

2006

Network design model to reduce latency on a multi-point frame relay network

Gretchen Trump
Lehigh University

Follow this and additional works at: <http://preserve.lehigh.edu/etd>

Recommended Citation

Trump, Gretchen, "Network design model to reduce latency on a multi-point frame relay network" (2006). *Theses and Dissertations*. Paper 920.

This Thesis is brought to you for free and open access by Lehigh Preserve. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of Lehigh Preserve. For more information, please contact preserve@lehigh.edu.

Trump, Gretchen
M.

Network Design
Model to Reduce
Latency on a Multi-
point Frame Relay
Network

January 2006

**NETWORK DESIGN MODEL TO REDUCE LATENCY
ON A MULTI-POINT FRAME RELAY NETWORK**

by

Gretchen M. Trump

A Thesis

Presented to the Graduate and Research Committee

of Lehigh University

in Candidacy for the Degree of

Master of Science

in

Information and Systems Engineering

Lehigh University

December 2005

This thesis is accepted and approved in partial fulfillment of the requirements for the Master of Science.

DECEMBER 9, 2005
Date

Thesis Advisor

Co-Advisor

~~Chairperson~~ of Department

Acknowledgements

I would first and foremost like to thank Dr. Rosemary Berger and Dr. Joseph Hartman for all of their assistance in this endeavour. Many thanks to Jeff Priester for serving as my industry advisor for this project and helping to define the problem. Also without the support and guidance from my friends and colleagues Bill Arey, Ruth Campolongo, John Fagan, Eric Kirstein, Harry Kubat, Steve Lilley, John Schadler, Bill Stevenson, and Bill Troxell I would have not been able to complete this project in the same manner.

Contents

Acknowledgements	iii
List of Tables	v
List of Figures	vi
Abstract	1
1 Introduction	2
1.1 Network Design Case Study	3
1.2 Improving Latency	4
2 Literature Review	6
3 Formulations	7
3.1 Parameters and Notation	7
3.2 Nonlinear Formulation	8
3.3 Linear Formulation	9
4 Computational Experiments	11
4.1 Data Collection	11
4.2 Base Model Results	12
4.3 Sensitivity Analysis	12
4.3.1 Change in Maximum Latency Value	12
4.3.2 Accounting for Percentage of Web Traffic	15
4.3.3 Adjusting the Circuit Costs	17
4.3.4 Restricting the Model to Use Larger Circuit Sizes	20
5 Conclusions	21
References	22
Appendix	23
Biography	28

List of Tables

1	Solution for Base Model.	12
2	Objective function values when altering maximum latency α_1	13
3	Configuration when $\alpha_i = 0.1$ for $j = 0$	13
4	Configuration when $\alpha_i = 0.15$	14
5	Configuration when $\alpha_i = 0.25$	14
6	Configuration when $\alpha_i = 0.3$	15
7	Configuration when $\gamma = 0.25$	16
8	Configuration when $\gamma = 0.5$	16
9	Configuration when $\gamma = 0.75$	17
10	Sensitivity of objective function value to circuit costs.	18
11	Configuration when $a = 100$	18
12	Configuration when $a = 120$	19
13	Configuration when $a = 160$	19
14	Configuration when $a = 180$	19
15	Base Model Solution with Circuit Size Restriction.	20
16	Original Data for Each Bandwidth Type.	23
17	Original Data Parameters.	23
18	Original Data for Locations 1 through 35.	24
19	Original Data for Locations 36 through 70.	25
20	Original Data for Locations 71 through 105.	26
21	Original Data for Locations 106 through 131.	27

List of Figures

1	Current Global Network Configuration.	4
---	-----------------------------------------------	---

Abstract

This thesis discusses a specific network design problem and solution approach to reduce latency in a global data network. The network can be described as a frame relay network with a centralized hub site and remote sites that connect directly to the hub. The work is supported by a case study in which a company would like to provide connectivity between the hub and the remote sites at minimum cost. Content engines offer one way in which to achieve this goal. The research investigates the feasibility, cost savings, and effect on network latency given the optimal location of content engines at remote sites.

1 Introduction

Efficient network design has become increasingly important over the last decade. The Internet boom has demanded faster, bigger, and more dependable networks, especially for large enterprises. Many companies rely on multiple computer applications to run all aspects of their business and, as a result, reliable, responsive computing environments are critical. Because of limited financial resources, companies need to provide these network services at minimum cost.

