Chapter 4 Virtual Supercomputer Using Volunteer Computing

Rajashree Shettar *R. V. College of Engineering, India*

Vidya Niranjan R. V. College of Engineering, India

> V. Uday Kumar Reddy CA Technologies, India

ABSTRACT

Invention of new computing techniques like cloud and grid computing has reduced the cost of computations by resource sharing. Yet, many applications have not moved completely into these new technologies mainly because of the unwillingness of the scientists to share the data over internet for security reasons. Applications such as Next Generation Sequencing (NGS) require high processing power to process and analyze genomic data of the order of petabytes. Cloud computing techniques to process this large datasets could be used which involves moving data to third party distributed system to reduce computing cost, but this might lead to security concerns. These issues are resolved by using a new distributed architecture for De novo assembly using volunteer computing paradigm. The cost of computation is reduced by around 90% by using volunteer computing and resource utilization is increased from 80% to 90%, it is secure as computation can be done locally within the organization and is scalable.

DOI: 10.4018/978-1-5225-2785-5.ch004

Copyright © 2018, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

1. INTRODUCTION TO NEXT GENERATION SEQUENCING

Modern quantitative biology has changed the perspective of data rich genomic sequencing technology. Large scale genomic data analysis requires the need for a new computational framework supported by High Performance Computing. One such application is the Next Generation Sequencing (NGS), which deals with terabytes or petabytes of genome data requiring high computational power.

Next Generation Sequencing (NGS) (Wilson et al., 2002; Narzisi et al., 2011) is a technique of sequencing the exact order of nucleotides which form the basic building blocks of Deoxyribonucleic Acid (DNA). NGS with a market size of over 2.7 billion dollars has diverse uses in fields of biological sciences ranging from identification of diseases in human beings to invention of sequence for novel species. Traditionally sequencing was done by treating DNA chemically and identifying nucleotides using color codes, but this technique of sequencing is not suitable for organisms with just thousands of nucleotides. Earlier, the cost of producing base pair information stored as 'reads' was limited to wet laboratory techniques and was very expensive. Hence the rate of production of data was very slow, but new sequencing technologies combined with wet lab techniques and information technology started producing millions to billions of short 'reads' quickly. The traditional assembly tools used earlier was incapable to handle this huge data.

To overcome these problems a number of assembly technologies have been invented that uses computations performed by computer, also known as the *In silico* approach. These assemblers started with small datasets and were effective. As the size of 'reads' increased, the assemblers required either a single computer with very large amounts of memory and computing resources or the data to be sent to third party for execution such as cloud computing which might lead to security concerns. These constraints make the analysis of huge amount of genomic data a tedious task.

An alternate solution to Cloud and Hadoop is to use volunteer computing which is proposed and explained in this chapter. In particular emphasis is on recommending a solution to Next Generation Sequencing (NGS) which uses an open source grid middleware namely Berkeley Open Infrastructure for Network Computing (BOINC) designed to handle various applications that require high computational power, data storage or both. This will be a great enabler for bioinformatics scientists to create applications that use public computing resources.

1.1 Importance of Big Data and Cloud in Sequencing

Bioinformatics domain has brought in lot of challenges with respect to management of enormous amount of genomic data that is growing exponentially. Modern Biology

23 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: <u>www.igi-</u> <u>global.com/chapter/virtual-supercomputer-using-volunteer-</u> computing/188124

Related Content

Realm Towards Service Optimization in Fog Computing

Ashish Tiwariand Rajeev Mohan Sharma (2019). *International Journal of Fog Computing (pp. 13-43).* www.irma-international.org/article/realm-towards-service-optimization-in-fog-computing/228128

From Cloud Computing to Fog Computing: Platforms for the Internet of Things (IoT)

Sanjay P. Ahujaand Niharika Deval (2018). *International Journal of Fog Computing* (pp. 1-14).

www.irma-international.org/article/from-cloud-computing-to-fog-computing/198409

Fairness-Aware Task Allocation for Heterogeneous Multi-Cloud Systems

Sanjaya Kumar Panda, Roshni Pradhan, Benazir Nehaand Sujaya Kumar Sathua (2015). *Advanced Research on Cloud Computing Design and Applications (pp. 147-170).*

www.irma-international.org/chapter/fairness-aware-task-allocation-for-heterogeneous-multicloud-systems/138503

Predictive Modeling for Imbalanced Big Data in SAS Enterprise Miner and R

Son Nguyen, Alan Olinsky, John Quinnand Phyllis Schumacher (2018). *International Journal of Fog Computing (pp. 83-108).*

www.irma-international.org/article/predictive-modeling-for-imbalanced-big-data-in-sasenterprise-miner-and-r/210567

E-Vote: A Cloud-Based Electronic Voting System for Large-Scale Election

Ionut-Mihai Posea, Marius Ion, Florin Pop, Decebal Popescuand Nirvana Popescu (2016). *Cloud Computing Technologies for Connected Government (pp. 188-213).* www.irma-international.org/chapter/e-vote/136878