

# CloudQTL: Evolving a Bioinformatics Application to the Cloud

John Allen<sup>1</sup>, David Scott<sup>2</sup>, Malcolm Illingworth<sup>2</sup>, Bartek Dobrzelecki<sup>2</sup>, Davy Virdee<sup>2</sup>, Steve Thorn<sup>3</sup>, Sara Knott<sup>1</sup>

<sup>1</sup>Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9 3JT,

<sup>2</sup>EPCC, University of Edinburgh, JCMB, Edinburgh, EH9 3JZ,

<sup>3</sup>Information Systems, University of Edinburgh, JCMB, Edinburgh, EH9 3JZ

## Abstract

A timeline is presented which shows the stages involved in converting a bioinformatics software application from a set of standalone algorithms through to a simple web based tool then to a web based portal harnessing Grid technologies and on to its latest inception as a Cloud based bioinformatics web tool. The nature of the software is discussed together with a description of its development at various stages and the resulting successful increase in the user base. A discussion is then made detailing the latest idea to achieve a paid for service using Cloud technologies.

## Introduction

A quantitative trait is a phenotype or organism characteristic with continuous measurement such as product yield and quality in agricultural species or risk factors for disease in animal and human populations. It is usually complex in that it is influenced by the actions and interactions of many genes and environmental factors and geneticists are interested in identifying and understanding the role of the genes involved.

Quantitative trait locus mapping is a statistical modelling approach to identifying regions of the genome known as QTLs (Quantitative Trait Loci) that are involved in the control of the trait and is an essential tool for understanding the genetic basis of complex traits. It involves the use of molecular markers to follow inheritance of specific genome locations from parent to offspring and combines information from these with pedigree and trait records to look for associations between genotype and phenotype.

### 1990s to 2005 – Standalone Application to the World Wide Web.

Production and release of *QTL Express* [1], a user-friendly, web-accessible analysis tool, involved converting QTL mapping algorithms [2] initially written in Fortran into Java servlets. *QTL Express* allowed users to send data and receive output in series for simple QTL mapping analyses using moderately sized data of the order of kilobytes. It has seen wide use for the analysis of experimental data for QTLs, and it has received almost 500 citations.

### 2005-2010 - e-Science push - Grid Portal technologies

The advent of microarray technologies that produce high-density multiple trait gene expression datasets and the availability of dense gene marker maps for thousands of individuals increased the dimensionality and complexity of QTL analyses requiring computationally intensive and more advanced QTL mapping tools. This led to a push for more computational power, a need to develop more complex QTL algorithms as well as the ability to accommodate more users using larger data sets of the order of megabytes as the QTL community grew.

*GridQTL* [3] & [4] provided an expanded and improved QTL analysis tool from *QTL Express* in a user friendly web portal environment, harnessing Grid technologies to deal with these increased computational demands and offering data persistence, parallel submission and retrieval of data with access via a user login to a personal data space for reviewing results. Work started in 2005 and involved collaboration with the Institute of Evolutionary Biology (IEB), Roslin Institute, National e-Science Centre (NeSC), and EPCC. The web portal was based on GridSphere [5] that acted as a container to the QTL algorithms that had evolved once more into JSR 168 compliant Java portlets [6]. The portal uses the power of the NGS [7], ECDF [8] & [9] and, for very large data sets Hector [10] in the computational Grid. Grid middleware from the *Globus Toolkit* [11], and *Enabling Grids for e-Science project, EGEE* [12] were used for job-submission and querying methods as well as for management tools for the authentication and authorisation processes involved in the use of the Grid resources. A typical view of the portal during an analysis run is shown in Figure 1.

*GridQTL* was first released in the autumn of 2006 and demonstrated at the UK e-Science All Hands conference of that year [13]. To date nearly 500 individual users have performed near to 100000 analyses in their QTL studies and are now using around 2 cpu years of computation time on our Grid per year. Around 50 users a month use *GridQTL* in every continent of the world; a map detailing the location of our users who have cited *GridQTL* is available from our website [4] and is shown in Figure 2.

QTL Studies performed with *GridQTL* to date have included: birth weight and fleece quality in sheep; growth in young cattle; fatness in pigs; harvest traits in salmon; domesticity studies in foxes; obesity in mice; growth in broiler chickens; wood quality of eucalyptus trees; scale quality in crocodiles and airway obstructions in thoroughbred racehorses.

### **2010 and onwards – reaching for the Clouds.**

A further tranche of funding gave us the ability to include new QTL models in the portal as well as to investigate areas of Cloud computing. The *GridQTL* portal has so far given users access to the QTL algorithms and the computational resources free of charge; however, there is no way of sustaining this once the project funds run out.

Our view of Cloud Computing is in line with the view presented in [14]. Cloud Computing brings together Software as a Service (SaaS) and Utility Computing where Utility Computing is a service made available in a pay-as-you-go manner by the Cloud Provider. One can distinguish several classes of Utility Computing amongst the current Cloud computing offerings. The difference is based on the level of abstraction presented to the programmer wanting to access virtualised resources. For example the

Google AppEngine [15] provides automatic scaling and load balancing but enforces the programmer to use a predefined application structure and a fixed API. On the other side of the coin is Amazon's EC2 [16] which allows the author to control nearly the entire software stack but at the same time is not providing any help in automatic scalability or failover. There is also the middle ground represented by Microsoft's Azure platform [17] that supports general purpose computing but requires applications to be compiled to the specific runtime. *GridQTL* uses complex backend applications to perform calculations, and it was deemed to be too expensive to port these to new runtime environments. Only the fully virtualised model, similar to Amazon's EC2, was practical for moving the existing portal to Cloud infrastructure.