This thesis examines the network design problem for a company that operates at hundreds of locations throughout the world. Each location is connected to a hub site, most often via a leased frame relay connection. In the current configuration, each site has an assigned bandwidth (committed information rate) based upon the number of users and the amount of network traffic at the site. At some sites, users are experiencing slow response times. The company is investigating methods to improve the service at these sites. However, any improvements must be made in a way that minimizes the cost of providing network service subject to providing a reasonable overall customer experience.

The overall customer experience is measured in terms of latency, which is the round trip time that it takes a packet to travel between a remote site and the hub. For the user, this provides some measure of response time for a web page request or sending of an e-mail, as these tasks are generally decomposed into multiple packets for subsequent transmission. Latency is a common measure of network performance, and it is a function mainly of distance, bandwidth, and traffic arrival rate. This research investigates the potential impact of changing the bandwidth assigned to a remote site as well as the introduction of equipment that may potentially reduce the traffic arrival rate. Optimization models are developed to simultaneously determine the proper circuit bandwidth for each site and the best equipment locations in order to minimize network costs subject to a maximum latency limit. We refer to this problem as the Network Design Model to Reduce Latency (NDMRL).

This thesis is organized as follows. The details of the company's network design and background information on the proposed approaches for improving latency are included in the remainder of Section 1. The relevant literature is reviewed in Section 2. The formulation of the optimization problem is presented in Section 3. A description of the computational experiments and the results are provided in Section 4. The conclusions of this study are presented in Section 5 along with ideas for future research.

1.1 Network Design Case Study

The company's global network is built around a core of three hub sites. The world headquarters and the corporate network hub are located in the United States; the other two major hubs are located in Europe and Asia. The corporate hub in the United States is also the main hub for all sites in North America. The bulk of the network data resides at the corporate hub, but both the European hub and the Asian hub have similar core infrastructures that house servers for critical applications, such as email, DNS (domain name services), software distribution, and domain controllers. This study focuses on the North America network but can be extended easily to the other two hubs and other telecommunications applications. The study considers 141 remote sites in the the North America network that have dedicated point-to-point frame relay connections to the hub.

The following networking terms will be used throughout the thesis:

Wide Area Network (WAN): a computer network that spans a large geographical area

Bandwidth (or link capacity): the maximum amount of data that can be transmitted across a WAN link, usually measured in the number of bits that can be transmitted in one second

Utilization: the percentage of bandwidth that is actually being used on the WAN link

Latency: the round trip time for one packet to travel from a remote site to the hub

Congestion: the slow response that happens when the amount of traffic exceeds the link capacity (usually measured in terms of the number of packets dropped)

Packet Loss: the dropping of packets that occurs when a link is congested, resulting in lost data

Kbps: kilobits per second, which is equivalent to 1024 bits per second

Mbps: megabits per second, which is equivalent to 1,048,576 bits per second

Figure 1 shows a basic overview of the company's network. There is a main WAN router that resides in the data center at the hub site (Hub eVPN router). The total link capacity to this WAN router from the WAN provider's network is 44 Mbps. Network traffic from the US hub is sent to its destination via the eVPN router. Network traffic originating at a field site exits the remote WAN site via the remote site WAN router and then travels to

its destination, which may be a server on the hub LAN or one of the Internet circuits (also located at the hub). All Internet traffic from the remote sites is routed back to the hub via a dual redundant circuit including one 22 Mbps circuit and one 19 Mbps circuit.

Most congestion in the network occurs either on the remote site WAN leased line or on the Internet circuit. This study will be concerned with the congestion from the remote sites to the hub, since this is where network statistics are monitored and where bandwidth changes can be made.

