Scientific
Research
Publishing

# Unified Platform for AI and Big Data Analytics

**Sik Kim, Yongjin Kwon***

Department of Industrial Engineering, College of Engineering, Ajou University, Suwon, South Korea
Email: *yk73@ajou.ac.kr

## Abstract

This paper describes an integrated platform for machine learning and big data analysis. The integrated platform is configured in a way that builds a large distributed data processing environment in the computing environment that makes up the NVIDIA AI platform. In addition, this paper describes the background of this idea selection and the use of the software to build the unified platform. The technical details are shown in terms of how to create the proposed platform. In the anlaysis section, the methodology is provided and also the steps are explained as to how to use this integration platform. Finally, the expected effects are elaborated in the conclusion section.

## Keywords

Integrated Platform, Hadoop Eco System, Ambari, Virtual OS, Jetson TX-1, Dev Box, SSH

## 1. Introduction

In general, artificial intelligence modeling requires a high-end computing environment. In particular, the modeling of AI, this is based on graphic processing capabilities such as NVIDIA, and requires the combination of high-end GPUs as well as CPUs [1] [2]. The inefficiency exists, however, in this high-end computing environment, if the computing power is only used for machine learning purposes. The valuable computing power can be better utilized, if a virtual server is made within the computer and used for the analysis of big data. In recent years, we have witnessed the massive increase of data streams generated by the unmanned systems, such as drones and autonomous vehicles. Those systems are increasingly integrated with machine learning algorithms, while generating and transmitting a large amount of data (*i.e.* image data, system parameter data, text data, and so on) in real-time. The onboard computers are conducting the processing for machine learning algorithms. However, the big data streams generated by the system itself also needs to be processed and analyzed simulta-

neously. In this regards, an integrated platform is proposed in this study that can efficiently and simultaneously perform the big data analysis as well as the machine learning processing (this function is our research purpose). This is achieved by creating distributed computing environment with the use of Hadoop EcoSystem [3] [4]. The details are explained in the following sections.

## 2. Idea Extraction Process

### 2.1. Idea Generation

When data storing of information gathered from drones or autonomous driving cars is made, the artificial intelligence modeling (that is, machine learning algorithms) and the big data processing have been performed on different platforms. This concept is illustrated in Figure 1.

This process depicts the inefficiency because the processing of big data is performed on a different computing platform. On the other hand, Figure 2 is depicting the unified platform that runs both machine learning algorithms and big data processing within a single PC.

### 2.2. Why Should Build a Distributed Computing Environment?

Figure 3 illustrates the performance between RDBMS and Distributed Computing Environment (Node 3 - Node 5). As one can see, the distributed computing environment is faster than RDBMS. Also, if there are more nodes, the data processing time becomes shorter.

Figure 4 illustrates the comparison of performance when using the Hadoop (distributed computing environment) and the SPARK, a more advanced technology that is called "in-memory" system.

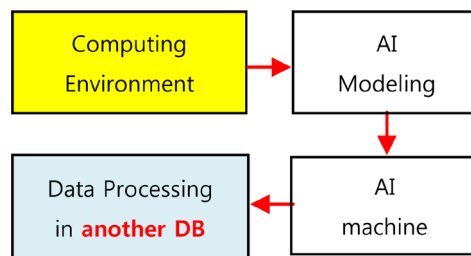As one can see from the figure, the Hadoop's Map Reduce performs better

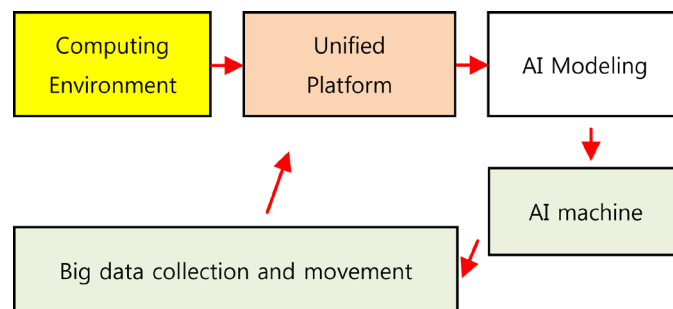Figure 1. Previous AI modeling and big data processing.
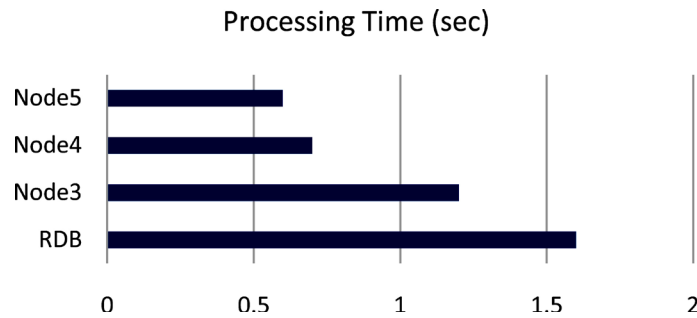
Figure 2. Unified platform model.

