Modelling Website Infrastructure using B-Node Theory G. kohli, D. Veal & S.P. Maj School of Computer and Information Science Edith Cowan University Perth, Western Australia {g.kohli,d.veal,p.maj}@ecu.edu.au

#### Introduction:

A large volume of business is conducted via the Internet (Schneider and Perry 2000). However, this has resulted in increased transaction delay times as systems and computer networks become overloaded (Devlin, Gray. J et al. 1999). Surveys and studies indicate that slow downloading time is the most often cited reason that an online customer leaves a site and searches for another vendor's site (Bakos 1998). According to Shklar: "Sites have been concentrating on the right content". Now, more of them specially e-commerce sites realize that performance is crucial in attracting and retaining online customers."(Shklar 1998). The performance of an Internet site is dependant not only upon the behavior of end users using that site but also the performance of the technologies employed. Currently there are a number of different models for defining e-business web sites performance, such as the 'business', 'functional', 'customer', and 'resource' models (Menasce, Virgilio et al. 2000). The Customer Behavior Modeling Graph (Inverardi & Wolf, 1995) (Union, 1996) and Client /Server Interaction Diagrams (CSIDs) (Stohr and Kim 1998) are techniques that can be used to capture the navigation patterns of customers during site visits and hence obtain quantitative information on workloads. However, although these models attempt to predict user behavior, they do not provide information about the actual load on the systems running on such sites. In the final analysis user load must be translated to hardware requirements thereby allowing performance bottlenecks to be identified. However, such problems associated with infrastructure design are non-trivial. According to Fenik "Being able to manage hit storms on commerce sites requires more then just buying more plumbing." (Fenik 1998). In order to predict the workload characteristics of e-commerce sites, effective modelling needs to be undertaken to determine key bottlenecks within the system. It is therefore necessary to investigate different types of models which could be used to model the infrastructure of e-commerce web sites.

## **Model Requirements:**

Models are used not only as a means of communication and controlling detail but may also form the basis of a conceptual understanding of a system. According to Cooling there are two main types of diagram: high level and low level (Cooling 1991; Booch, James Rumbaugh et al. 1999). High-level diagrams are task oriented and display overall system structure and major sub-units. Such diagrams describe the overall function of both the design and interactions between both the sub-systems and the environment. The main focus is upon finding answers to the question what does the system do? According to Cooling, "Good high-level diagrams are simple and clear, bringing out the essential major features of a system" (Cooling 1991). By contrast, low-level diagrams are solution oriented and must be able to handle considerable detail. The main emphasis is 'how does the system work'. However, all models should have the following characteristics: diagrammatic, self-documenting, easy to use, control detail and allow hierarchical top down decomposition. For example, the Data Flow Diagram (DFD) (Hawryszkiewycz 2001) model enables a complex system to be partitioned (or structured) into independent units of an amenable size so that the entire system can be more easily understood. It is possible, therefore, to examine a system in overview and with increasing levels of detail, whilst maintaining links and interfaces between the different levels. DFD's are not only simple, but also graphical; hence they serve not only as documentation but also as a communication tool (Pressman 1992). DFD's are therefore a top-down diagrammatic representation of information flow within a system, and are a means of defining the boundaries and scope of the system being represented, checking the completeness of the analysis and providing the basis for program specifications. This technique is relatively simple to use, yet powerful enough to control complexity during the analysis and design of both small and large systems. It is recognized that communication with end users is especially important as this helps to validate a model for correctness. There are various high level models that are used to evaluate web site performance.

## Web Performance

There are various well established methods for evaluating Internet site performance (Menasce, Almendia et al. 1994). The Customer Behaviour Modeling Graph (CBMG) can be used to measure aggregate metrics for web sites (Menasce and Almendia 1999). Using this modeling technique it is possible to obtain a wide variety of different performance metrics that include: hits/s, unique visitors etc. However, when using this technique it is not possible to relate these metrics to hardware specifications. The Client /Server Interaction Diagrams (CSIDs) (Stohr and Kim 1998) can be used to capture the navigation patterns of customers during site visits and hence obtain quantitative information on workloads. However, it does not provide any insight into how the workload will affect the underlying infrastructure.

Furthermore, the World Wide Web (WWW) has some unique characteristics that distinguish it from traditional systems (Mogul 1995; Almedia, Bestravos et al. 1996; Arlitt and Williamson 1996; Almeida, Virgilio Almeida et al. 1997). Firstly, the number of WWW clients is in the range of tens of millions and rising. Secondly, the randomness associated with the way users visit pages makes the problem of workload forecasting and capacity planning difficult. (Menasce, Almendia et al. 1994).

