

# Toward a Virtual Grid Service of High Availability

Xiaoli Zhi and Weiqin Tong

School of Computer Engineering and Science, Shanghai University,  
Shanghai 200072, P.R. China  
{xlzhi, wqtong}@mail.shu.edu.cn

**Abstract.** A new regulation approach is proposed to obtain a virtual resource service of high availability and service capacity on the basis of resources of low availability and small capacity. Some regulation algorithms with distinct characteristics are introduced.

## 1 Introduction

With the widespread proliferation of Grid services, quality of service (QoS) will become a significant factor in distinguishing the success of service providers. In OGSA, anything providing some functions to the public can be treated as a virtual service. QoS of a service refers to its non-functional properties such as performance, reliability, availability, etc [1]. E.g., QoS of a processor can be measured by availability, computation capacity (in MIPS). QoS of a storage can be measured by availability, storage capacity (as in Terabytes) and so on. This paper will pay more interest in availability and service capacity (computation capacity or storage capacity).

Grid computing borrowed its term from the “Electric Power Grid” [2]. We got the idea of regulation from comparison of computational grids with the electrical power industry. Regulation of grid services is to achieve more stable, high available service from unstable source services just as the rectifier in the power grid get direct current from alternating current.

## 2 Regulation Algorithms

In a regulation system, several services, termed as source services, are organized into the ‘backend’ of a regulated service, as shown in Fig.1. The main task for a regulated service is to delegate service request to some appropriate source services according to its regulation algorithm. A regulated service appears nothing special externally just as a normal resource service on which a grid resource scheduler or broker can act.

The regulated service and its source service form a service aggregation. Different from the classic purpose of service aggregation, this aggregation is to serve as a buffer or stabilizer to passivate the sensitivity of service failure or variation.

The most important thing in service regulation is the regulation algorithm. There can be various regulation algorithms under different conditions to meet different requirements. Here presents only a small yet typical and important directory of regulation paradigms. The directory is open to evolve.

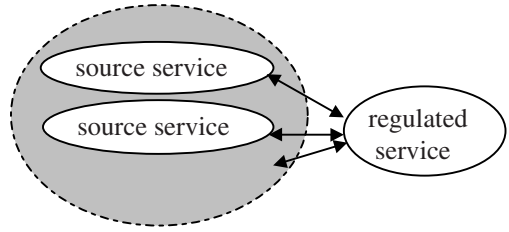


Fig. 1. Diagrammatic view of a regulated service

**2.1 Paradigm 1: Multiple-to-One**

Paradigm 1 is used to integrate service capacity or increase service availability.

**Paradigm11: Heaping.** The regulated service treats source services as a heap of service power and uses up the source services one by one. It can add up source services’ capacity to make a service of bigger capacity and higher availability than a single component service.

Assume there have M source services. Source services can provide normal capacity  $C_1, C_2, \dots$  with the probabilities  $p_1, p_2, p_3, \dots$ , that is, availabilities for source services with normal capacity. Unavailabilities or downtime probabilities are  $1-p_1, 1-p_2, 1-p_3, \dots$ , respectively. For simplicity, we presume source services only have two states: normal work or downtime/repair state. Then availability for the regulated service with service capacity  $c_r$ :

$$\text{Pro}(c_r = \sum_{i \in \text{subset of } \{1,2,\dots,M\}} C_i) = \prod_{i \in \text{subset of } \{1,2,\dots,M\}} p_i \tag{1}$$

The unavailability of the regulated service:

$$\prod_i (1 - p_i) \tag{2}$$

**Paradigm12: Stripping.** In this paradigm, workload is down into parts and each part is assigned to a separate source service. The service capacity and some performance parameters such as I/O speed for data transfer is greatly improved. However, it does offer a few disadvantages. It is more vulnerable than its source service. Availability for the regulated service with service capacity  $c_r$ :

$$\text{Pro}(c_r = \sum_i C_i) = \prod_i p_i \tag{3}$$

The unavailability of the regulated service:

$$1 - \prod_i p_i \tag{4}$$

**Paradigm13: Fault tolerant configuration.** The source services in this paradigm are configured as an active-active fault tolerant system. And workload is mirrored identically in every source service. This paradigm has the highest availability among paradigms introduced in this paper. Availability of the regulated service:

$$1 - \prod_i (1 - p_i) \tag{5}$$

### 2.2 Paradigm 2: Multiple-to-Multiple

Paradigm 1 normally demotes the utilization coefficient of source services although it can provide an integrated service of bigger capacity or higher availability. This paradigm will promote the utility factors as well as availability by grasping idle source service to serve in the place of a failed component service. A regulated service in Paradigm 2 has a designated main source service (Note every source service is the main service of a regulated service). The regulated service is just a transparent broker when its main source service runs normally. It acts what its main source service acts. But it will draft an idle other source service for the incoming task using some scheduling algorithm when its main source service is down.

Assume idle coefficients for source services (nothing to do while the service is ready to do something) are  $q_1, q_2, q_3, \dots$ . Then Availability for the regulated service  $k$  (assume the number of its main source service is also  $k$ ):

$$p'_k = p_k + (1 - p_k)(1 - \prod_{j \neq k} (p_j \times (1 - q_j) + (1 - p_j))) \tag{6}$$

Under the assumption of random choosing idle source services, idle coefficient for the source service  $k$  after it participates into a regulation system:

$$q''_k = q_k - \sum_j ((p'_j - p_j) \times (1 - q_j)) \times \left( \frac{p_k \times q_k}{1 - q_k} \right) / \sum_j \left( \frac{p_j \times q_j}{1 - q_j} \right) \tag{7}$$

The above formula is an approximation due to the highly intricacies of the accurate computation of  $q''_k$ . Actually, the idle coefficient for the source service after regulation is too difficult to be formulated mathematically when the scheduling algorithm is not a random one.

### 2.3 Paradigm 3: One-to-One

In the power industry, a rectifier has element of energy storage and transformation (e.g., capacitance or inductance) to regulate a fluctuating current into a smoother one. Paradigm 3 adopts a task buffer to serve as something like ‘energy storage and transformation’. The buffer revokes the source service when available to process buffered tasks and return tasks’ results at proper time. This paradigm is suitable for batch processing, asynchronous applications but not areas with high time demand.

Assume availability of the source service is  $p$ ; the service capacity for the source service is  $C$ ,  $B$  for the regulated service ( $B < C$ ), then availability for the regulated service:

$$p' = \frac{Cp}{B} > p \quad (8)$$

### 3 Discussions

A regulation service is composed of several source services as a service aggregation is. But they implement different purposes with completely disparate techniques. A service aggregation is to finish a task through cooperation of its components while regulation is targeted for improving services' QoS. Component services of a service aggregation are usually playing different roles with different functions, while those in regulation have similar service capability.

In some sense, regulation seems like a resource broker. But a regulated service actually distinguishes itself by different status in the grid. A regulated service is just a resource service. A resource broker is a part of or a base service in the grid middleware. And the scheduling (if any), interaction of a regulated service with its source services are hidden completely from the grid user while a resource broker doesn't. In addition, the scheduling algorithm and communication protocol, implemented in a regulated service, tend to be more likely proprietary, while they may not be adopted in a resource broker for out of standardization.

In summary, advantages of regulation are the following:

- ◆ Partly the resource management function will be spread around to virtual regulated services. This will alleviate the burden of resource management and enhance the grid middleware's reliability with higher available low-level resource services.
- ◆ Various proprietary resource scheduling algorithm or composition technique can be utilized in regulation within a relatively local scope.
- ◆ No additional complexity to the grid system middleware will be induced.

### References

1. A. Mani, A. Nagarajan: Understanding quality of service for Web services---Improving the performance of your Web services. January 2002, [http://www-900.ibm.com/developerWorks/cn/webservices/ws-quality/index\\_eng.shtml](http://www-900.ibm.com/developerWorks/cn/webservices/ws-quality/index_eng.shtml)
2. I. Foster and C. Kesselman: The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, San Fransisco, CA, 1999.
3. G. Mateescu: Quality of service on the grid via metascheduling with resource co-scheduling and co-reservation. International Journal of High Performance Computing Applications, vol.17, no.3, 2003, p 209-218