*Mini Review*

# Cloud Computing at a Discount: A Comparative Genomics Tool Case Study Summary

## Endre Aldington*

Department of computer science Gambia

*Corresponding Author's E-mail: Endre_al@7gmail.com

## Abstract

Scale and complexity of comparative genomics resources like ortholog detection tools and repositories are rapidly rising. In bioinformatics, cloud computing could be a useful alternative to large computational tools because it enables researchers to dynamically build a dedicated virtual cluster. Roundup, a large-scale comparative genomics resource, is optimized for cloud computing in this manuscript. We also detail important cost-effective procedures for ensuring maximum computation at minimal costs, as well as the proper operating principles needed to achieve cloud computing efficiency. On Amazon's Elastic Compute Cloud, we computed ortholog among 902 fully sequenced genomes using the comparative genomics tool Roundup as a case study. We devised a plan to deploy the web service Elastic Map Reduce, make the most of the cloud, and cut costs while simultaneously managing the ortholog processes. Based on the size and complexity of the genomes being compared, we developed a model that, in advance, determines the best order in which jobs should be submitted.

**Keywords:** Cloud computing, Elastic computing cloud, Roundup, Comparative genomics, High performance computing, Amazon, Orthologs

## INTRODUCTION

Roundup is one of many bioinformatics applications that predict the evolutionary relationship between genes, organisms, and biological functions by computationally comparing hundreds—and soon thousands—of genomes (Rosenthal A, 2009). The reciprocal smallest distance algorithm (RSD), which uses estimates of evolutionary distance and global sequence alignment to find orthologous genes between organisms, is at the heart of Roundup. The RSD calculation utilizes a computational pipeline that incorporates nearby homology worldwide arrangement (clustalw) and most extreme probability evaluations of developmental distances (paml) to find successions in a couple of genomes that have equally close transformative distances (Chenna R, 2003). Compared to reciprocal BLAST procedures, RSD's use of evolutionary distance and reciprocity can more effectively recover putative orthologs. Anyway this improvement in quality comes at the expense of expanded calculation. The demand for computational power increases proportionally to the square of the number of genomes for comparative genomics resources like Roundup (Altschul SF, 1990). According to a survey of the comprehensive resource for protein sequences, the number of fully sequenced genomes has nearly doubled annually since 2003, but processing speed has lagged, continuing to grow in accordance with Moore's Law, which states that it should double every 18 to 24 months. As a result, Roundup's computational requirements have increased eight times faster than processor capabilities. This disparity is consistent with other complex bioinformatics computational tasks and highlights the necessity of optimizing code even when massive compute capacity is available and looking for alternative computational processing strategies in general (Wall DP, 2011). Code reproducibility, pay-as-you-go pricing, rapid scaling, and the potential for 100 percent utilization are all features of loud computing. When contrasted with local clusters, which have highly variable usage patterns and the standard management and logistical issues that are typical of local compute farms, the demand pricing

model of cloud services is especially appealing. However, if the cloud's "elastic" nature is not fully utilized, its value is greatly diminished (Wall DP, 2007). Utilizing Amazon's cloud services and the comparative genomics tool Roundup, we demonstrate fundamental principles for optimizing large-scale bioinformatics operations in this manuscript. In order to cut costs, we emphasize the significance of optimizing code and explain the steps that must be taken to estimate costs on standard cloud infrastructures prior to launch. We demonstrate the way that our methods can significantly decrease distributed computing costs and examine how the standards we have created can be summed up to different kinds of bioinformatics undertakings (Altschul SF, 1990).

# METHODS

### Implementation of clouds

The cloud-based pipeline makes use of a combination of Simple Storage Service (S3), Elastic Compute Cloud (EC2), and Elastic Map Reduce (EMR), all of which are provided by Amazon Web Services. We previously published a series of methods to port the reciprocal smallest distance algorithm of Roundup to the Amazon cloud13. A pay-per-use model governs all of the services (Cost Breakdown). Roundup's code base, all genomes (Genome Data), executable programs linked to Roundup (BLAST, Clustalw, and codeml), and all computed results were stored using S3's persistent storage. We used EC2 instances with 8 cores, 7 GB of memory, and 1.7 TB of local storage for computations. We used EMR, Amazon's Hadoop on EC2 implementation, for workflow management. Hadoop is set up to run Java programs by default (Abdullah A, 2022). As a result, we came up with a method to take advantage of Map Reduce while still executing any program from the command line by utilizing the streaming mode of EMR. We were able to use our existing codebase without making any significant changes thanks to this strategy, which demonstrated that it could be easily adapted to any complicated task and codebase.

### Optimization of code

We optimized RSD's stability and speed for the purposes of this case study and to ensure maximum cost savings on the cloud (Jouini M, 2015). Explicit upgrades included diminishing circle I/O, expanding in-memory storing, and wiping out extreme summons of the Impact executable On a selection of 334 new genomes from the Roundup database, which contained 568 genomes at the time of writing, we tested the optimized configuration of RSD. A wide variety of sizes and sequence complexity sufficient to test the cloud's speed and cost were obtained from our selection of genomes, which included genomes from all branches of life. With an average of 4,415 sequences, the number of amino acid sequences in each genome varied by more than two orders of magnitude. The total amount of space required to host all 334 genomes was 4.3 GB.

