# Cloud-Based Big Data Analytics in Bioinformatics: A Review

Cephas MAWERE[1], Kudakwashe ZVAREVASHE[2], Thamari SENGUDZWA[3], Tendai PADENGA[4]

[1]*Harare Institute of Technology, School of Industrial Sciences & Technology Biotech. Dept., P.O. Box BE277 Belvedere, Harare, Zimbabwe*

*Tel: +263 772 384985, +263 4 741 422-58, Fax: +263 4 741 406*

*Email:cmawere@gmail.com*

[2]*Harare Institute of Technology, School of Information Science IT Dept., P.O. Box BE277 Belvedere, Harare, Zimbabwe*

*Tel: +263 733 784352, +263 4 741 422-36, Fax: +263 4 741 406*

*Email:kzvare@gmail.com*

[3]*Harare Institute of Technology, School of Industrial Sciences & Technology Biotech. Dept., P.O. Box BE277 Belvedere, Harare, Zimbabwe*

*Tel: +263 774 088979, +263 4 741 422-57, Fax: +263 4 741 406*

*Email: thamaryc@gmail.com*

[4]*Harare Institute of Technology, Dean of School of Information Science, P.O. Box BE277 Belvedere, Harare, Zimbabwe*

*Tel: +263 775 958368, +263 4 741 422-36, Fax: +263 4 741 406*

*Email: tepadenga@gmail.com*

## Abstract

The significant advances in high- throughput sequencing technologies over the last decade have led to an exponential explosion of biological data. Consequently, bioinformatics is encountering unprecedented challenges in data storage and analysis, memory allocation, computational power and time complexity. It is therefore becoming increasingly disturbing for small laboratories and some large institutions to establish and maintain infrastructures for data processing. This has forced bioinformatics to take a leap-forward from in- house computing to cloud computing to address these issues. In this paper, the authors thus explore the current applications of big data analytics in bioinformatics on the cloud. The findings in each in each use case reveal that cloud computing remarkably improve processing and storage of the continuously growing data in bioinformatics. Though there are some bottlenecks to solve, cloud computing promises to provide a lightweight environment to develop pipelines for prognosis, diagnosis, drug discovery and personalized medicine. The knowledge generated by this review will help scientist who are facing challenges with big data to resort to use cloud computing as an alternative solution and get results faster and without the need to install and maintain or update expensive software. This paper also serves to bring awareness to scientists on the current technologies that are used to manage data.

**Keywords:** High-throughput Sequencing Technologies; Bioinformatics; Data Processing; In- House Computing; Cloud Computing; Big Data Analytics.

# 1. Introduction

In the last seven years, more scientific data has been generated than in the entire human history. Hide, W. (2012). This has been triggered by the recent advances in high throughput technologies such as Next Generation Sequencing (NGS), genotyping, Single Nucleotide Polymorphism (SNP) discovery, gene expression, proteomics as well as virtual screening Blankenburg, L, et al (2009). For example, GenBank has seen DNA sequencing data doubling every 14 months in the last decade as attested by Figures 1-3. As on 21 October 2014, sequence data has been recorded to have reached $10^{8}$ bytes while the number of bases has catapulted to $10^{11}$ and $10^{12}$ bytes for GenBank WGS data respectively GeneBank. (2014).
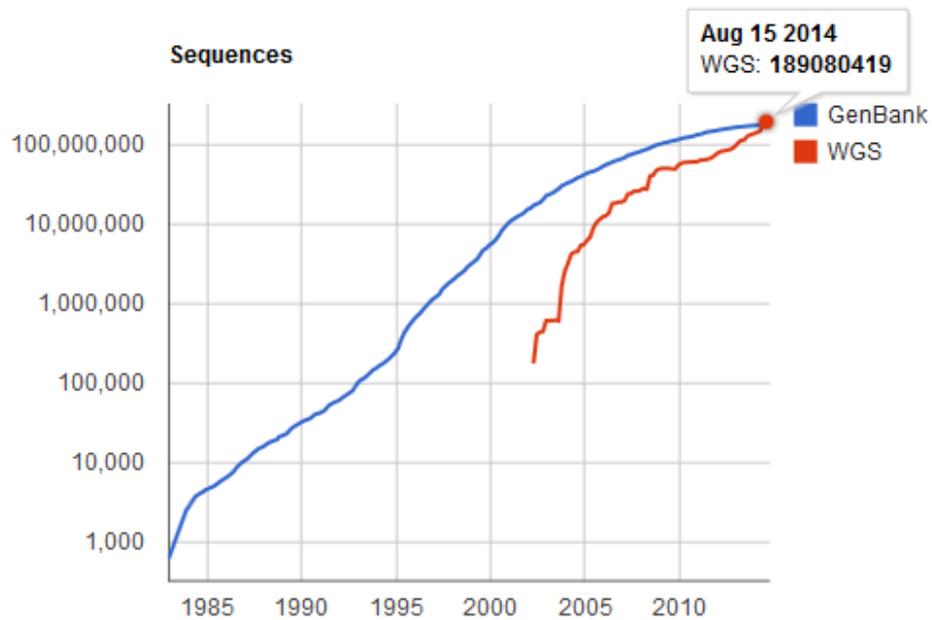