When developing *CloudQTL* we first sought the Amazon route via Eucalyptus [18] and OpenStack [19] middleware, both of which implement subsets of EC2 API, using a prototype local Cloud provided by the Edinburgh University ECDF Cloud; this would enable eventual Cloudbursting to similar Clouds implementing the EC2 API. Development of *CloudQTL* has however been considered with other Cloud scenarios in mind (e.g. OpenNebula [20] and OCCI [21]) so as not to tie ourselves to one specific access route to Cloud systems. In partnership with EPCC an initial version of *CloudQTL*, incorporated into the 3.0.x release of *GridQTL*, has been written and is undergoing tests prior to release to a selection of our clients. When the product proves robust a cost model accounting system based on EPCC's SAFE project [22] will be then considered for implementation.

## References

1. Seaton G, Haley CS, Knott SA, Kearsey M, Visscher PM: *QTL Express*: mapping quantitative trait loci in simple and complex pedigrees. *Bioinformatics* 2002, **18**:339-340.
2. Haley CS, Knott SA, Elsen JM : Mapping quantitative trait loci in crosses between outbred lines using least squares. *Genetics* 1994, **136**:1195-1207.
3. Seaton G., Hernandez J., Grunchev J.A., White I., Allen J., De Koning D.J., Wei W., Berry D., Haley C., Knott S. (2006) GridQTL: A Grid Portal for QTL Mapping of Compute Intensive Datasets. *Proceedings of the 8th World Congress on Genetics Applied to Livestock Production, August 13-18, 2006. Belo Horizonte, Brazil*. ISBN: 85-60088-01-6
4. GridQTL: <http://www.gridqtl.org.uk>
5. Novotny J, Russell M, Wehrens O: GridSphere: an Advanced Portal Framework. *Euromicro* 2004, 412-419.
6. The Java Community Process Program [<http://www.jcp.org/en/jsr>].
7. UK National Grid Service (NGS) [<http://www.ngs.ac.uk>].
8. <http://www.ecdf.ed.ac.uk>. 1 August 2007. U of Edinburgh. 7 Oct. 2008.
9. O Richards, M Baker. "GridPP and the Edinburgh Compute and Data Facility or How a general purpose cluster bore the weight of Atlas on its shoulders". Proc. UK All Hands Meeting 2008 (AHM2008)
10. HECToR [<http://www.hector.ac.uk>]
11. Foster I, Kesselman C: Globus: A Metacomputing Infrastructure Toolkit. *Intl J. Supercomputer Applications* 1997, **11**(2):115-128.
12. EGEE: Enabling Grids for E-science [<http://public.eu-egee.org>]
13. Proc. UK All Hands Meeting 2006 (AHM (2006))

14. Armbrust M., Fox, A., Griffith R., Joseph A.D., Katz R.H., Konwinski A., Lee G., Patterson D.A., Rabkin A., Stoica I., Zaharia M. Feb 2009. Above the Cloud: A Berkley view of Coud Computing Tech Rep U/CB/EECS-2009-28, EECS Department, /university of California Berkeley.
15. Google appEngine: <https://developers.google.com/appengine/>
16. Amazon EC2: <http://aws.amazon.com/ec2/>
17. Microsoft WindowsAzure: <http://www.windowsazure.com/en-us/>
18. Eucalyptus: <http://www.eucalyptus.com/>
19. OpenStack: <http://openstack.org/>
20. OpenNebula: <http://opennebula.org>
21. OCCI: <http://occi-wg.org/>
22. Project SAFE: <http://www.epcc.ed.ac.uk/projects/grid-safe>

## Figures

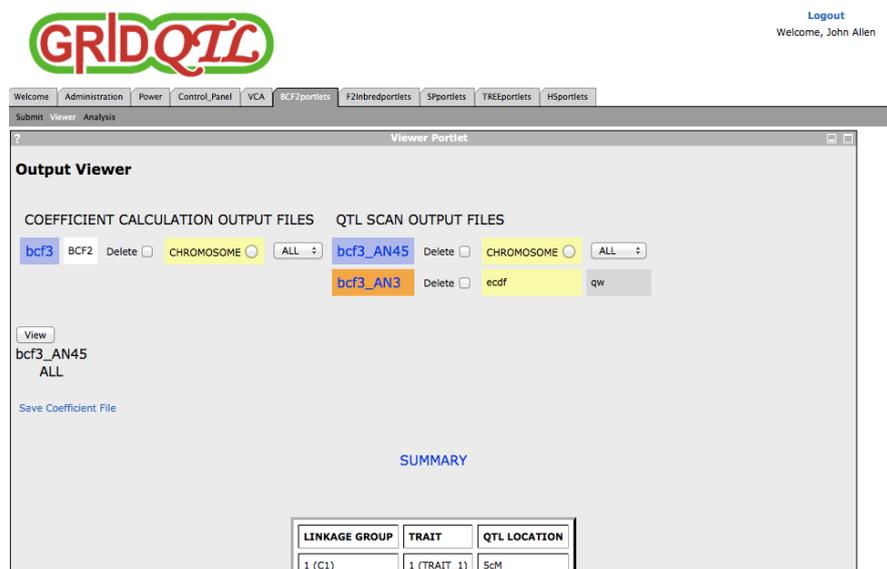


Figure 1 –Analysis Screen from the *GridQTL* Portal.



Figure 2 – GridQTL user citations by country.