

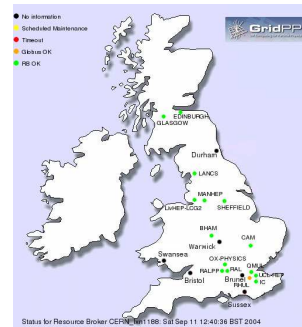
GridPP: Meeting the Particle Physics Computing Challenge

D. Britton (Imperial College), A.J. Cass (CERN), P.E.L. Clarke (University of Edinburgh), J.C. Coles (CCLRC), **A.T. Doyle** (University of Glasgow), N.I. Geddes (CCLRC), J.C. Gordon (CCLRC), R.W.L. Jones (Lancaster University), D.P. Kelsey (CCLRC), S.L. Lloyd (QMUL), R.P. Middleton (CCLRC), S.E. Pearce (QMUL), D.R. Tovey (University of Sheffield)
on behalf of the GridPP Collaboration.

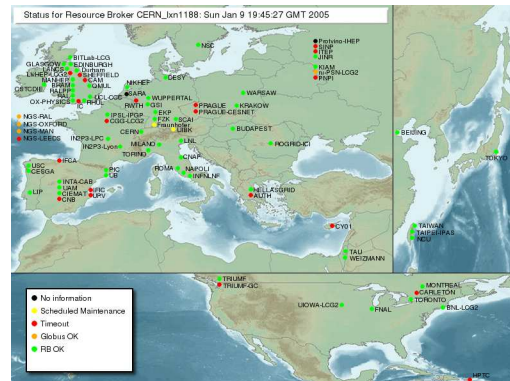
In 2007, following more than ten years of preparatory work, the Large Hadron Collider (LHC) at CERN, Geneva will start to collide protons with an energy equivalent of 7 trillion volts to recreate the conditions that prevailed in the universe at the earliest moments of the “Big Bang”. A billion interactions will be generated every second and amongst them perhaps one event will involve the production of a Higgs boson, responsible for the intrinsic mass of all the other fundamental particles. This will rapidly decay and its daughter particles will be detected in massive twenty metre high detectors with up to twenty million readout channels. The data will be efficiently filtered using dedicated electronics, but the data rates will still be enormous with around 10 PetaBytes of data produced each year from each of the four experiments. Thousands of physicists from all around the world will be eager to analyse the first data. Tens of millions of lines of analysis code will be written and the required processing power will be more than 100,000 processors operating continuously over many years. This is the nature of the LHC Computing Challenge.

The Grid is the chosen technology, a hardware and software infrastructure that provides dependable, consistent, pervasive and inexpensive access to high-end computational capabilities. The system must allow sharing of data between thousands of scientists with multiple interests; link major and minor computer centres across the globe; ensure all data is accessible anywhere, anytime; grow rapidly, yet remain reliable for more than a decade; cope with different management policies at different centres; ensure data security; and, be up and running routinely by 2007. The Grid is a practical solution to meet this and many other challenges in various areas of science, as part of the UK’s e-Science programme. Here we concentrate on the scale of the infrastructure, the tests that have been performed in recent months, and the plans for the second phase of the GridPP project, going from prototype to production.

In September 2001, the UK defined its aims in this area and the 3-year GridPP project commenced. The aim was to develop a highly functional prototype Grid deployed across the UK, consisting of 2000 CPUs, capable of accessing a PetaByte (PB) of data and linked to computer centres around the world.



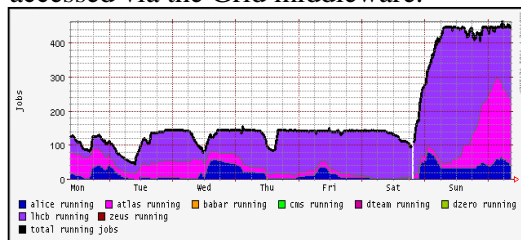
This was achieved by the GridPP collaboration via a Grid deployed at institutes across the UK. 2004 was a pivotal year, marked by extraordinary and rapid change with respect to Grid deployment, in terms of scale and throughput. The UK Grid is integrated seamlessly with the international LHC computing Grid project. The scale of the prototype is now significant with more than 10,000 CPUs linked and able to access 5 PB of data across more than 100 institutes worldwide.



This is illustrated in a recent snapshot of the LCG map, monitored via the Grid Operations Centre at the Rutherford Appleton Laboratory (RAL). Sites across Europe are linked to those in the Far East and North America. This prompted The Economist to declare the LCG the World’s Largest Grid in October 2004.

The UK is the biggest single contributor to the LCG, with more than a fifth of the Grid's processing power at its 16 sites. In addition 4 sites linked via the UK e-Science National Grid Service are using much of the same middleware to provide a service for a wide range of academic research projects. Middleware is the key to a successful Grid: the initial software stack of more than one million lines of code was developed and tested in conjunction with the European DataGrid and Enabling Grids for E-science projects, building upon earlier work by the US-based Globus and Condor projects. This enables a user in Glasgow, armed with his/her virtual passport (a digital certificate) to submit an analysis job to a resource broker. As a member of a recognised experiment (or Virtual Organisation, VO), the job will run a few minutes later on any of the sites with pre-installed VO-specific software and the required hardware for the job. This is a powerful generic methodology, applicable in many branches of science and industry, where access to large-scale computing resources are required, typically with demanding timescales, across the globe.

The system was intensively tested during 2004 via a set of planned stress tests by each of the 4 LHC experiments. Individual experiments accumulated up to 400 CPU years' worth of test data, with individual jobs running for up to a day and a peak load of almost 6,000 simultaneous jobs achieved in August 2004. The power of the Grid resource discovery method was illustrated in July 2004 when 300 new CPUs at RAL were brought online and automatically discovered and fully used within 8 hours, accessed via the Grid middleware.



Following these stress tests, acceptable (>90%) throughput was possible for each

of the experiments during 2005, but the inherent complexity of the system is apparent and many operational improvements are required to establish and maintain a production Grid of the required scale.

The GridPP project reached its halfway mark in September 2004 and started its second phase. Numerous issues have been identified that are now being addressed as part of GridPP2 planning in order to establish the required resource for particle physics computing. Further site and middleware validation tests are needed in order to improve the overall Grid efficiency. Each experiment needs to develop an individual application interface: in this regard it is noteworthy that the system developed for the large LHC experiments has been shown to work effectively for other less resource-intensive applications. Addressing analysis group computing within the experiments, developing distributed file and database management systems, installing and validating experiment software, setting up production accounting systems, and creating an environment where everyone can share resources are all areas requiring further development.

In GridPP1, 88% of the 190 milestones were completed and all 44 monitoring metrics were within specification at the end of the project. No milestones were missed and most of those outstanding were superseded by more detailed planning in GridPP2. The aim of GridPP2 is to deliver a "Production Grid": a robust, reliable, resilient, secure, stable service delivered to end-user applications. The Collaboration aims to develop, deploy and operate a 10,000 processor, multi-PetaByte Grid in the UK. To this end, the new project is more complex with 253 high-level milestones and 112 monitoring metrics currently defined. Tests of the prototype have given us some confidence that the final system will be capable of providing the required resource for LHC and other experiments' data analysis - in this way we plan to meet the Particle Physics Computing Challenge.