Benchmarks are the standard metrics used in defining the scalability and performance of a given piece of hardware or software. For example the Adaptive Computing System (Sanjaya, Chirag et al. 2000) is a collection of benchmarks that focus upon specific characteristics from the start of a computation until its completion. Benchmarks evaluate the ability of a configurable computing infrastructure to perform a variety of different functions. SPECWEB and TPC-C (Smith 2000) are notable benchmarks in the e-business environment. These benchmarks come close to representing the complex environment of an e-business workload. Benchmark programs are used for evaluating computer systems. Different end user applications have very different execution characteristics; hence there exists a wide range of benchmark programs. The four main categories are: science and engineering (examples include; Whetstones, Dhrystones, Livermore loops, NAS kernels, LINPACK, PERFECT club, SPEC CPU), Transaction Processing (for example; TPC-A, TCP-B, TPC-C), server and networks (examples include; SPS/LASSIS, SPEC web) and general use (examples include; AIM Suite III, SYSmark, Ziff-Davis PC

Benchmark). However none of these benchmarks are directly relevant to Ecommerce web transactions. The Transaction Processing Council introduced the TPC-W that simulates the workload activities of a retail store Web sites (Smith 2000). In the TPC-W standard the products are books and the user is emulated via a remote Browse that simulates the same HTTP traffic as would be seen by a real customer using the browser. E-business sites have transient saturation so it is hard to use these benchmarks to get the correct idea about the actual load generated on the web servers. Benchmarks currently in use fail to measure the web performance characteristics, whilst others may be incorrectly interpreted (Humphrey 1990; Lilja 2000). According to Skadron, "Research cannot pursue futuristic investigation when they are limited to systems for which no benchmark programs are available. The current short coming in computer systems evaluation could ultimately even obstruct the innovation that is driving the information technology revolution" (Skadron, Martonosi et al. 2003). Furthermore the basic problem still remains. Using benchmarks it is not possible directly relate the technical specification to the metrics used in the service level agreements. The difficulties of developing effective models for large networks are becoming greater as noted by Clark, "As networks grow to connect millions of nodes, and as these nodes all communicate in unpredictable patterns, the resulting behaviour becomes very difficult to model or predicts" (Clarke and Pasquale 1996). The question was then asked, what methods do IT web site managers use to design and manage the performance of an Internet web site?

## **Commercial Practices:**

A questionnaire was distributed to several small to medium size companies in Western Australia. The results indicated that infrastructure requirements were typically based upon past experience and also purchasing the highest performance equipment within budget constraints (Maj and Kohli 2002). Alternatively companies outsourced this problem to vendors. These approaches are arguably entirely unsatisfactory as they relegate IT systems analysis to conventional 'wisdom' and mythology (Maj and Kohli 2002). None of the companies analysed employed any techniques for modelling infrastructure performance. The scope of the above survey is currently being extended both within Australia and internationally. From the data gathered to date it can further be concluded that most companies where not aware of an effective model that could applied to effectively model e-commerce website workloads. Hence there is current need for such a new model in this area.

## Modelling infrastructure using B-Nodes:

Computer and network equipment is complex. Furthermore they use a wide range of heterogeneous technologies with different performance metrics. By example hard disc drive performance is often quoted in rpm; electronic memory performance is quoted in nanoseconds; microprocessor performance is quite in MHz etc. This results in two problems. Firstly the performance of a web site (server with switches, hubs etc) depends upon the speed of the slowest device. It is not possible, using these metrics, to easily determine the relative performance of each device. Is 10ns electronic memory faster or slower than a 1GHz microprocessor? Secondly it is difficult to relate the technical performance metrics to user requirements defined in the Service Level Agreement. Can a hard disc drive operating at 5,000 rpm deliver 100 web pages per

minute? The B-Node model was proposed to address these problems (Maj and Veal 2000). The B-Node model:

- Can be used to model a wide range of computer and network technology equipment
- Is diagrammatic, self documenting and easy to use
- Uses recursive decomposition, hence can be used to model both small systems (e.g. a server) or a larger system (e.g. an Intranet)
- Uses a common performance metric (Mbytes/s). Hence the performance of heterogeneous technologies can easily be compared
- Uses a common fundamental unit (Mbytes/s) allowing other units to be derived. Hence it is possible to define the performance of a wide range of different technologies using, for example, a common, derived metric such as web pages per second.