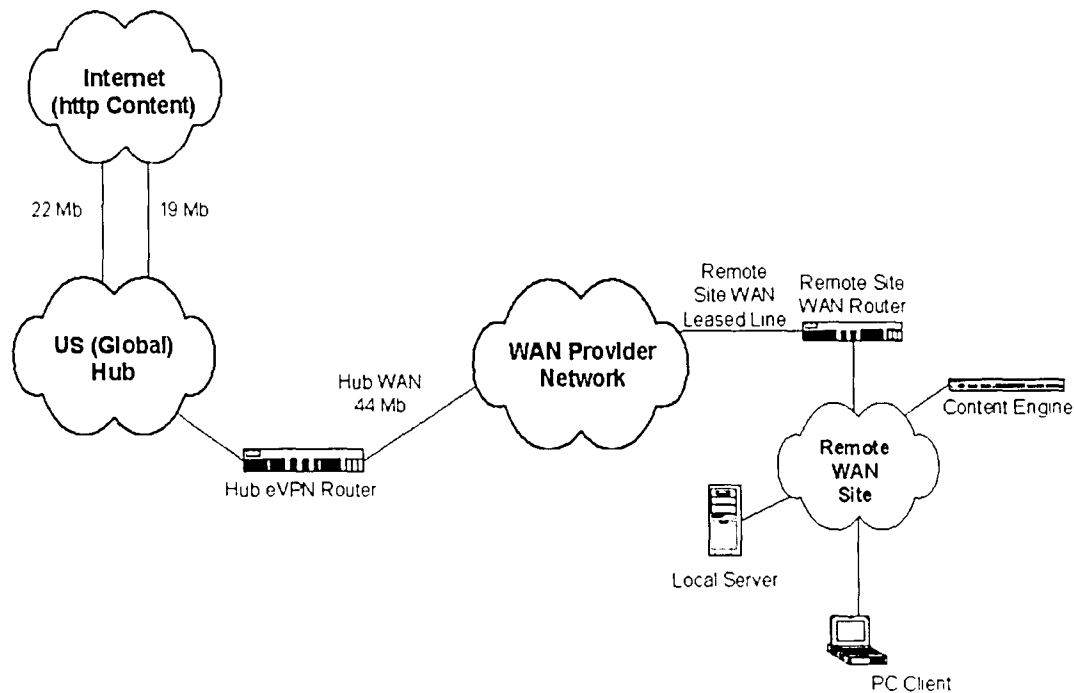


Figure 1: Current Global Network Configuration.

1.2 Improving Latency

In this study, the user's network experience is measured in terms of latency. Latency is widely recognized as an important measure of network performance. It is the single most important factor to address when optimizing application performance and the user experience according to a Gartner Research study by Fabbi [4]. Fabbi argues that while having enough bandwidth at each site is important, the bandwidth is well-utilized only if latency issues are considered. In his study, Fabbi reports that latency is responsible for more than 50 percent of the total application delay of a 128 Kbps network to 95 percent of a 1544 Kbps pipe, and he notes that globalization of networks is greatly affecting latency parameters be-

cause of increased dependency on the network; for example, using enterprise application software such as SAP on a global basis increases the network traffic flow significantly across the LAN and WAN.

Bartlett, Sevcik and Moore [1] discuss the importance of maximizing the utilization at WAN sites and argue that latency and packet loss need to be controlled in order to realize this goal. They point out that a delay slightly larger than 100 milliseconds is noticeable and that a delay exceeding 250 milliseconds would be 'objectionable'. In addition, they indicate that the decision of what size circuit to buy can be best determined by figuring out what the maximum utilization should be (based upon latency and packet loss). This concept will be studied as part of the model in this study.

Latency is a function mainly of distance, bandwidth, and traffic arrival rates. There are a number of approaches to improving network latency, each one directed at one or more of the components. Fabbi [4] notes that the best approach would be to improve applications by improving the actual interface itself. One class of approaches focuses on protocol improvements both for HTML and related protocols as well as for application layer protocols. However, the present study investigates the impact of approaches such as changing the bandwidth, distance, or arrival rate, and/or adding a *content engine* at a site.

A content engine is a network appliance that can be used to cache web traffic and store streaming media for on-demand viewing. Since both web and multimedia applications are bandwidth-intensive, locating a content engine at a remote site offers a way to alleviate some of the load on the circuit connecting the remote site to the hub. Instead of retrieving content from servers at the hub or from the Internet, users may find that the content engine has the cached web content or the stored streaming media that they need, thereby eliminating the need to request the content via the circuit. The hub site has a several 'root' content engines that store content intended for the remote sites and a content management device that 'pushes' the content to the remote sites with content engines. The company in this case study relies heavily on multimedia as a means to provide important employee training and pertinent information remotely. This information is (ideally) located on a content engine.

In the solution being studied, a remote site can house one content engine, which resides on the local LAN. The type of content engine being evaluated is a smaller model specifically designed for being deployed at remote sites. To date, the company has deployed several content engines as part of a pilot project. These content engines are being used at several of the larger remote sites. These sites were selected based on size; no formal selection process was used when they were installed.

The main objective of this research is to develop a framework for investigating the performance impacts of introducing content engines and selecting circuit sizes in the company's frame relay network. In particular, an optimization model is developed that determines where to locate content engines and what size circuit to lease at each site in such a way as to minimize total cost subject to maximum latency constraints. The total cost includes the costs of leasing the frame relay circuits and the costs of adding content engines to the network.