## Processing Time (sec)



**Figure 3.** Comparison between distributed computing environment and existing RDBMS.

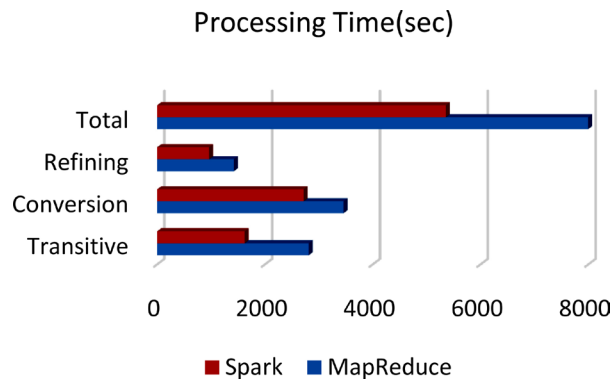## Processing Time(sec)



■ Spark ■ MapReduce

**Figure 4.** Comparison between SPARK and Hadoop (Map Reduce).

than RDBMS. The SPARK also performs better than the Hadoop's Map Reduce. This comparison shows the importance of using a distributed computing environment to handle big data.

## 3. NVIDIA Artificial Intelligence (AI) Platform

### 3.1. NVIDIA AI Platform—Host Server (Physical HW)

The host server is made up of HW equipment for AI machine learning such as NVIDIA's DevBox or Jetson TX1. These devices use the Linux Ubuntu 14.04 as the OS. The virtual server for the distributed environment also uses the same OS to enhance the compatibility. Overall, the distributed environment is made by using the host server and virtual server (built by using Oracle VirtualBox) [5] [6].

### 3.2. NVIDIA Digits (Image Training SW)

NVIDIA's representative image training software, Digits, is the SW that can be integrated with the autonomous vehicles and drones. It is basically SW that supports CUDA development environment developed by NVIDIA and is optimized for image training. The details are shown in **Figures 5-7** [7] [8] [9].

## 4. Configuration of Network between Host and Slave Servers

After building the slave servers, we used the SSH network configuration to make it possible to access the slave server from the host server without sharing the information through the network configuration between the servers.
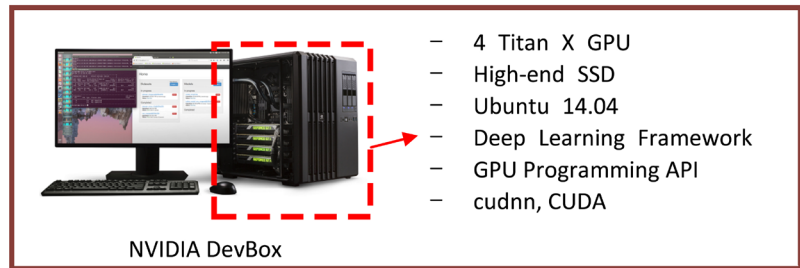
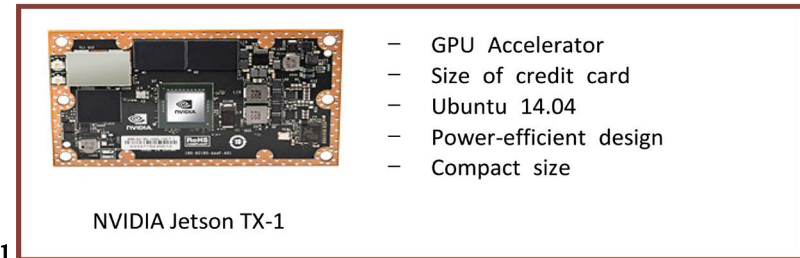**Figure 5.** NVIDIAI AI modeling "DevBox" (HW).