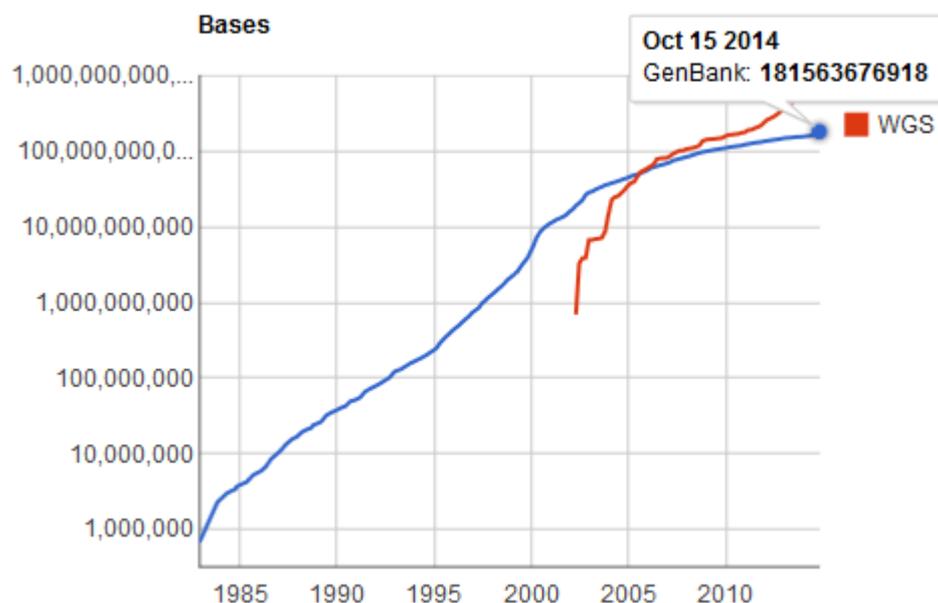Figure 1: Sequence data from GenBank statistics [3]

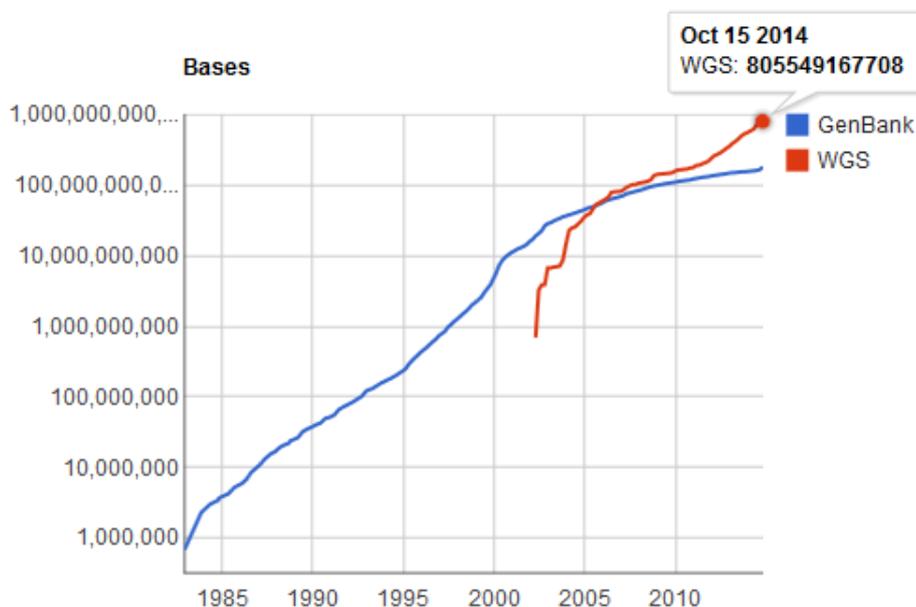Figure 2: GenBank statistics on number of bases [3]



Figure 3: Whole Genome Shotgun (WGS) base data statistics [3]

Inevitably, biological data has exploded beyond the capacity of computers to capture, store, analyse, search, share, visualize and transfer it for processing in the required time (Monica, N. & Kumar, R. (2013), Such big data has daunted small laboratories, hospitals and even some large institutions to establish and maintain infrastructures for data processing Dai, L., et al (2012). This has forced bioinformatics to leap forward from in-house computing to cloud computing which promises to address these issues.