2 Literature Review

The NDMRL problem involves two types of decisions – where to locate equipment (content engines) and how to size network links (circuits). As such, the problem combines aspects of facility location and network design, both of which have been extensively studied in the literature. However, the NDMRL problem has some unique features not well captured by models in the facility location or network design literature.

In facility location models, if a facility is located at a site, then it is assumed that any demand at that site (and any others assigned to that site) is served *completely* by the facility at that node. In the NDMRL model, however, a content engine located at a site serves only *some* of the traffic requests from that site. The site also requires service from the hub node. Therefore, existing facility location models are not appropriate for the NDMRL problem.

If we ignore the location decision, the NDMRL problem can be viewed as a special case of an access network design problem. The key questions that arise in access network design are how to connect the remote sites to a central location and, in some cases, what size to select for the links. The NDMRL problem is a special case in that a star topology already has been specified, so only the link sizing decisions remain. In a star topology, there is a direct connection between each site and the central location.

The latency constraints are a key feature of the NDMRL model. Latency includes time spent in the buffers at the routers and switches; time spent in the buffers depends in part on the service rate, thus on the circuit size. Since circuit size is a decision variable in the model, the time in the queue cannot be calculated independently of the circuit size. To capture the dynamic and stochastic nature of the system, we model the network equipment at each site as a queueing system and we explicitly model the waiting time calculation in the constraints of the problem. Our work is in the spirit of a class of problems known as facility location problems with congestion. Most of the research on facility location problems with

congestion has been developed in the context of locating emergency services. See Berman and Krass [3] for a review of work in this area. More closely related to our research is the paper by Berger and Raghavan [2] in which the authors incorporate a queueing model into a discrete hub location model for an access network design problem. Integrating a queueing model into a facility location or network design model tends to give rise to an optimization model with nonlinear constraints; as a result, most of the models are solved via heuristics. In this research, the latency constraints give rise to a set of nonlinear constraints, however, an equivalent linear formulation can be developed for the model.

3 Formulations

In this section, optimization models for the NDMRL problem are developed. Given a hub and a set of remote sites, the objective of the NDMRL problem is to determine where to locate content engines and what circuit size to lease between each site and the hub so as to minimize total cost subject to maximum latency constraints.

3.1 Parameters and Notation

We define the following notation for the optimization models:

Sets

I = set of remote sites

K = set of circuit sizes

Parameters

λ_i = packet arrival rate at location i without a content engine, $\forall i \in I$

$\hat{\lambda}_i$ = packet arrival rate at location i with a content engine, $\forall i \in I$

c_i = cost of installing a content engine at location i , $\forall i \in I$

$m_{i,k}$ = monthly cost of leasing a circuit of size k at location i , $\forall i \in I, \forall k \in K$

α_i = maximum latency allowed at site i , $\forall i \in I$

d_i = distance (in meters) between location i and the hub, $\forall i \in I$

$E[L]$ = expected packet length

β_k = bandwidth of circuit size k (in bits per second), $\forall k \in K$

μ_k = $\beta_k/E[L]$, maximum service rate of a circuit of size k , $\forall k \in K$

T = total bandwidth of circuit at hub site

s = speed of light in a fiber optic cable (2×10^8 meters per second)

3.2 Nonlinear Formulation

We now define our first formulation, which requires the use of two classes of binary variables.

Decision Variables

$$x_i = \begin{cases} 1 & \text{if a content engine is placed at site } i, \forall i \in I \\ 0 & \text{otherwise} \end{cases}$$

$$z_{ik} = \begin{cases} 1 & \text{if a circuit of size } k \text{ is leased for location } i, \forall i \in I, \forall k \in K \\ 0 & \text{otherwise} \end{cases}$$