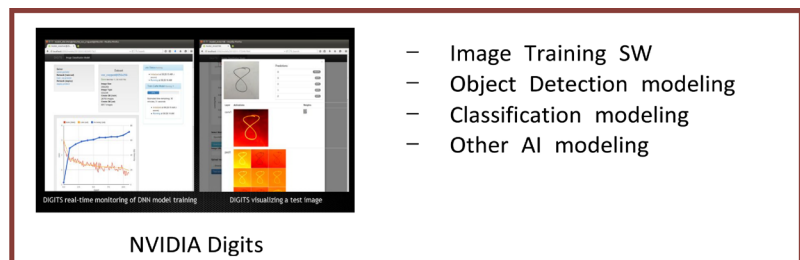**Figure 6.** NVIDIAI AI modeling machine "Jetson TX-1" (HW).



**Figure 7.** NVIDIAI AI modeling SW "Digits".

## 4.1. SSH Key Generation and Share

SSH network configuration allows access from a host server to a virtual server (slave). However, a password is required to access another server in its current state. Therefore, it is necessary to create an SSH shared key and share it between the servers so that communication can be smoothly performed without a password. In other words, setting up SSH network allows the autonomous communication between the servers without a password. The Linux command code is given in Table 1 and Table 2.

## 4.2. After Network Configuration between Host and Slave

Figure 8 shows that the SSH key is generated. The SSH key acts as the connection identity of the platform that needs to be connected. For example, one system is given a specific key structure and this key is only unique to this system. By examining the SSH key, one can identify the each individual system. Once the key is created and given to each system (*i.e.*, server) and the connection is established, the host server and the slave server can freely access each other's resources. Then, the computing of bid data can be performed on each other's platform. Figure 9 shows the remote connection between the host server and the

**Table 1.** Linux command code for SSH network configuration in host server.

| In Host Server | |
|---|---|
| generation of ssh key | root@client:~# ssh-keygen |
| confirmation of ssh key | root@client:~# ls − al ~/.ssh |
| firewall set-up | root@client:~# chmod 700 ~/.ssh<br>root@client:~# chmod 600 ~/.ssh/id_rsa<br>root@client:~# chmod 644 ~/.ssh/id_rsa.pub<br>root@client:~# chmod 644 ~/.ssh/authorized_keys<br>root@client:~# chmod 644 ~/.sshknown_hosts |
| copy ssh public key to Slave Server | root@client:~# scp~/.ssh/id_rsa.pub root@slave:id_rsa.pub |

**Table 2.** Linux command code for SSH network configuration in slave server.

| In Slave Server | |
|---|---|
| move ssh public key to ".ssh" dirictory | root@slave:~# cat id_rsa.pub >> ~/.ssh/authorized_keys |
| firewall set-up | root@slave:~# chmod 700 ~/.ssh<br>root@slave:~# chmod 644 ~/.ssh/authorized_keys |



**Figure 8.** Generations of SSH key.



**Figure 9.** Remote access to virtual server.

slave server. It is shown that the IP address is different.

## 5. Creation of Hadoop Cluster on NVIDIA AI Platform

This section explains the building of a Hadoop cluster (distributed environment) on the AI platform. The framework called "Apache Ambari" has been used to build a Hadoop cluster.

### 5.1. Creation of Hadoop Cluster

Asone can see in Figure 10, Apache Ambari UI can be easily installed in Linux environment. After that, the NVIDIAAI platform (Client.com) is connected and the slave server (Slave.com) to the Hadoop cluster using the Ambari framework is constructed. In this way, we have created an environment for analyzing not only machine learning algorithms but also big data on the same platform. The Ambari allows you to install SPARK-like SWs in Hadoop clusters. You can also install and uninstall the SW after the Hadoop cluster is completed.

