

# Towards Predictable Overload Control of Large Real-Time Server Systems First Year Report

Thiemo Voigt  
Computer and Network Architectures Laboratory  
Swedish Institute of Computer Science  
thiemo@sics.se

May 22, 2001

## 1 Project Plan and Achieved Results

In the project “Towards Predictable Overload Control of Large Real-Time Server Systems” we focus on overload protection mechanisms for Web servers.

As stated in the project plan we invested first in a mechanism that dynamically reduces the number of active server processes. The aim of this mechanism is to allow fast recovery from overload, in particular overload caused by lack of memory, without discarding already accepted connections. As shown [1], we were able to show the effectiveness of the idea.

During November/December Thiemo visited IBM TJ Watson to extend the paper originally submitted to Infocom 2001. We conducted experiments that proved that our overload mechanisms implemented in the kernel were much more efficient and scalable than mechanisms implemented in user space. The result of this work was paper [4]. At the same time we also ported our architecture to the latest AIX version.

After porting the overload protection architecture from AIX to Linux we worked on the next item discussed in the project plan. We extended the architecture to protect web servers from overload caused by persistent connections. Persistent connections are a challenging problem since the resource consumption of the requests on a persistent connection is unknown at the time the admission control decision has to be made. In an overload situation caused by resource demands of persistent connections our architecture aborts persistent connections that are regarded as less important (e.g. users just browsing a web site), while keeping alive important persistent connections (e.g. a user has placed an item into a shopping bag and

thus the likelihood of a purchase is high). The importance of connections is determined by the cookies found in the HTTP header. Our experiments have shown that this approach can prevent uncontrollable server overload, provide service differentiation under high server load and has low overhead. The result of this work was paper [3]. This paper was one of only six papers accepted for a workshop that will be held in conjunction with SIGMETRICS 2001. A slightly different version of this paper will be presented to the ARTES network at Real-Time in Sweden 2001 [2]. Due to the encouraging comments by the reviewers we are discussing an extension of this work.

Currently, we are building an adaptation entity that dynamically adapts the rate control policies to changing workload conditions. We explore an approach that uses control theory, in particular a PI-controller. Our very early results look promising.

As stated in the project plan we plan to collaborate with Lars Albertsson using his extension of SIMICS to explore the behavior of a server system under high load in more detail. Due to the implementation intensity of this task, we are currently looking for a thesis worker for the bulk of the implementation.

Furthermore, I presented paper [5] describing a first version of our overload architecture at a European workshop/summer school in Holland.

In summary, Thiemo has been able to follow the project plan with only small discrepancies. As stated in the project plan, we expect that Thiemo will present his doctoral dissertation in spring 2002.

## 2 Industrial Participation

As mentioned above, Thiemo has visited IBM TJ Watson Research Center. Together with researchers from IBM paper [4] has been written. According to the last information from my colleagues at IBM (May 18, 2001) there are proposals to include this work in IBM's AIX operating system. Thiemo has also discussed the other research items with researchers both at IBM TJ Watson and Ericsson SARC, at SARC in particular with our industrial contact person, Dr. Lars Björnfot.

## References

- [1] Thiemo Voigt and Per Gunningberg. Dealing with memory-intensive web requests. Technical report, Uppsala University, 2001.
- [2] Thiemo Voigt and Per Gunningberg. Handling persistent connections in overloaded web servers. In *Real-Time in Sweden 2001*, August 2001.

- [3] Thiemo Voigt and Per Gunningberg. Kernel-based control of persistent web server connections. In *Performance and architecture of web servers, PAWS 2001*, June 2001.
- [4] Thiemo Voigt, Renu Tewari, Douglas Freimuth, and Ashish Mehra. Kernel mechanisms for service differentiation in overloaded web servers. In *Usenix Annual Technical Conference 2001*, June 2001.
- [5] Thiemo Voigt, Renu Tewari, and Ashish Mehra. In-kernel mechanisms for adaptive control of overloaded web servers. In *Sixth Eunice Open European Summer School*, Twente, Holland, September 2000.

Paper [3] will also be published in a special issue of the ACM Sigmetrics publication "Performance Evaluation Reviews".