But what is big data? Well, the term big data is defined as "a collection of data sets so large, varied, frequently updated and complex that it becomes difficult to process using on-hand data management tools" Movirtu. (2014a) Movirtu. (2014b). The process of research into such data to reveal hidden patterns and secret correlations is thus called 'big data analytics' Implementa. (2014). This research therefore seeks to explore the current applications of big data analytics in bioinformatics on the cloud. In each of the applications the authors will highlight how cloud computing has improved the processing of such exponentially growing biological data.

The paper will thus be composed of 5 categories. Section 2 will look at the overview of cloud computing and its application in bioinformatics. Consequently section 3 elaborates the applications of bioinformatics on the cloud for each area of high throughput technology and section 4 brings out the advantages and challenges of big data analytics on the cloud. Finally the conclusion and recommendation sums it up in section 5 and 6 respectively.

## 2. Cloud Computing In Bioinformatics

### 2.1 Overview of cloud computing

The buzzword 'cloud computing' was inspired by the cloud symbol that is often used in reference to the Internet in flowcharts. However, cloud computing is a form of computation

that exploits the full potential of multiple computers and delivers dynamically allocated virtual resources via the internet Telco Review (2014).  These hosted resources include storage, computation, applications, servers and network. Such services can be accessed on-demand (in a pay-as-you-go model) by any user via Web Application Programming Interfaces (API), without the need of the user knowing where the services are hosted or how they are delivered.

Such virtualization of data in the cloud enables 'snapshots' of large data sets to be moved from point to point within the cloud at high transfer rates without associating particular machines or storage devices at either source or destination of the data transfer Butte, A. J. & Dudley, J. T. (2010). This has necessitated Internet-based companies like Google, Windows and Amazon, whose cloud architectures can harness petabyte scales of data, to offer on-demand services to tens of thousands of users simultaneously Lin, Y et al (2013).

### 2.2 Bioinformatics resources on the cloud

  Bioinformatics clouds thus involve a large variety of services for big data analytics. Generally, the clouds fall into four categories (Figure 4) which are Data as a Service (DaaS), Software as a Service (SaaS), Platform as a Service (PaaS) and Infrastructure as a Service (IaaS) Stanoevska-Slabeva, K. & Wozniak, T. (2010).  DaaS enables dynamic data access on demand and provides up-to-date that are accessible by a wide range of devices that are connected over the Web. For instance, Amazon Web Services (AWS) which provides a centralized repository of public data sets like 1000 Genomes.
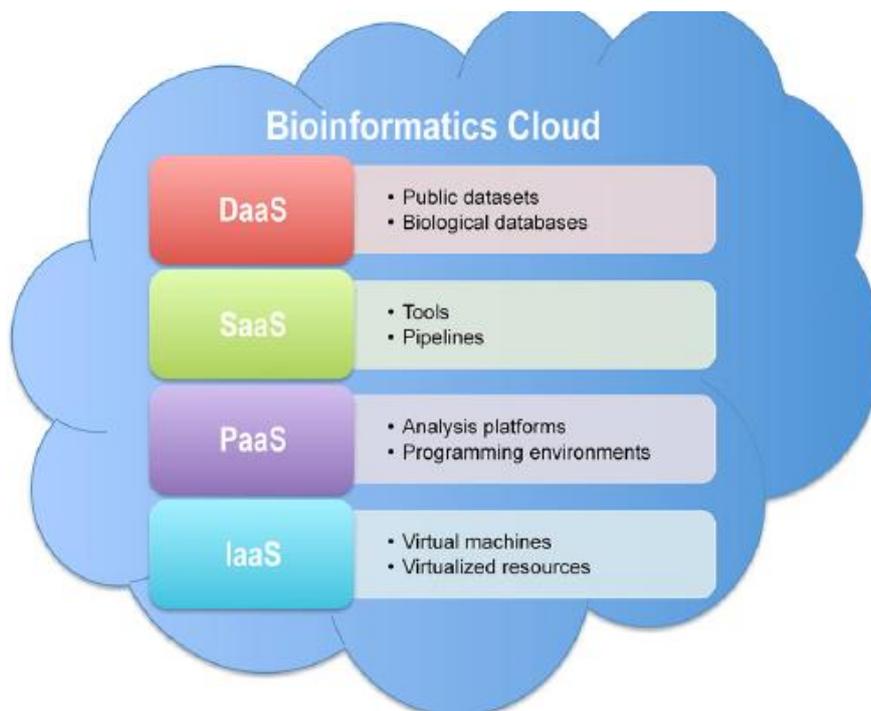


Figure 4: Illustration of bioinformatics resources on the cloud Dai, L (2012), Stanoevska-Slabeva, K. & Wozniak, T. (2010)

IaaS offers full computer infrastructure in form of virtualized resources like CPUs and operating systems via the Internet. On the other hand, PaaS provides an environment for users to develop, test and deploy cloud applications where computer resources automatically scale

to match application demand. The platforms commonly used thus include the Amazon Elastic Compute Cloud (EC2), Windows Azure and the Google app engine. This research however focuses on the use of SaaS which delivers software services online and facilitates remote access to available bioinformatics software tools via the Internet.

