

Distributed Computing for Small Experiments

Daniela Bauer^{1,a}

¹*Imperial College London*

Abstract. The large Large Hadron Collider experiments have successfully used distributed computing for years. The same infrastructure yields large opportunistic resources for smaller collaborations. In addition, some national grid initiatives make dedicated resources for small collaborations available. This article presents an overview of the services available and how to access them, including an example of how small collaborations have successfully incorporated distributed computing into their workflows.

1 Introduction

The amount of data produced by most modern research almost always requires a degree of automated processing. The Large Hadron Collider (LHC) experiments chose a distributed computing model, where the data processing is shared among the collaborating institutes. To facilitate this they developed a dedicated software infrastructure. This infrastructure has evolved over the past decade, and is now readily usable by other communities, with a number of funding bodies stipulating resource allocation for non-LHC communities.

GridPP [1] is a consortium of 19 UK universities, which provides resources and grid specific computing expertise for any experiment with a UK affiliation. It currently allocates 10% of its resources to small collaborations while also allowing opportunistic use of its resources beyond this.

2 Distributed Computing

Most users will be familiar with a batch system, where they submit jobs to a local queue, which will then get processed in batch on a computer in their institute. Distributed computing expands on this by amalgamating resources across multiple institutions. In this case the users submit their tasks to a resource broker which in turn submits them to a suitable queue, either at their own institution or elsewhere (see figure 1). The Grid [2] is a specific implementation of a distributed computing model.

2.1 The Grid

Authenticating users across multiple domains is a challenge. The Grid's unified authentication model allows a user access to all grid resources using a single set of credentials. Each user is issued with a X.509 certificate [3] by their home institute. In practice this consists of a pair of encrypted text files the user installs in their home directory.

^ae-mail: daniela.bauer@imperial.ac.uk

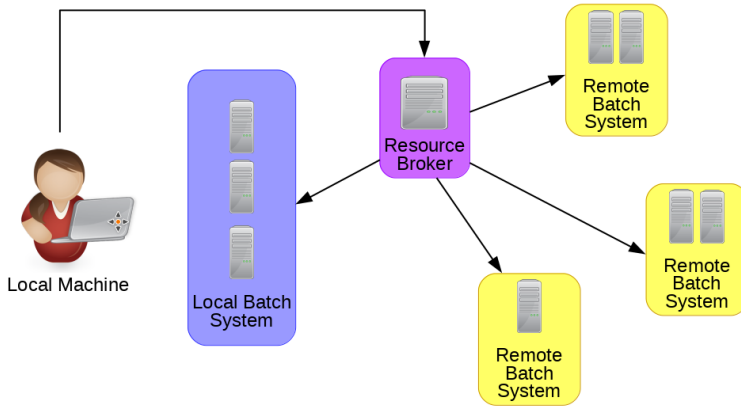


Figure 1. Diagram of a resource broker aided distributed computing model.

Most resources, especially data storage, which might hold confidential data, have a requirement of only being accessible by certain groups of researchers. The Grid uses the concept of a Virtual Organisation (VO) to accomplish this. Often VOs correspond to an actual experiment, examples include the CMS, ATLAS and LZ VOs. There are also VOs that cater to specific communities, e.g computational biology. These umbrella VOs also serve as incubators by reducing the overheads required for new projects to evaluate grid computing or for projects that require a one off use of additional computing resources. A user typically joins one or more VOs that match their research interests to gain access to the appropriate resources.

To avoid having to manually install VO specific software at all participating institutions, The Grid uses a software distribution system called CVMFS [4]. There the software is uploaded to a central server, where it is catalogued and hosted on a standard HTTP based server. When a client accesses a software repository, the required catalogue and file data are downloaded to the node via a series of web caches. As generally multiple similar tasks require the same software in the same place, the web caches ensure low latency access to the software repositories.

Tasks submitted by users are assigned to compute resources through dedicated resource brokers. There are a number of different types of resource brokers available, e.g. glideinWMS, BigPanda and DIRAC. GridPP provides a central DIRAC service which is available to all experiments with UK participation. These resource brokers typically already exist for all research communities, and requesting access for a new community is generally encouraged.

3 Case study: The LZ experiment

The LZ experiment [5] is a dark matter detector located in Sanford Lab, South Dakota. As the experiment is under construction, their current computing effort focuses on Monte Carlo production and analysis. Their data processing model also includes making their entire data accessible via The Grid, once data taking has started. The collaboration had no previous experience in distributed computing.

Initially they used the DIRAC Python API to manually submit batches of jobs. While this was suitable for low throughput Monte Carlo production, it became apparent that a more streamlined interface would be required for large scale production. A production interface [6] was developed by one of the participating institutes and successfully used for the first LZ Mock Data Challenge, which

produced a total of 133 TB of MC Data spanning 732288 files entirely on grid resources. LZ plans to use this production model for the foreseeable future.

4 Conclusions

Distributed computing has been successfully used by a wide variety of small collaborations with limited computing support to meet their data processing requirements. It is hoped that participation by small research communities will increase over time to encourage the shared use of existing resources.

References

- [1] D. Britton *et al.*, Phil. Trans. R. Soc. A **367** 2447-2457 (2009)
- [2] R. Brun, F. Carminati, *From the Web to the Grid and Beyond: Computing Paradigms Driven by High-Energy Physics* (Springer, Heidelberg, 2012)
- [3] M. Cooper *et al.*, RFC 4158
- [4] J. Blomer, P. Buncic, T. Fuhrmann, DOI 10.1145/2110217.2110225 (2011)
- [5] B.J. Mount *et al.*, arXiv:1703.09144
- [6] D. Bauer, S. Fayer, J. Phys.: Conf. Ser. **898** 052003 (2017)