

# Scalable and Cost-Effective Framework for Continuous Media-On-Demand

**Dang Nam Chi Nguyen**

A Thesis presented for the degree of  
Doctor of Philosophy



Department of Computer Systems

Faculty of Information Technology

University of Technology Sydney

Australia

2006

# Certificate of Authorship

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

**Candidate Name: Dang Nam Chi Nguyen.**

**Signature of Candidate.**

# Acknowledgements

There are many people I would like to thank for their help and support over the years. Firstly, I would like to thank my supervisor, Prof. Doan Hoang, for his thoughtful input, and constructive feedback throughout. Without his wise guidance, help and encouragement, it is quite certain none of this would have been possible. His advice has been truly wonderful and I am forever in his debt.

I would also like to thank my former supervisor, Prof. Antonios Symvonis who had helped me greatly in the early parts of my journey. His continuous support and advice, despite the great distance separating us, is very much appreciated.

I would like to thank my family - my mother, aunts and uncles - for their faith and belief in me. I would especially like to thank my grandparents who have always loved, cared and persevered with me. Without them I would not have been in a position to undertake this endeavor and it is for them that I dedicate this thesis.

Finally, I would also like to thank all my friends, Bushar, Ming, Hanh, and Joe, who, beside putting up with me in the same room, have proved to be wonderful company as well collaboration partners. Your friendships have got me through many hard times and I will miss getting my doses of second-hand smoke from our daily coffee breaks!

# Abstract

This dissertation was motivated by the exponential growth in bandwidth capacity of the Internet, coupled with the immense growth of broadband adoption by the public. This has led to the development of a wide variety of new online services. Chief amongst the emerging applications is the delivery of multimedia contents to the end users via the network on-demand. It is the “on-demand” aspect that has led to problems which, despite the advances in hardware technology and network capacity, have hampered wide scale adoption of multimedia delivery. The focus of this dissertation was to address these problems, namely: scalability, cost-effectiveness, and network quality of service for timely presentation of multimedia contents.

We proposed an architecture, which we referred to as “Delayed-Multicast”, to address the scalability problem. The new architecture introduced buffers within the network to reduce demands on core network bandwidth and server load. A feasibility study of the architecture was conducted through the use of a prototype. It was found that such a system is within reach by demonstrating the prototype using cheap, common-of-the-shelf (COTS) components, and with help of freely available system software such Linux with real-time support.

The introduction of buffers within the network led to the requirement of how to minimize buffer space. We developed an optimal algorithm for allocating buffer space in a single level caching layout (i.e. only one buffer in the transmission path from the server to the end user).

For the case of multi-levels network caching, we thoroughly examined different optimization problems from an algorithmic perspective. These problems included how to minimize total system memory, and minimize the maximum memory used per node. We proved that determining the optimal buffer allocation in many of these

---

cases is an NP-complete problem. Consequently, we developed heuristics to handle multi-level caching and showed through simulations that the heuristics greatly help in minimizing buffer space and network bandwidth requirement.

An important aspect of the heuristics was how to handle the case when the arrival times of client requests were not known *a priori*. For these “online” problems we also proposed heuristics that can significantly reduce overall system resource requirements. If the cost of buffer space was also taken into account along with the cost of network bandwidth, a different optimization problem was how to minimize the total system cost. Here, we also proposed heuristics, which in simulations show that the total system cost can be significantly reduced.

Besides the problems associated with resource allocation, in terms of buffer space and bandwidth, we also examined the problem of how to provision the necessary network quality of service on-demand. Most current networks rely on best-effort delivery which is ill suited for the delivery of multimedia traffic. We proposed a solution which relied on the use of a programmable network plane, that is present in many current routers, to dynamically alter the priority of flows within the network in real-time. We also demonstrated the effectiveness of the flow prioritization on an actual Nortel router.

Finally, we examined the problem of how to admit and achieve fair bandwidth allocation for the end-users within a Differentiated Service (DiffServ) network. DiffServ is an IETF standard that aims to provide a “better than best-effort” network in a scalable manner, and is used widely, especially within the same autonomous domain for prioritization different classes of traffic. However, there are open problems on how to provide fair bandwidth allocation amongst competing flows. We proposed an edge-aware resource discovery loop, which as the name suggests, sent packets to gather information about the internal states of the core network. With this information, we proposed a price-based admission control algorithm for use within the DiffServ network that would allow fair admission, effective congestion control, and fair bandwidth allocation amongst different traffic flows.