## 3. Applications of bioinformatics on the cloud

In particular, this paper reviews the use of cloud computing to address big data generated in bioinformatics areas such as (i) sequence mapping, (ii) sequence alignment, (iii) sequence analysis, (iv) peak calling for ChIP-seq data, (v) identification of epistatic interactions of SNPs as well as (vi) drug discovery. These applications are run on Amazon and Windows Azure platforms, with a parallelization implementation of Hadoop's MapReduce as illustrated by the example of CrossBow in Figure 5.

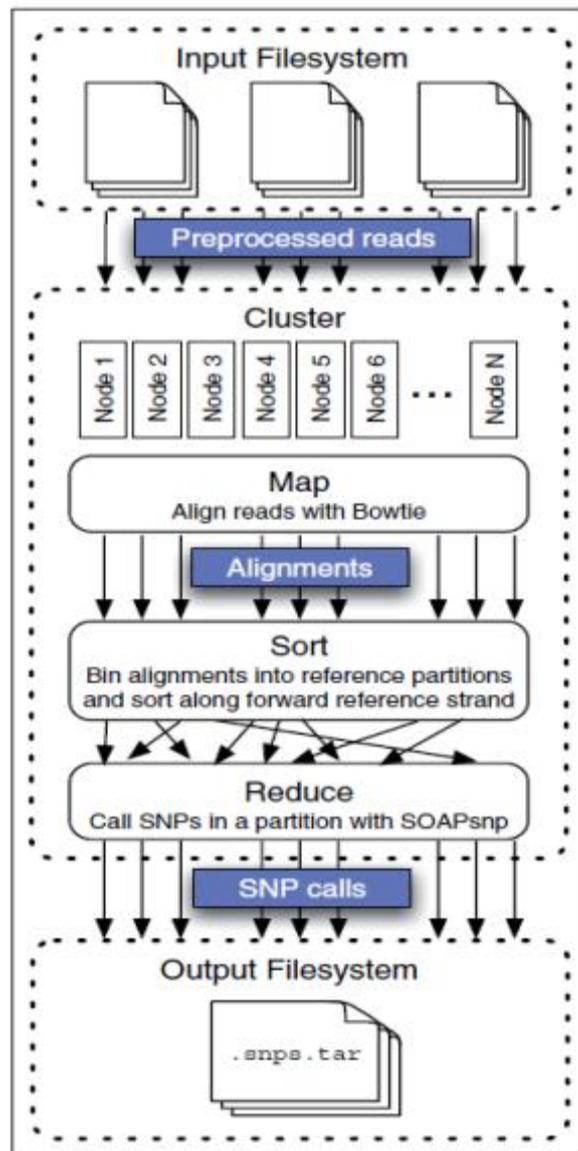### 3.1 CrossBow- for SNP calling and read mapping

Figure 5: Calling SNPs with CrossBow Russell, J. (2012)

The MapReduce on Crossbow thus includes the Map step, the Sorting step and the Reduce step. The Map step uses the Bowtie algorithm (named after the Burrows Wheeler Transform, BWT) to map sequence reads (which are short DNA sequences about 25-250 base pairs (bp) long) in parallel to the reference genome. The sequence alignments generated are then aggregated so that all alignments on the same chromosome or locus are grouped together and sorted by position. The sorted alignments are then scanned to identify SNPs within each region using the SOAPsnp (Short Oligonucleotide Analysis Package for SNP) algorithm. The results are stored via the Hadoop Distributed File System (HDFS), and then archived in SOAPsnp format Russell, J. (2012)

In the benchmark set on Amazon EC2 cloud, 2.7 billion reads were genotyped in less than 3hrs using a 320 CPU cluster for a cost of $85. Below is a snapshot of Crossbow on the cloud showing how one can genotype or call SNPs from a genome.



Figure 6: A snapshot of CrossBow on the Amazon EC2 (Gurtowski, J., et al, 2013). On this graphical user interface (GUI) any user can have his job processed on a pay-as-you go service by specifying the processing details in the textboxes.