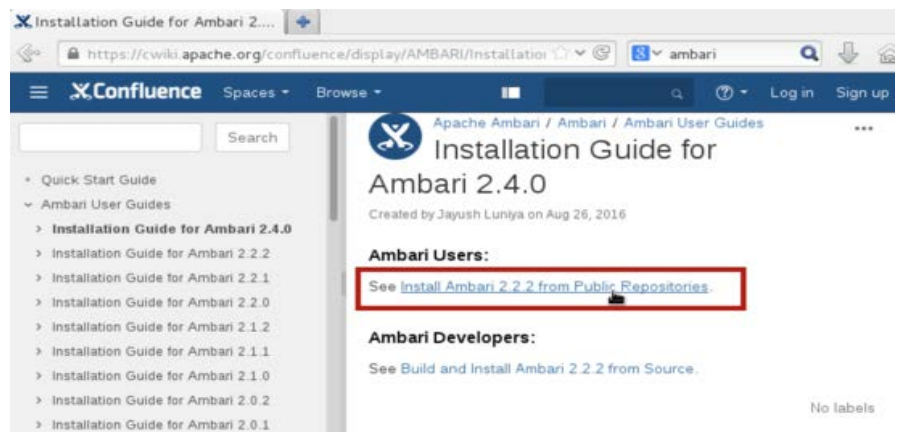Figure 11 shows the completion of Hadoop Cluster and the completed inte-



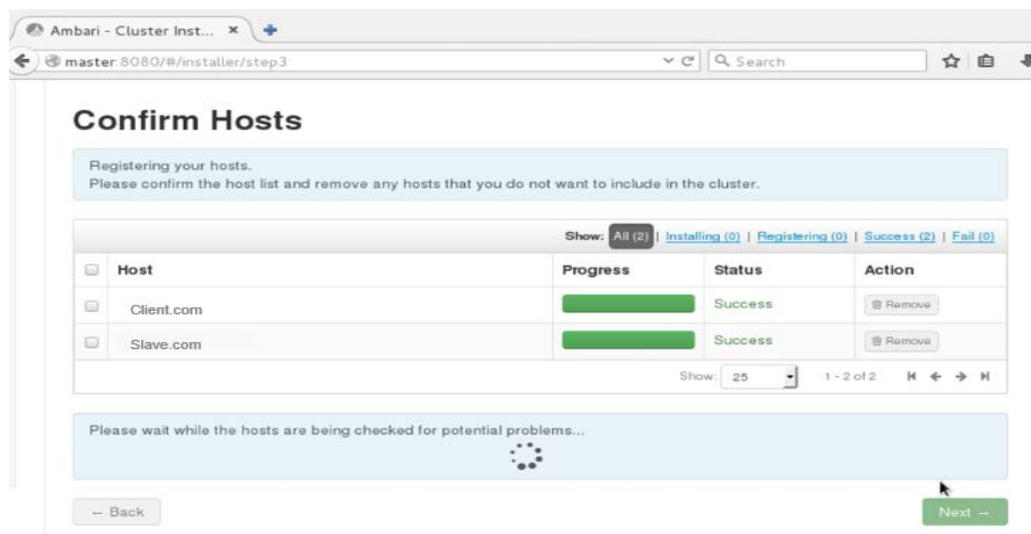**Figure 10.** Apache Ambari installation guide.



