

## Chapter 8.9

# PheGee@Home: A Grid-Based Tool for Comparative Genomics

**Bertil Schmidt**

*Nanyang Technological University, Singapore*

**Chen Chen**

*Nanyang Technological University, Singapore*

**Weiguo Liu**

*Nanyang Technological University, Singapore*

**Wayne P. Mitchell**

*Experimental Therapeutics Centre, Singapore*

### ABSTRACT

*In this chapter we present PheGee@Home, a grid-based comparative genomics tool that nominates candidate genes responsible for a given phenotype. A phenotype is the physical manifestation of the interplay of genetic, epigenetic and environmental factors. Our tool is designed to facilitate the discovery and prioritization of candidate genes controlling or contributing to the genetically determined portion of a specified phenotype. However, in order to make reliable nominations of candidate genes from sequence data, several genome-size sequence datasets are required. This makes the approach impractical on traditional computer architectures leading to prohibitively long runtimes. Therefore, we use a computational architecture based on a desktop grid environment and commodity graphics hardware to significantly accelerate PheGee. We validate this approach by showing the deployment and evaluation on a grid testbed for the comparison of microbial genomes.*

### INTRODUCTION

High-throughput techniques for DNA sequencing have led to an enormous growth in the amount of publicly available genomic data. As of Febru-

ary 2008, 716 complete genome sequences are available and another 2,756 genome-sequencing projects are in progress ([www.genomesonline.org](http://www.genomesonline.org)). As the sequences of more and more genomes become available, we have reached a critical mass where, instead of focusing on a subset of

DOI: 10.4018/978-1-4666-0879-5.ch8.9

sequences, we can use entire genome data sets to derive global inferences and metadata. Comparative genomics refers to the study of relationships between the genomes of different species or strains. It is currently being used for ortholog detection (Itoh, Goto, Akutsu & Kanehisa, 2005) bacterial pharmacogenomics (Fraser, et al., 2000), clustering of similar protein sequences (Itoh, Akutsu & Kanehisa, 2004), etc. Unfortunately, comparative genomics applications are highly computationally intensive tasks due to the large sequence data sets involved and typically take a few months to complete. These runtime requirements are likely to become even more severe due to the rapid growth in the size of genomic databases.

The objectives of this chapter are therefore two-fold:

1. The presentation of a new comparative genomics tool called PheGee (*Phenotype Genotype Explorer*). PheGee nominates candidate genes responsible for a certain phenotype  $\pi$  given genomic sequence datasets of phenotype positive ( $\pi^+$ ) and phenotype negative ( $\pi^-$ ) species.
2. The proposition of a hybrid computational grid platform to accelerate PheGee.

The proposed hybrid grid architecture efficiently combines desktop grid computing with GPGPUs (General-Purpose computation on Graphics Processing Units). The driving force and motivation behind this architecture is the price/performance ratio. Using desktop grids as in the volunteer computing approach is currently one of the most efficient and simple ways to gain super-computer power for a reasonable price. Installing in addition massively parallel processor add-on boards such as modern computer graphics cards within each desktop can further improve the cost/performance ratio significantly. We show how this architecture can be used to accelerate PheGee efficiently. Moreover, the proposed grid approach is

flexible and is therefore applicable to a variety of compute-intensive genomics applications.

## BACKGROUND

### Biological Motivation

PheGee exploits the simple observation: Given a set of species,  $A$ , which evince a particular phenotype,  $\pi$ , and a second set of species,  $B$ , which do not evince  $\pi$ , then in general genes responsible for  $\pi$  will be present in  $A$  and absent from  $B$ . Given incomplete knowledge, in the universe of sequenced organisms there will be a third set,  $C$ , generally the largest of the three, comprising those organisms undefined in terms of  $\pi$ . To cast this in symbolic terms, if a particular process is the outcome of three genes  $x \rightarrow y \rightarrow z$ , then in general homologs of  $x$ ,  $y$ , and  $z$  will be found in  $A$  but not in  $B$ . Those genes found in all members of  $A$  and not found in  $B$ , are therefore the best candidates for genes responsible for the process resulting from  $x$ ,  $y$ ,  $z$ . However, a perfect concordance between phenotype and genotype in sets of organisms is not guaranteed. Apart from epigenetic and environmental factors that can shape phenotype, the following genetic phenomena might disrupt the relationship.

- *Gene Replacement* (Achtmann & Wagner, 2008): Gene  $y$  can be lost in some members of  $A$ , say  $A^*$ . In  $A^*$  the function of  $y$  may have been replaced by some other gene,  $w$ . Thus, although the phenotype is preserved in all members of  $A$ , and although  $y$  is a gene that contributes to the  $A$  phenotype,  $y$  will nevertheless not map to every member of  $A$ . In particular,  $y$  is not found in  $A^*$ .
- *Horizontal Gene Transfer* (Achtmann & Wagner, 2008): Genes  $x$ ,  $y$ , or  $z$  may be adopted by a new species through the process of “horizontal gene transfer”. In this scenario some members of  $B$ , say  $B^*$ , will

17 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:

[www.igi-global.com/chapter/phegee-home-grid-based-tool/64572](http://www.igi-global.com/chapter/phegee-home-grid-based-tool/64572)

## Related Content

---

### Rough Entropy Clustering Algorithm in Image Segmentation

Dariusz Malyszko and Jarosław Stepaniuk (2010). *Novel Developments in Granular Computing: Applications for Advanced Human Reasoning and Soft Computation* (pp. 285-304).

[www.irma-international.org/chapter/rough-entropy-clustering-algorithm-image/44708](http://www.irma-international.org/chapter/rough-entropy-clustering-algorithm-image/44708)

### Mind Genomics With Big Data for Digital Marketing on the Internet

Jakob Salom (2021). *Handbook of Research on Methodologies and Applications of Supercomputing* (pp. 282-289).

[www.irma-international.org/chapter/mind-genomics-with-big-data-for-digital-marketing-on-the-internet/273407](http://www.irma-international.org/chapter/mind-genomics-with-big-data-for-digital-marketing-on-the-internet/273407)

### Predictive File Replication on the Data Grids

Chen Han Liao, Na Helian, Sining Wu and Mamunur M. Rashid (2010). *International Journal of Grid and High Performance Computing* (pp. 69-86).

[www.irma-international.org/article/predictive-file-replication-data-grids/38979](http://www.irma-international.org/article/predictive-file-replication-data-grids/38979)

### Sketch-Based 3D Model Retrieval Using Attributes

Haopeng Lei, Guoliang Luo, Yuhua Li, Jianming Liu and Jihua Ye (2018). *International Journal of Grid and High Performance Computing* (pp. 60-75).

[www.irma-international.org/article/sketch-based-3d-model-retrieval-using-attributes/205504](http://www.irma-international.org/article/sketch-based-3d-model-retrieval-using-attributes/205504)

### Resource Scheduling for Energy-Aware Reconfigurable Internet Data Centers

Mohammad Shojafar, Nicola Cordeschi and Enzo Baccarelli (2016). *Innovative Research and Applications in Next-Generation High Performance Computing* (pp. 21-46).

[www.irma-international.org/chapter/resource-scheduling-for-energy-aware-reconfigurable-internet-data-centers/159038](http://www.irma-international.org/chapter/resource-scheduling-for-energy-aware-reconfigurable-internet-data-centers/159038)