## 3.2 CloudBurst

Cloudburst is a new highly sensitive parallel read- mapping algorithm optimized for mapping next-generation sequence data to reference genomes. Also, it uses the open-source Hadoop implementation of MapReduce to parallelize execution using multiple compute nodes (Schatz, M. C., 2009). In a case study done by Schatz in 2008, it was noted that CloudBurst scales linearly as the number of reads increases, and with near linear parallel speedup as the size of the number of processors increases (Schatz, M. C., 2009). In a 24- processor core configuration, CloudBurst is up to 30 times faster than RMAP executing on a single core, while computing an identical set of alignments.
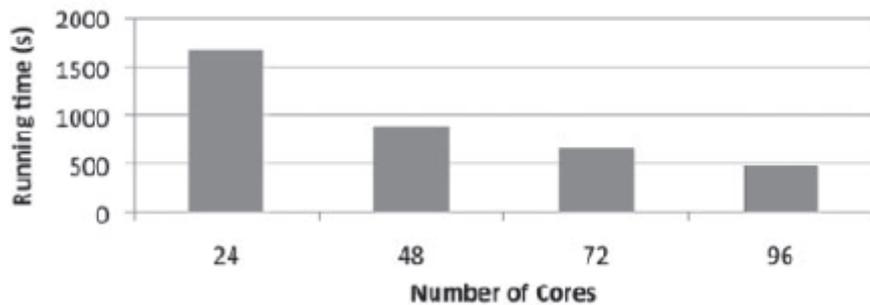


Figure 7: Running time on EC2 High-Medium Instance Cluster. Comparison of CloudBurst running time while scaling the size of the cluster for mapping 7million reads to human chromosome 22 with at most four mismatches on the EC2 Cluster. The 96- core cluster is 3.5X faster than the 24- core cluster Schatz, M. C. (2009)

Using a larger remote EC2 with 96 cores, CloudBurst improved performance by more than 100 fold. The running time was also reduced from hours to mere minutes for typical jobs involving mapping of millions of short reads to the human genome as shown by Figure 7.

### 3.2.1    Genome-Wide Epistatic Interactions

Following high-throughput SNP genotyping, the complex relationship between genotype and phenotype can now be unravelled on the cloud. As the cardinality of SNPs combinations gets larger, a dynamic clustering approach is now used for detecting high- order genome- wide epistatic interactions. In one case study, the speed-up evaluation of the clustering method (DCHE) was performed on Windows Azure platform on datasets with different sizes ranging from 10 000 to 100 000. The number of nodes varied from 1 to 40 with 5 minor units (Figure 8).
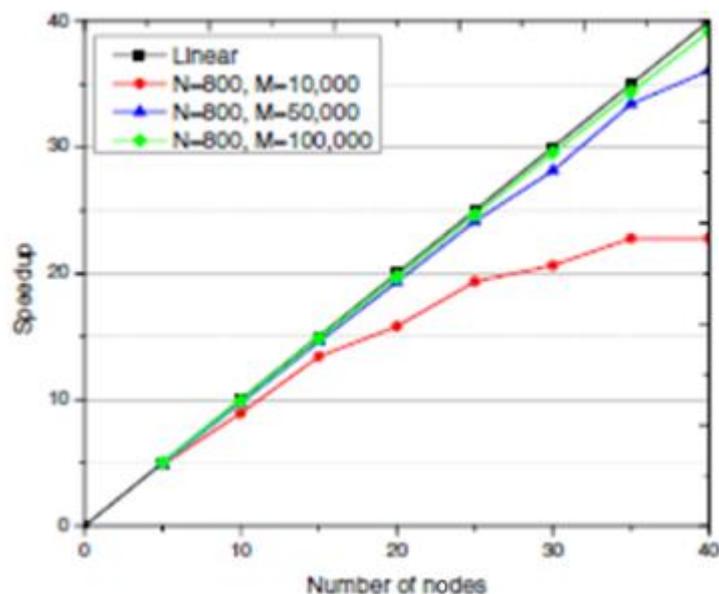
Figure 8: Genome-Wide Epistatic Interactions. Computing nodes are sampled from 1-40 with 5 as interval. The red, blue, cyan and grey curve show functions of speed-up of M=10 000, M= 50 000, and M= 100 000, with size fixed to 800 Guo, X 2014).

The results showed that when there are a limited number of SNPs in the dataset, for instance, M= 10 000, it has a lower speed-up curve. This is because dataset with such size along SNP dimension can be easily processed by a stand-alone version program alone in a couple of minutes (Guo, X, 2014). However, as the size of datasets increases, speed-up performs better (green and blue curves). Therefore, the cloud implementation of DCHE has a very good performance with respect to the speed-up.

