

# Designing High Performance, Reliable, and Energy-Efficient Networked Computing Systems for the Future



**Chen-Zhong Xu inside the Earth Simulator, the then fastest computer.**

As computer systems become more and more networked and complex, new foundations are needed for understanding and controlling their integral properties. Cheng-Zhong Xu, a professor of Electrical and Computer Engineering and the Director of Center for Networked Computing Systems (CNC), is dedicated to the investigation, establishment, and experimental evaluation of new theoretical foundations and system artifacts to significantly improve the systems performance, reliability, security, and even energy efficiency. The systems of particular interest to his research team include distributed systems and the Internet, parallel computers, and networked embedded systems. The novelty of their solutions led to four funding awards by the National Science Foundation in the past two years alone in support of their research and development activities in related areas.

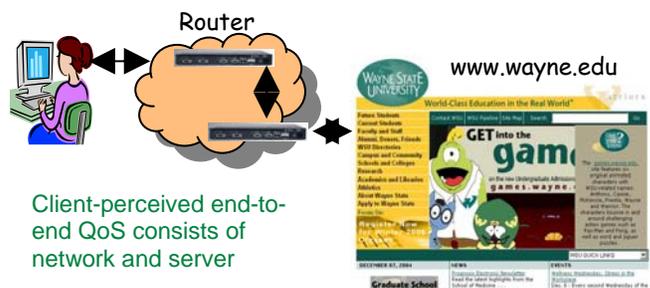
## Quality of Service Assurance in Internet Servers

In the past decade, we all have witnessed an explosive growth of Internet services. Today's Internet services are largely run in a best-effort model, without providing service quality assurance. In this model, clients can be serviced to their satisfaction when a server is lightly loaded. A critical problem arises when the server becomes heavily loaded and has insufficient resources to meet the needs of all clients. The stress condition

may be due to flash crowd of legitimate access. In particular, during big, often unforeseeable events such as stock market roller coaster rides and terror attacks, there is a surge of Internet traffic that can quickly saturate the server capacity and inhabit the services. The stress condition may be even caused by attacks mounted by using the sheer volume of compromised machines to mimic the web browsing behavior of normal clients.

“To enhance the service availability and resilience to stress, a server should be able to run in a service quality assurance (QA) model so as to provide different levels of quality to different types of requests when the server becomes heavily loaded,” says Dr. Xu. “For example, in an e-commerce web site, such a QA model would provide heavy buyers with guaranteed service quality to maximize profits. At the same time, the needs of occasional visitors would not be over-compromised because their requests could also turn out to be profitable. It can also help protect servers from flash crowd-like attacks by assuring the quality of legitimate requests and meanwhile downgrading those in suspicion.”

The concept of quality of service, while not new, is being taken to uncharted territory in the work of Dr. Xu and his research team. Their research is particularly applicable to the design of next generation of stress-resilient servers. Its main thrust is to develop application-level autonomic resource management technologies for multi-class quality assurance with respect to client-perceived end-to-end response time in Internet servers. The technologies are fused into a transparent reverse proxy deployed in the front of a web site. For given



service quality requirements for multi-class requests, the proxy will monitor the service quality in real-

time and adjust the amount of resources allocated to different classes adaptively to keep their quality in desired levels. The proxy and the back-end web site form a closed feedback control loop to deal with dynamics and uncertainty of the servers.

For service differentiation, the resource manager first classifies the incoming requests according to rules defined by service providers. For each classified request, admission controller blocks requests selectively to prevent the back-end servers from illegitimate access or being overloaded. It also rejects requests that are most likely to be malicious when flash crowd attacks are detected. Admission control is the first measure of stress resilience. But it has no control over the quality of admitted multi-class requests. It must be complemented by QoS-aware resource allocation to treat high priority requests preferentially and grant other requests degraded quality when the servers become heavily loaded. In the case of flash crowd attacks, it guarantees the quality of legitimate access and slowdowns suspicious ones.

The admitted requests in different classes are served in a work-conserving weighted fair queuing discipline to guarantee fair sharing of the processing rate of a server between different classes. A challenge is to determine appropriate processing rates for different classes dynamically so as to control their quality in a coordinated way and provide service quality guarantees at desired levels. This is achieved by other two key subsystems: online real-time performance monitor and feedback controller for resource allocation.

Dr. Xu and his team recently developed a method to measure client-perceived end-to-end pageview response time of both unsecured and secured Internet services in a non-intrusive manner at the server side. The method was presented in 2006 USENIX Annual Technical Conference.

The quality controller adjusts the processing rate of each class by responding to measured deviations from the desired quality level. The quality monitor, controller, resource manager, and the back-end web site together form a closed control loop in support of autonomic resource management. Preliminary results on PlanetLab

using industry benchmarks are very encouraging. Dr. Xu's latest book, "Scalable and Secure Internet Services and Architecture" (Chapman & Hall/CRC, 2005), presented a comprehensive treatment of the issues.