The B-Node model has been used to model an E-commerce server and hence identify hardware bottlenecks. It has also been used to evaluate the performance of different E-commerce serves (Web server, payment server etc) (Maj, D.Veal et al. 2001). However, the use of bandwidth as a sole indicator of performance may be

problematic. According to McComas notes there are problems due to bandwidth and latency (McComas 2001), as does Buzen and Shum, (Buzen and Shum 1996). This point is also made with respect to network technology by Openhiemer

"It is possible to improve throughput such that more data per second is transmitted, but not increase goodput, because the extra data transmitted is overhead or retransmissions ...more data is transmitted per time, but the user sees worse performance. ...most end users are concerned about throughput rate for applications. Marketing materials from some networking vendors refer to application-layer throughput as 'goodput'. Calling it goodput sheds light on the fact that it is a measurement of good and relevant applicationlayer data transmitted per unit time" (Oppenheimer 2001).

In effect, it is possible to have higher bandwidth but it is not being used effectively to transfer data. Hennessy also notes the: "... *pitfall of using bandwidth as the only measure of network performance.* ... this may be true for some applications such as video, where there is little interaction between the sender and the receiver, but for many applications such as NFS, are of a request-response nature, and so for every large massage there must be one or more small messages ... latency is as important as bandwidth" (Hennessy and Patterson 1996). In spite of this the B-Node model has many potential advantages and it may be possible to address the latency issue.

## Using B-Nodes to measure Network Technology performance:

A wide range of different files were transferred between two PCs using a simple cross over cable using FTP. This represented the base line performance. Then a range of different networking technologies were introduced and the performance measured. In order to address the problem of latency the authors have subsumed the effects of latency under a definition of bandwidth. Namely bandwidth = the size of the file in mbytes / total time to send that file:

# $\mathbf{B} = \mathbf{L}_1 / \mathbf{T}_{\mathrm{T}}$

Where  $T_T = t_1 + t_L$  B = (bits passed)/(time taken to pass those bits)  $L_1 = Length of the files in Mbytes.$   $t_1 = Time required to transfer the file$  $t_L = Latency measured in (msec)$ 

Table 1 shows a summary of various devices with respect to bandwidth using File Transfer Protocol (FTP):

Technology	Bandwidth (Mbytes/s)	
PC crossover cable PC	11.5	
PC switch PC	11.5	
PC router PC	7.5	

Table 1 B-Node performance figures (Mbytes per second)

The crossover cable between two PCs can be modelled as a B-Node with a performance of 11.5Mbytes/sec. A switch can be modelled as a B-Node with a performance also of 11.5Mbytes/sec. In effect a switch works at 'wire speed' and has no measurable affect on performance. A router modelled as a B-Node gives a performance of 7.5 Mbytes/sec. The use of common fundamental units allows two different technologies (layer 2 switches and layer 3 routers) to be compared. Furthermore, common derived units can be used. Assuming the messages in a web transaction are 10 Kbytes each and the load is 1000 per second. It can then we concluded from Table 2 by introducing a Router will create a bottleneck in the system as the utilisation is more then 100%. Additionally it is possible to identify, using meaningful metrics the relative performance of each technology which can then help network designer to better design web sites infrastructure.

Technology	Bandwidth (Mbytes/s)	Transaction size (Kbytes)	Load	Utilization
PC crossover cable PC	11.5	10	1000	86%
PC switch PC	11.5	10	1000	86%
PC router PC	7.5	10	1000	133%

Table 2 B-Node performance figures (Transactions per second)

The authors are further developing the experiment by taking into account different protocol like HTTP and HTTPS and the use of Access control list (ACL).

## **Conclusion:**

The performance of Network application affects the productivity in many areas: ecommerce, a model base approach provides a good foundation for developing solutions to these problems. The B-Node model is simple, diagrammatic and selfdocumenting modelling technique, it use common fundamental units which can help the network designer to identify the key bottlenecks within the system. The B-Node model is undergoing development and testing in an attempt to model infrastructure from the bottom up to enable top down conceptual understanding of workload characteristics.