### 3.3 Virtual Screening

Virtual Screening is an established time and cost- effective alternative to in vitro experiments to enable initial hit finding in the early stage of drug discovery (Certara. 2013). A virtual screening experiment using Surflex- Dock has been performed in the Amazon Compute Cloud to find lead structures that bind to the target under investigation. It consisted of 375 000 ligand-receptor dockings which were then run on 160 cores. Accordingly, the experiment was completed in < 9hrs for a total of less than $20.
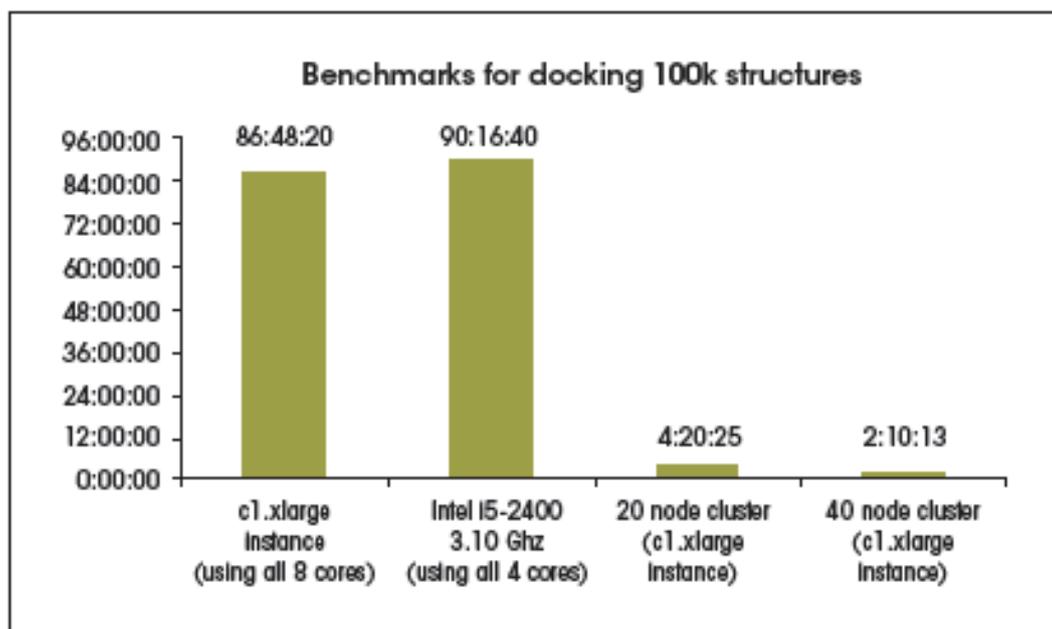
Figure 9: Surflex-Dock on Amazon EC2.One C1.xlarge EC2= modern computer with 4 cores (-90hrs).

An EC2 cluster of 40 c1.xlarge instances will thus perform in 2 hours Certara. (2013)

Prior to this, a benchmark study for docking 100 000 structures into a single receptor conformation was done (Figure 9). It was realised that a single c1.xlarge EC2 instance with its eight cores delivers nearly the same performance as a reasonably modern desktop computer with four cores (-90 hours), even though equipped with less CPU power per core. An EC2 cluster utilizing 20 and 40 c1.xlarge instances is able to perform the same task in 4 and 2 hours, respectively (Certara, 2013).

## 3.4 Peak calling

Chromatin immunoprecipitation (ChIP), coupled with massively parallel short-read sequencing (seq) is used to probe chromatin dynamics (Feng, X.et al ,2011.PeakRanger is an algorithm that is used to resolve closely-spaced peaks on both punctuate sites, such as transcription factor binding sites, and broad regions like histone modification marks.