### **High Performance and Reliable Servers**

Due to the unprecedented scale of the Internet, Internet services have become an important class of driving applications for parallel computing. Replicating a service in multiple hosts to form a server farm provides a scalable solution to keeping up with ever-increasing application load. A crucial issue to parallel efficiency is task scheduling that distributes tasks between the servers for the objective of load balancing and task locality preserving. Xu and his team developed stochastic scheduling and optimization algorithms for the purpose. Their algorithms are completely analyzable in theory. Their beauty is more than mere tractability. They have since applied the techniques with success in a number of real-world problems, such as computational fluid dynamics, molecular dynamics, and radiation treatment planning. Dr. Xu's widely cited book, "Load Balancing in Parallel Computers: Theory and Practice" (Kluwer Academic/Springer Verlag, 1995), was dedicated to this issue. Their research on this topic has received continual support from NSF for more than a decade.

Networked computing systems research focused on performance and speed in the past. Nowadays, people are concerned more about systems reliability as the systems continue to grow in scale and in the complexity of their components and interactions. It is because component failures in these systems become norm instead of exception. For example, a recent IBM internal study on the ASC White machine with 512 nodes installed in LLNL showed that the mean time to failure of a node is about 160 days. It implies there are 3 to 4 nodal failures every day. If the same failure model is applied to IBM's latest BlueGene machines, there would be more than a hundred nodal failures per hour! As a result, a long running job on a large number of nodes may find it difficult to make progress due to frequent failures. It was because of the system availability concern, BlueGene/L in LLNL had to disable L1 cache in each node when jobs larger than 4 hours was

running because the cache was found prone to failure.

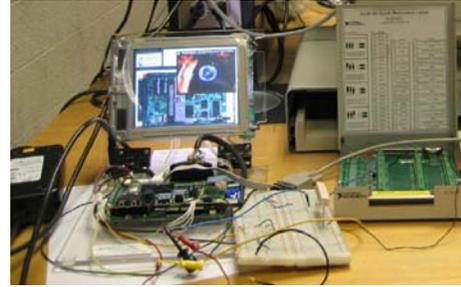
Dr. Xu is leading an NSF-funded 3-year multi-disciplinary project, in collaboration with Prof. L. Wang in Electrical Engineering and Prof. George Yin in Mathematics, to stochastic modeling, optimization, learning, and scheduling technologies for adaptive, highly reliable, and self-manageable systems.

Check-pointing is a conventional approach for fault tolerance. Because of its high overhead, frequent periodic check-pointing often prove counter-effective. Instead, failure prediction based on the analysis of past failure data could lead to a deep understanding of emergent, system-wide phenomena and self-managing resource burdens. Dr. Xu and his team recently developed formal models to quantify the temporal correlation among failures in different timescales and the spatial correlation of failures in different scopes. The co-existence of both spatial and temporal correlation may lead to failure propagation from node to node at different times. The models facilitate the study of failure propagation and its impact on prediction accuracy. The model accuracy was verified by using the systems logs of LANL, as well as year-long online testing data at Wayne State University's grid environment. The results have been accepted for presentation in the IEEE Supercomputing Conference in Fall 2007.

### **Energy-Efficient Resource Management**

Energy consumption has long been a crucial issue in battery-powered embedded systems. It has become an increasingly important concern in servers because of their mass deployment nowadays. Power consumption now accounts for nearly 40% of a data center's operating budget. In December 2006, Congress passed an Eshoo Computer Server Energy Efficiency Bill (H. R. 5646), "calling for additional research to reduce energy costs and electricity consumption by computer servers and data centers."

Dynamic voltage and frequency scaling (DVFS) technology lays a foundation for CPU power saving. However, down-scaling voltage and frequency does not necessarily mean system-wide



energy saving, when other components like memory and I/O are considered.

Dr. Xu and his team developed a rigorous model for the first time for system-wide energy optimization. Although the problem is proven NP-hard for a given set of real-time tasks, they developed pseudo-polynomial algorithms for the derivation of their critical CPU speeds. The results are to be published in ACM Transactions on Embedded Computing Systems.

### **Students Training and Experience**

Dr. Xu's research has attracted a group of high caliber students and research associates to work with him over the years. All of his PhD students and research associates graduated from his group joined tenure-track faculty in academy or research staff in leading IT industry. Most recent ones include Xiaobo Zhou at University of Colorado at Colorado Spring, Haiying Shen at University of Arkansas, Song Fu at New Mexico Tech., Jianbin Wei at Yahoo! Technical, and Xiliang Zhong at Microsoft.

Dr. Xu taught courses at both graduate and undergraduate levels in the areas of networking and the Internet, distributed and parallel systems, and computer architecture. He conveyed current knowledge and design principles effectively while at the same time preparing students who were considering further study and practice in the field to learn engineering trade-offs. His classes earned a reputation of challenging but rewarding in the department. Because of my innovations and excellence in teaching, he was awarded the university's most prestigious "*President's Award for Excellence in Teaching*" in 2003.