Objective Function

$$\text{Minimize } \sum_{i \in I} c_i x_i + \sum_{i \in I} \sum_{k \in K} m_{ik} z_{ik} \quad (1)$$

$$\text{subject to } \sum_{k \in K} z_{ik} = 1 \quad \forall i \in I \quad (2)$$

$$\sum_{i \in I} \sum_{k \in K} \beta_k z_{ik} \leq T \quad (3)$$

$$\frac{d_i}{s} + x_i \cdot \frac{1}{\sum_{k \in K} \mu_k z_{ik} - \hat{\lambda}_i} + \quad (4)$$

$$(1 - x_i) \cdot \frac{1}{\sum_{k \in K} \mu_k z_{ik} - \lambda_i} \leq \alpha_i \quad \forall i \in I \quad (5)$$

$$x_i \in \{0, 1\} \quad \forall i \in I \quad (6)$$

$$z_{ik} \in \{0, 1\} \quad \forall i \in I, \forall k \in K \quad (7)$$

The objective function (1) minimizes the sum of the content engine costs and the circuit lease costs. Constraints (2) require that a circuit be leased for each location i . Constraint (3) restricts the aggregate bandwidth to be at most the total bandwidth available at the hub site. Constraints (4) are the latency constraints for the sites. Constraints (5) and (6) restrict the variables to be binary. There are 987 binary variables and 283 constraints in this model.

Latency includes three components: propagation delay, transmit delay and queuing delay. The propagation delay is the amount of time for one bit to propagate (travel) the length of a circuit. The rate at which a signal propagates depends on the medium in use: for our model, we assume a fiber optic cable with a propagation rate of 2.0×10^8 m/s. To calculate the propagation delay associated with a circuit of length d_i meters, we divide d_i

by 2.0×10^8 m/s. The transmit time is the amount of time required to transmit packets; typically transmit time is measured with respect to average packet length. Then, transmit time can be calculated as the average packet length ($E[L]$) divided by the data rate of the circuit (β_k). The queueing delay captures the amount of time packets spend in the buffer waiting to be transmitted. The time spent in the queue typically is represented by the average time in the queue. To compute average time in the queue, we model the network at each site as an M/M/1/ ∞ queue. We assume a single class of traffic with a packet arrival rate of λ_i ($\hat{\lambda}_i$) and a single “server” with a maximum service rate of $\mu_k = \beta_k/E[L]$ if a circuit of size k is installed. The installed circuit size is a decision variable, so $\sum_{k \in K} \mu_k z_{ik}$ is the expression that represents the bandwidth of the installed circuit. For an M/M/1/ ∞ queueing model with an arrival rate of λ_i and a service rate of $\sum_{k \in K} \mu_k z_{ik}$, the expected waiting time in the system is $\frac{1}{\sum_{k \in K} \mu_k z_{ik} - \lambda_i}$; the time in the system includes both time in the queue and transmit time. The latency constraint (4) follows easily from here.

As formulated, the model is a 0-1 integer program with a set of nonlinear constraints, which make the problem difficult to solve. The next section describes an equivalent formulation in which all of the constraints are linear.

3.3 Linear Formulation

Instead of explicitly including the latency calculation as a constraint, as in (4), an alternative formulation can be developed in which only feasible configurations are included as possibilities in the model. A configuration at a site specifies a circuit size and whether or not a content engine is installed. A configuration is *feasible* if the corresponding latency is no larger than the maximum latency value. Since a configuration specifies a particular circuit size and a content engine or not, its associated cost can be easily calculated. Using this idea, an equivalent formulation can be developed in which the latency calculation is performed as a preprocessing step. Moreover, rather than having two decision variables (i.e. one for placing a content engine and one for determining the circuit size), we can use one decision variable to determine the optimal configuration at each location. Of course, this comes at a price, as we must enumerate these combinations a priori. In the new model, $j = 0$ denotes a configuration without a content engine and $j = 1$ denotes a configuration with a content engine. Several new parameters for the new formulation are required:

Parameters

$$f_i = \text{fixed cost of circuit at site } i, \forall i \in I$$

v_k = variable cost of circuit of size k , $\forall k \in K$

g_{ik} = total cost of service at site i without a content engine and with bandwidth size k , $\forall i \in I, \forall k \in K$

h_{ik} = total cost of service at site i with a content engine and with bandwidth size k , $\forall i \in I, \forall k \in K$

Decision Variables

$$y_{ijk} = \begin{cases} 1 & \text{if a content engine is placed at site } i \text{ with bandwidth } k, \forall i \in I, \\ & \forall k \in K, \forall j \in J \\ 0 & \text{otherwise} \end{cases}$$

Objective Function

$$\text{Minimize } \sum_{i \in I} \sum_{k \in K} g_{ik} y_{i0k} + \sum_{i \in I} \sum_{k \in K} h_{ik} y_{i1k} \quad (8)$$

$$\text{subject to } \sum_{i \in I} \sum_{j \in J} \sum_{k \in K} \beta_k y_{ijk} \leq T \quad (9)$$

$$\sum_{j \in J} \sum_{k \in K} y_{ijk} = 1 \quad \forall i \in I \quad (10)$$

$$y_{ijk} \in \{0, 1\} \quad \forall i \in I, \forall j \in J, \forall k \in K \quad (11)$$

The objective function (7) minimizes the total cost based upon the cost of not placing a content engine (when $j = 0$) or placing a content engine (when $j = 1$) plus the total cost of leasing the circuits. Constraint (8) ensures that the total bandwidth does not exceed T , the bandwidth at the hub. Constraints (9) ensure that each site has a circuit and is thus connected to the hub. Again, the decision variables are binary, as in (10). The problem has 1692 variables and 142 constraints, all linear.