In a use case carried out by Feng et. al., (2011) the MapReduce version of PeakRanger demonstrated that on a fixed number of  nodes with increasing data set sizes, the execution time was shorter and increased more slowly, than the regular single- processor version. For example, the cloud version processed 14 Gb dataset of 192 million reads in less than 5 minutes, more than 10 times faster than the original PeakRanger (Figure 10a).
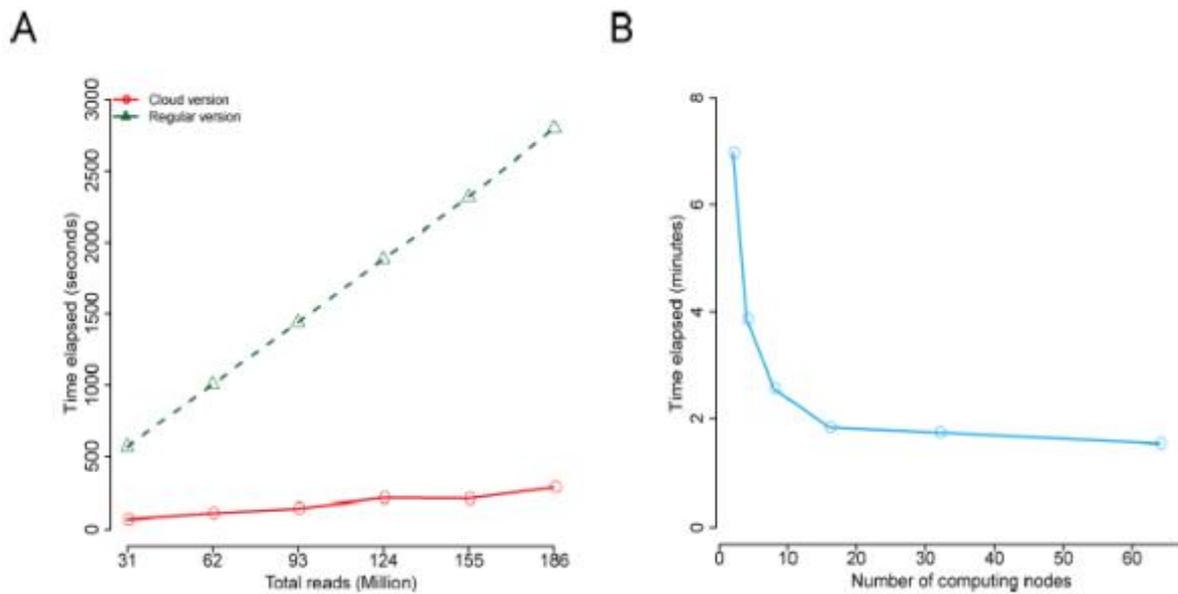
Figure 10: Performance of PeakRanger in cloud parallel computing. A) test with fixed number of nodes and data sets of increasing size; B)test with increasing numbers of nodes and data sets with fixed sizes (Feng, X., et al 2011)

In a second test, the scalability of time against the number of nodes was performed. As expected, the runtime decreased rapidly until the number of nodes equals the number of chromosomes. Further addition of nodes did not provide further benefit (Feng, X., et al 2011). Thus, PeakRanger on the cloud is provides marked improvements in run time and is therefore suitable for high-throughput environments.

## 4. Advantages and Challenges of Cloud Computing

### 4.1 Benefits of using bioinformatics on the cloud

From the use cases just explored on cloud implementation of bioinformatics tools, cloud computing has the following advantages:

- ✓ It provides a unified, location- independent platform for data & computation. Thus, it is available even to small labs which can harness its power for big data analytics in bioinformatics.
- ✓ It uses the 'Pay-per-use' model which gives access to any user.
- ✓ It is flexible; virtual resources can be shared.
- ✓ The costs incurred are affordable and relatively cheap for processing of big bio data
- ✓ It provides linear scaling of software tools with increasing number of nodes
- ✓ It harnesses the power of parallel execution provided by MapReduce
- ✓ It is time effective; execution occurs faster than on a local in- house cluster
- ✓ It eliminates the need for local installation and eases software maintenances and updates, providing up-to-date cloud based services.

### 4.2 Bottlenecks of bioinformatics on the cloud

Nevertheless, the following are challenges commonly faced on the cloud for big data analytics:

- Security: - trustability issues. How confidential is the data? There is also limited control over remote storage. Security updates are installed by vendors and the data is stored in a redundant manner. Data can be encrypted to secure it as this ensures no eavesdropping during data transmission (Forer, L et al, 2012).
- Data transfer: Large data sets need to be moved to the cloud. However there is slow upload speed because of low bandwidth. Some vendors thus have offered to ship data on hard drives (Forer, L et al, 2012).  A promising solution is the use of high speed transfer technologies such as data compression, Peer-to-Peer data distribution , as well as Aspera's fasp high speed file transfer technology (which speeds file transfers over the Web and outperforms traditional technologies such as FTP and HTTP) [http://www.aspeerasoft.com/en/technology/fasp_overview_1/fasp_technology_overv iew_1].
- Scalability: - need for parallelization, and depends on the algorithm itself if MapReduce suits the use case & desired effects.
- Usability: - bioinformatics involves a pipeline of scripts and a command line interface. Setting up nodes in the cloud involves a lot of command line and some deeper understanding. Also adapting new applications to the cloud still requires some technical expertise.

## 5. Conclusion

The significant advances in high throughput technologies like NGS have led to an explosion of biological data. This has given rise to bottlenecks in storage and processing of this data. This has forced bioinformatics to shift its focus from in- house computing to cloud computing to address these big data challenges. Here, SaaS was the primary focus of our review with emphasis on software that can be used on the cloud to address big data challenges in bioinformatics. Cloud computing of big data in bioinformatics can be done on several platforms such as Amazon EC2, Windows Azure and Google app engine. Most of these platforms implement the open source MapReduce software for parallel execution which makes the cloud preferred over in- house computing in addition to its benefits. Based on finding obtained from the explored use cases, the cloud- based resources make the big data meaningful and usable upon processing.

