Dissertation

"Service Level Agreement aware Resource Management "

Matthias Hovestadt

Abstract

Next Generation Grids aim at attracting commercial users to employ Grid environments for their business critical compute jobs. These customers demand for contractually fixed service quality levels, ensuring the availability of results in time In this context, a Service Level Agreement (SLA) is a powerful instrument for defining a comprehensive requirement profile.

Numerous research projects worldwide already focus on integrating SLA technology in Grid middleware components like broker services. However, solely focusing on Grid middleware services is not sufficient. Services at Grid middleware may accept compute jobs from customers, but they have to realize them by means of local resource management systems (RMS). Current RMS offer best-effort service only, thus they are also limiting the service quality level the Grid middleware service is able to provide.

In this thesis the architecture and operation of an SLA-aware resource management system is described, which allows Grid middleware components to negotiate on SLAs. The system uses its internal mechanisms of application-transparent fault tolerance to ensure the terms of these SLAs even in case of resource outages. The main parts of this work focus on scheduling aspects and strategies for ensuring SLA compliance, respectively design aspects on implementation.

Scheduling strategies significantly determine the level of fault tolerance that the system is able to provide. After presenting requirements of Grid middleware components on service qualities and a description of operation phases of an SLA-aware resource management system, intra-cluster scheduling strategies are described. Here, the system solely uses its own resources and mechanisms for coping with resource outages.

For further increasing the level of fault tolerance, strategies for cross-border migration are presented. Beside a migration to other cluster systems in the same administrative domain, the system uses also Grid resources as migration targets. For ensuring the successful restart, mechanisms for describing the compatibility profile of a checkpointed job are presented.

The concept of the SLA-aware resource management system has been implemented in the scope of the EC-funded project HPC4U. We will describe design aspects of this realization and show results from system deployments at use-case customers.

Keywords

Grid Computing, Grid Middleware, Grid, Resource Management, Service Level Agreement, SLA, Negotiation, Fault Tolerance, Scheduling, Process Checkpointing, Storage Snapshot, Network Fault Tolerance, Quality of Service, QoS