### Model of computation time called "Bin"

To simulate the amount of time it would take to run RSD on a set of genomes in a cluster of a certain size, we created a model. We fit two linear models to account for the computational complexity of both stages of RSD, which is a two-stage computational process: first, BLAST is run on all pairs of genomes, and then ortholog detection is run using pre-computed BLAST results to calculate orthologs over a range of parameters. We first fit a linear model based on the product of the number of sequences in a pair of genomes and the computational time for BLAST in the RSD pipeline (Zhao S, 2013). Altschul et al demonstrated that the expected time complexity of BLAST is O (WN), where W is the linear length of a query sequence and N is the number of residues in the subject genome. This model was based on their findings. By extension, O (QN), where Q is the number of residues in the query genome, is used to run BLAST on all sequences in the genome. Based on the fit of the linear model ($r2 = 0.8$;), we used the number of sequences in a genome as a proxy for the number of residues. P = 2.2e-16), yielded estimations of time that were precise enough to order genome computations on a cluster to roughly minimize idle time. We calculated a linear relationship with time and the minimum number of sequences in a pair of genomes for the ortholog detection stage ($r2 = 0.3$; P = 2.2e-16). The ortholog detection stage, which internally performs O(1) work for each sequence in the genome with the fewest sequences, was the foundation for this (unpublished results). We used a preselected set of 30 genomes with sizes ranging from 75 to 66,710 sequences to fit these models. These genomes were a good representation of the phylogenetic diversity in our larger sample (Banu SS, 2018).

# DISCUSSION

The compensation per-utilize model of distributed computing gives an alluring option in contrast to superior execution registering calculations in bioinformatics. However, there aren't many published case studies that show how and when such infrastructure can be useful. The reciprocal smallest distance algorithm (RSD), a common ortholog detection tool, is used to demonstrate best practices for utilizing the Amazon Elastic Computing Cloud (EC2) in this manuscript. We explain in detail how the web service Elastic Map Reduce (EMR), which makes it easier to create and manage a large number of jobs in a cloud-based cluster of a predetermined size, can be used to run RSD on EC2. On Amazon's cloud, the EMR is a Map Reduce implementation based on Hadoop. The initial goal of Hardtop's development was to crawl, index, and process a lot of data for Internet search engines. Several businesses, including Yahoo!, Facebook and the New York Times have made excellent use of Hadoop. Job monitoring and flow management are just two examples of the many advantages offered by EMR's cluster computing model. However, batch-based algorithms like RSD do not work well with EMR or similar Hadoop-based computing frameworks.

# CONCLUSIONS

Cloud computing services like Amazon offer potential alternatives when confronted with computational bottlenecks brought on by the availability of local resources. The reciprocal smallest distance algorithm, or RSD, a common ortholog detection tool, was utilized in this manuscript to directly test Amazon's cloud's efficacy. The need to keep Roundup1's genome coverage up to date with the rapid pace of genome sequencing and our own local computing bottlenecks fuelled our work. We demonstrated how to make the most of Amazon's EMR web service for efficient job management at low costs. In particular, we planned a model to anticipate work runtime and costs for a variety of cloud group measures, and demonstrated the way that this model can be utilized to distinguish the ideal bunch size as well as the best system for requesting position preceding accommodation to the cloud.

# REFERENCES

1. Rosenthal A, Mork P, Li MH, Stanford J, Koester D, et al (2009). Cloud computing: a new business paradigm for biomedical information sharing. J Biomed Inform Apr. 43: 342–53.

2. Chenna R, Sugawara H, Koike T (2003). Multiple sequence alignment with the Clustal series of programs. Nucleic Acids Res. 31: 3497–500.

3. Altschul SF, Lipman DJ (1990). Protein database searches for multiple alignments. Proc Natl Acad Sci USA. 87: 5509–5522.

4. Wall DP, Fraser HB, Hirsh AE (2011). Detecting putative orthologs. Bioinformatics. 19: 1710–1721.

5. Wall DP, Deluca T (2007). Ortholog detection using the reciprocal smallest distance algorithm. Methods Mol Biol.396: 95–110.

6. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990). Basic local alignment search tool. J Mol Biol. 215: 403–413.

7. Abdullah A, Ajay DK (2022). Cost-Effective Spot Instances Provisioning Using Features of Cloud Markets. Int J Cloud Comput. 12: 1-27.

8. Jouini M (2015). Mean Failure Cost Extension Model towards Security Threats Assessment: A Cloud Computing Case Study. J Comput. 10: 184-194.

9. Zhao S, Prenger K, Smith L, Messina T, Fan H (2013). Rainbow: a tool for large-scale whole-genome sequencing data analysis using cloud computing. BMC Genomics. 14: 425.

10. Banu SS, Balasundaram SR (2018).Cost Effective Approaches for Content Placement in Cloud CDN Using Dynamic Content Delivery Model. Int J Cloud Comput. 8: 78-117.