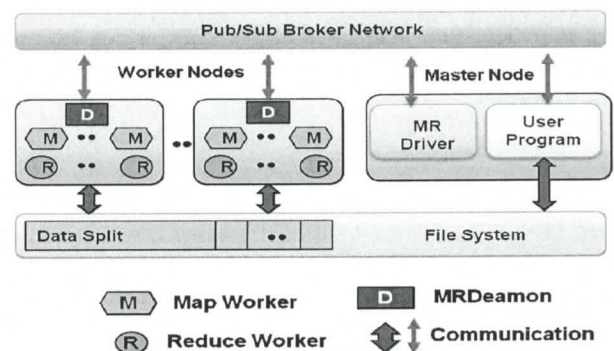


Fig. 1. Overall architecture of Twister

As shown in Figure 1, all communication and data transfers are performed via a pub/sub broker network via NaradaBrokering that is an open-source, distributed messaging infrastructure [9]. Twister uses a pub/sub messaging infrastructure to handle four types of communication needs; (i) sending/receiving control events, (ii) sending data from the client side driver to the Twister daemons, (iii) transferring intermediate data between map and reduce tasks, and (iv)

sending the outputs of the reduce tasks back to the client side driver to invoke the combine operation.

## B. Docker

Virtualization provides a way to abstract the hardware and system resources from an operation system, which is used to reduce the actual number of physical servers and to improve scalability and workloads in the cloud environment. In cloud computing environment, a VM is a computing platform that creates a virtualized layer between the computing hardware and the application.

Docker is an open platform released as a container project, originated by dotCloud, for developers and system administrators to build, ship and run distributed applications. It automates the deployment of applications inside software containers, by providing an additional layer of abstraction and automation of operating-system-level virtualization on Linux and enables applications to be quickly assembled from components and eliminated the friction. Especially virtualization technology is the key technology, though it, deployment of cloud computing system can be implemented.

## III. Experimental Results

### A. K-means clustering

K-means clustering is a method of cluster analysis, which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. Given a data set, D, of n objects, and k, the number of clusters to form, a partitioning algorithm organizes the objects into k partitions (k <= n), where each partition represents a cluster. The error function used is the sum of the distance that each point is from its cluster's centroid. In the k-means MapReduce application we built, it is used to determine the number of iterations in the data set. We describe an overall framework of a MapReduce based k-means with an initial starting configuration of this k-means clustering as input parameters. Figure 2 shows the designed framework.
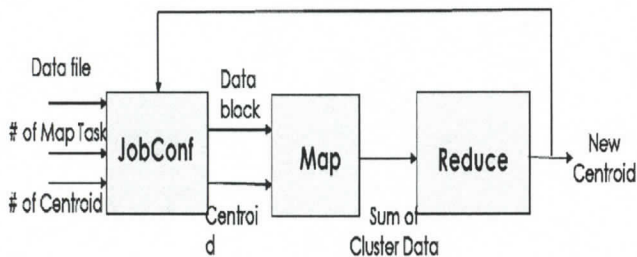


Fig. 2. Designed framework of a MapReduce based k-means

### B. Experiment Environment

The following shows the two Docker containers generated from the image ubuntu:twister. Each of these containers has a network interface with IP addresses 172.17.0.9 and 172.17.0.10 respectively.



The following shows the result of k-means MapReduce application. To evaluate performance of a k-means MapReduce application, we define a configuration of this experiment. This configuration is used to cluster 3,000 data points with 3 mapper tasks and 1 reducer task. The elapsed time is 15.646 second with 87 iterations. In the perspective of elapsed time, there is no difference between physical machine and virtual machine.



## IV. Conclusions

This paper described the overall process of building the experimental environment for Twister applications. For this work we have chosen this specific tool due to the increasing popularity of MapReduce and cloud container technologies such as Docker. Basically this is the first attempt towards automatic configuration of Twister jobs on container-based cloud platform VM for many workloads. This paper aims at automatically configuring Twister workloads for container-driven clouds. Basically this is the first attempt towards automatic configuration of Twister jobs on container-based cloud platform VM for many workloads. In the perspective of elapsed time, there is no difference between physical machine and virtual machine.

REFERENCES

[1] Dean, J., Ghemawat, S.: MapReduce: A Flexible Data Processing Tool. CACM 53, 72-77 (2010)

[2] Morton, K., Friesen, A., Balazinska, M., Grossman, D.: Estimating the Progress of MapReduce Pipelines. IEEE 26th International Conference on Data Engineering (ICDE), 2010, pp. 681 - 684, Long Beach, CA (2010)

[3] Dean, J., Ghemawat, S.: Mapreduce: Simplified data processing on large clusters. CACM 51, 107-113 (2008)

[4] Ekanayake, J., Li, H., Zhang, B., Gunarathne, T., Bae, S.-H., Qiu, J., Fox, G.: Twister: A Runtime for Iterative MapReduce. The First International Workshop on MapReduce and its Applications (MAPREDUCE'10) - HPDC2010, (2010)

[5] Geoffrey Fox and Shrideep Pallickara. Deploying the NaradaBrokering Substrate in Aiding Efficient Web & Grid Service Interactions. Invited paper for Special Issue of the Proceedings of the IEEE on Grid Computing. Vol 93, No 3, 564-577( 2005).