## **References:**

- Almedia, V., A. Bestravos, et al. (1996). <u>Characterization reference locality in the</u> <u>WWW</u>. IEEE-ACM PDIS.
- Almeida, J. M., Virgilio Almeida, et al. (1997). "WebMonitor: a Tool for Measuring World-Wide Web Server Performance."
- Arlitt, M. and C. Williamson (1996). <u>Web server Workload Characterization</u>. ACM SIGMETRICS Conference of Measurement of Computer Systems, Philadelphia, PA.
- Bakos, Y. (1998). "The Emerging Role of Electronic Market Places on the Internet." <u>ACM</u> **41**(8).
- Booch, G., James Rumbaugh, et al. (1999). <u>The Unified Modeling Langauge User</u> <u>Guide</u>, Addision Wesley longman, Inc.
- Buzen, J. P. and A. N. Shum (1996). <u>Beyond Bandwidth: Mainframe Style Capacity</u> <u>Planning for Networks and Windows NT.</u> Proceedings of the 1996 Computer Measurement Group Conference., Olando Florida USA.
- Clarke, D. and J. Pasquale (1996). "Strategic Directions in Networks and Telecommunications." <u>Computing Surveys</u> **28**(4): 679 690.
- Cooling, J. E. (1991). <u>Software Design for Real-Time Systems</u>. Cornwall, Chapman and Hall.
- Devlin, B., Gray. J, et al. (1999). Scalability Terminology: Farms, Clones, Partitions and Packs: RACS and RAPS, Microsoft Research.
- Fenik, H. (1998). Zona Research. Lan Times.
- Hawryszkiewycz, I. (2001). Introduction to Systems Analysis & Design. Malaysia, Prentice Hall.
- Hennessy, J. L. and D. A. Patterson (1996). <u>Computer Architecture A Quantitative</u> <u>Approach</u>. San Francisco, Morgan Kaufmann.
- Humphrey, W. S. (1990). <u>Managing the software process</u>. MA, Addison-Wesley, Reading.
- Lilja, D. J. (2000). <u>Measuring Computer Performance A Practitioner's Guide</u>. Cambridge, Cambridge University Press.
- Maj, S. P., D.Veal, et al. (2001). <u>A New Method for controlling Information</u> <u>Technology Hardware Complexity</u>. American Conference on Information Systems, Boston, USA.
- Maj, S. P. and G. Kohli (2002). <u>B-Node: A bridge between the business model</u>, <u>information model and infrastructure model</u>. Seventh Caise/IFIP-WG8.1 Internation Workshop on Evaluation of Modeling Methods in System Analysis and Design, Toronto.
- Maj, S. P. and D. Veal (2000). "Computer Technology Curriculum A New Paradigm for a New Century." <u>Journal of Research and Practice in Information</u> <u>Technology</u> **32**(August/September): 200-214.
- McComas, B. (2001). Can the Pentium 4 Recover?, InQuest Market Research. 2001.
- Menasce, D. A. and V. A. F. Almendia (1999). "A metholody for workload characterization for e-commerce server." <u>ACM conference in Electronic Commerce, Denver, CO, ACM</u>.
- Menasce, D. A., V. A. F. Almendia, et al. (1994). <u>Capacity Planning and Performance</u> <u>Modeling: From Mainframes to Client-Server Systems</u>. Upper Saddle Rive, NJ, Prentice Hall.
- Menasce, D. A., V. A. F. Almendia, et al. (1994). <u>Capacity Planning and Performance</u> <u>Modeling</u>. N J, Prentice Hall.

- Menasce, D. A., A. F. Virgilio, et al. (2000). <u>Scaling for E-Business: Technologies</u>, <u>Models, Performance, and Capacity Planning</u>.
- Mogul, J. C. (1995). <u>Operating system support for busy Internet servers</u>. Proceedings of the fifth workshop on Hot Topics in Operating Systems.
- Oppenheimer, P. (2001). <u>Top Down Network Design: A systems analysis approach to</u> <u>enterprise network design</u>. Indianapolis IN, Cisco Press.
- Pressman, R. S. (1992). <u>Software Engineering A practitionar's Approach</u>. Singapore, McGRAW-HILL International Editions.
- Sanjaya, K., N. Chirag, et al. (2000). "A Benchmarking suite for evaluating configurable computing systems- Status, Reflections, and Future Directions." <u>Eighth International Symposium on Field- Programmable Gate Arrays,</u> <u>Monterey, California</u>.

Schneider, G. P. and J. T. Perry (2000). <u>Electronic Commerce</u>, Course Technology. Shklar, G. (1998). The New York Times. New York.

Skadron, K., M. Martonosi, et al. (2003). Challenges in Computer Architecture Evaluation. <u>Computer</u>. **36:** 30-36.

Smith, W. D. (2000). TPC-W: Benchmarking An Ecommerce Solution, TPC.

Stohr, E. A. and Y. Kim (1998). <u>A Model for Performance Evaluation of Interactive Systems</u>. 31st. Annual Hawaii Internternational Conference on System Sciences.