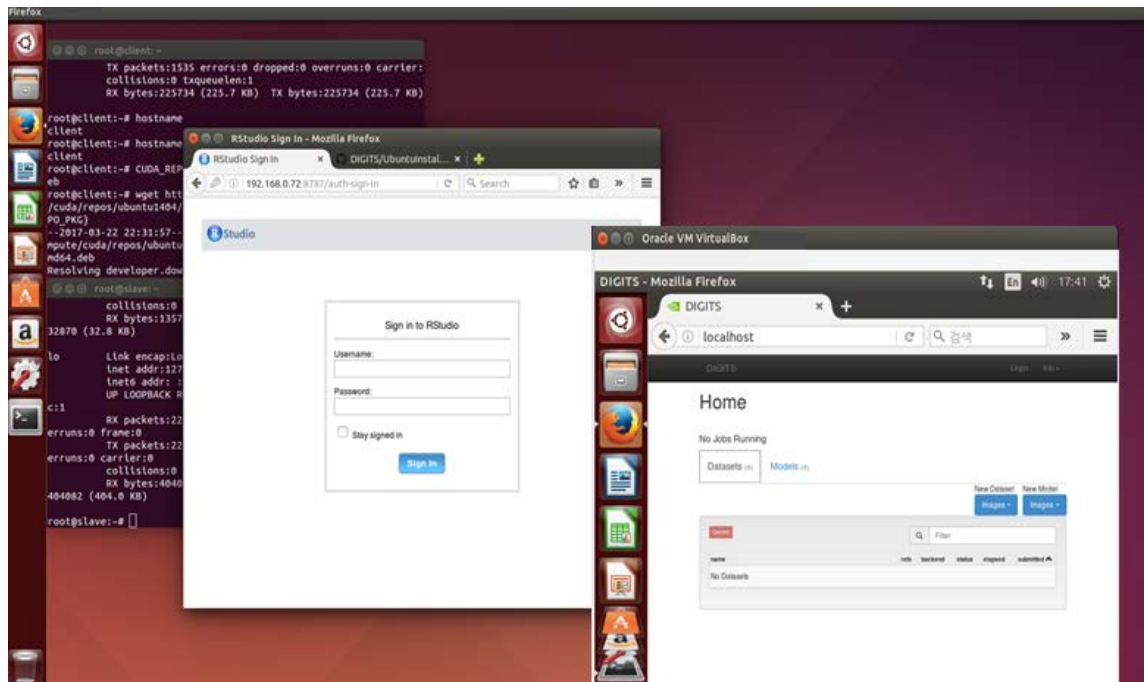**Figure 11.** Completion of making Hadoop cluster.

**Figure 12.** Unified platform for AI and big-data processing.

grated platform. The figure describes the completion of distributed computing environment by using Ambari UI. If one wants to use another SW (*i.e.*, SPARK and PIG) to process big data, one can also download them after the Hadoop Cluster is made.

## 6. Conclusions

In Figure 12, one can see the finished unified platform for AI and big data analytics. The UI represents the Digits SW accessed from the slave server (NVIDIA AI platform). The Digits SW is mainly for image training and machine learning purposes. The UI also represents the R-studio connected to the host server. The R-studio mainly processes the statistical analysis of big data. By creating this platform, the performance of the computing speed and the processing time can be significantly improved, as opposed to the conventional system that has been explained in Figure 1.

Therefore, in this integrated platform, it is possible to process the big data as well as the artificial intelligence algorithms using the same GPU accelerator. The development of this platform maximizes the utilization of the AI platform. Then a high-performance computing environment will improve efficiency. Therefore, it is not necessary to add additional computers for the big data analysis for the information gathering devices, such as drones and autonomous vehicles. This kind of technology will be very useful in the near future, where we expect the introduction of huge amount of autonomous devices.

## Acknowledgements

## References

[1] Kim, J.W., Kim, J.H. and Kim, I. (2015) SPQUSAR: A Large-Scale Qualitative Spatial Reasoner Using Apache Spark. *KIISE Transactions on Computing Practices*, **21**, 774-779. https://doi.org/10.5626/KTCP.2015.21.12.774

[2] Chen, L., Ko, J.H. and Yeo, J.M. (2014) Performance Comparison of DW System Tajo Based on Hadoop and Relational DBMS. *KIPS Tr. Software and Data Eng*, **3**, 349-354 https://doi.org/10.3745/KTSDE.2014.3.9.349

[3] Lim, Y.-H. (1995) ComBiStation: A Computer Platform for a Distributed Multimedia Computing Environment. *Journal of KIISE*, **2**, 160-181.

[4] Lee, W.-H. and Lee, B.-H. (2012) Service Delivery Time Improvement Using HDFS in Desktop Virtualization. *Journal of KIICE*, **16**, 913-921.
https://doi.org/10.6109/jkiice.2012.16.5.913

[5] NVIDIA (2017) Jetson.
http://www.nvidia.com/object/embedded-systems-dev-kits-modules.html

[6] NVIDIA (2017) DevBox.
https://developer.nvidia.com/devbox

[7] NVIDIA (2017) Digits.
https://developer.nvidia.com/digits

[8] ORACLE (2016) Virtual Box.
https://www.oracle.com/virtualization/virtualbox/index.html

[9] Apache Ambari (2016).
https://ambari.apache.org/