## 6. Recommendations and Future Works

High speed transfer technologies can be used in future to aid big data transfer. Despite the challenges existing in cloud computing, future efforts must also be devoted to open and publicly accessible bioinformatics clouds to the whole scientific community to accelerate diagnosis, prognosis, drug discovery and personalized medicine. Bioinformatics clouds should integrate both data and software tools and a lightweight programming environment should be provided to help in development of pipelines for data analysis.

## References

Blankenburg, L., Haberland, L., Elvers, H-D., Tannert, C. & Jandrig, B. (2009) 'High-Throughput Omics Technologies: Potential Tools for the Investigation of Influences of EMF on Biological Systems' *Curr Genomics*. 10 (2) pp. 86-92.

Butte, A. J. & Dudley, J. T. (2010) '*Ín silico* Research in the Era of Cloud Computing' *Nature Biotechnology*. 28 (11) pp. 1181-1185.

Certara. (2013) 'Cloud Computing Enables Cost Effective Virtual Screening' pp. 1-4. http://www.certara.com/ [accessed 29 August 2014].

Dai, L., Xiao, J., Zhang, Z., Gao, X. & Guo, Y. (2012) 'Bioinformatics Clouds for Big Data Manipulation' *Biology Direct*. 7 (43) pp. 1-7.

Feng, X., Stein, L. & Grossman, R. (2011) 'PeakRanger: A Cloud- Enabled Peak Caller for ChIP-seq Data' *BMC Bioinformatics*. 12 (139) pp. 1- 11.

Forer, L., Schonherr, S., WeinBensteiner, H., *et. al*. (2012) 'Cloud Computing' *Computational Medicine*. pp. 27-36.

GeneBank. (2014) 'GeneBank Statistics'http://www.ncbi.nlm.nih.gov/genbank/statistics. [accessed on Oct 21, 2014]

Guo, X., Meng, Y., Yu, N., & Pan, Y. (2014) 'Cloud computing for Detecting High-Order Genome-Wide Epistatic Interaction via Dynamic Clustering' *BMC Bioinformatics*. 15 (102) pp. 1471-2105.

Gurtowski, J., Schatz, M., & Langmead, B. (2013) http://www.bio-cloud-1449786154.us-east-l.elb_amazonaws.com/cgi-bin/crossbow.pl/ [accessed 8 November 2014]

Hide, W. (2012) 'The Promise of Big Data'. HSPH News. http://www.hsph.harvard.edu/news/magazine/spr12-big-data-tb-health-costs. [accessed 31 August 2014].

Implementa. (2014) http://www.implementa.com/solutions [accessed 10 June 2014]

Lin, Y., Yu, C., & Lin, Y. (2013) 'Enabling Large-Scale Biomedical Analysis in the Cloud' *BioMed Research International*. Volume 2013 pp. 1-6.

Monica, N. & Kumar, R. (2013) 'Survey on Big Data by Coordinating MapReduce to Integrate Variety of Data' *International Journal of Science & Research* . 2 (11). pp. 236-242.

Movirtu. (2014a) *Movirtu Virtual SIM Solutions*. http://www.movirtu.com/#!manyme/c1aq3 [accessed 10 September 2014]

Movirtu. (2014b) *Press Release*. http://www.movirtu.com/#!090614-airtel-movirtuagreement/c140h [accessed 10 September 2014]

Russell, J. (2012) 'Hadoop's Rise in Life Sciences' *Bio-IT World*. pp. 1-7.

Schatz, M. C. (2009) ÇloudBurst: Highly Sensitive Read Mapping with MapReduce' *Bioinformatics*. 25 (11) pp. 1363-1369.

Stanoevska-Slabeva, K. & Wozniak, T. (2010) 'Çloud Basics- An Introduction to Cloud Computing' In: Dai, L., Xiao, J., Zhang, Z., Gao, X. & Guo, Y. (2012) 'Bioinformatics Clouds for Big Data Manipulation' *Biology Direct*. 7 (43) pp. 1-7.

Telco Review. (2014) *Will the virtual sim break the telco strangle hold?*http://techday.com/telco-review/news/will-the-virtual-sim-break-the-telco-strangle-hold/195914/ [accessed 28 October 2014]

## Biographies

**Cephas Mawere**: Attained his BTech degree in Biotechnology at CUT, Zimbabwe in 2010. Completed his MTech in Bioinformatics at JNTUH, India in September 2014. He is a HIT Lecturer and Researcher. His research interests encompass the fields of Big Data, Cloud Computing, Information Security and Computer-Aided Drug Design

**Kudakwashe Zvarevashe:** Attained his BSc degree in Information Systems at MSU, Zimbabwe in 2010. Completed his MTech in Information Technology at JNTUH, India in September 2014. He is a HIT Lecturer and Researcher. His research interests are in the areas of Big Data, Information security, Cloud Computing and Web Services.