As latency varies depending on whether or not a content engine is placed at the location, the total cost is dependent on two parameters: latency at the site with a content engine and latency at the site without a content engine. The maximum latency value is checked against these values, and if the latency without a content engine is less than the maximum, the total cost is equal to the cost of a content engine plus the cost of the circuit at that location. Otherwise, the latency without a content engine is checked against the maximum latency value, and the total cost equals the cost of the circuit. If neither latency value is less than the maximum value, then total cost is set to a large value. Mathematically, this pre-processing check can be stated as follows:

For all $i \in I, k \in K$

$$\begin{aligned} \text{if } j = 0 & \quad \text{if } \frac{d_i}{s} + \frac{1}{\mu_k - \lambda_i} \leq \alpha_i \\ & \quad \text{then } g_{ik} = f_i + v_k \\ & \quad \text{else } g_{ik} = \infty \end{aligned}$$

$$\begin{aligned} \text{if } j = 1 & \quad \text{if } \frac{d_i}{s} + \frac{1}{\mu_k - \lambda_i} \leq \alpha_i \\ & \quad \text{then } h_{ik} = c_i + f_i + v_k \\ & \quad \text{else } h_{ik} = \infty \end{aligned}$$

4 Computational Experiments

In this section, we describe the methods for data collection and the experimental design. We report the results of the computational experiments on the base model and the sensitivity analysis.

4.1 Data Collection

The data for this research was collected via existing monitoring tools. Some items of note:

- The cost of a content engines is based upon list price.
- The sample size is 141 sites in North America.
- Average arrival rate is based on average measurements over a 6-month period from May 2005 until October 2005.
- Average arrival rate after deploying a content engine is based upon statistics that are collected from test content engines currently in the field.
- Maximum allowed latency, or α_i , is a fixed value of 200 milliseconds for all sites i , as determined by the service level agreements at the company.
- Total cost of each circuit is based upon a predetermined value with the network vendor, converted to an annual cost.
- All data, except for the confidential values for f_i , is given in the Appendix.

- The original data for $E[L]$ was modified to 512 if the data collected for that circuit was less than 512, since the minimum ethernet packet length is 512.

4.2 Base Model Results

The linear formulation was modeled using the AMPL modeling language and solved with the CPLEX solver, version 9.1 [5].

The optimal objective function value for this base model was computed to be \$1.24 million. The optimal configuration is given in Table 1. As shown in the table, 13 content engines were placed in the network with all five circuit sizes in use, although no site required the maximum bandwidth circuit size. The locations chosen to have content engines were either 512 Kbps sites or 128 Kbps sites originally, indicating that the use of a content engine could allow them to use a lower bandwidth of 56 Kbps. The largest sites at 1544 Kbps seemed to get a reduction in bandwidth size but no content engine assignment.

Table 1: Solution for Base Model.

Circuit size	Locations without content engines
1	6-9,15,17-21,26,28,29,32,33,36,41,42,44,48,52,54-57,60,64,70,74-78,80,82,85,88,91-94,97,102,103,106,107,110,112,113,116-121,123,125-127,131,132,134,136,140
2	4,14,22,34,35,37,38,43,46,51,53,86,105,109,122,124,128-130,133,138
3	1-3,5,12,13,16,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,135,137,139,141
4	10,11,31,39,45,69,89,90,104,
5	65
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,114,118

4.3 Sensitivity Analysis

A number of factors are likely to influence the model's solution. In particular, we expect that the locations selected for content engines and the circuit sizes will vary with the maximum latency value, the arrival rates, and the costs. We have designed a set of experiments to investigate the impact of these factors.

4.3.1 Change in Maximum Latency Value

In the first experiment, the maximum latency values were varied to evaluate the effect on the optimal solution. The α_1 values used were 0.1, 0.15, 0.25, and 0.3 seconds. Note that

the original value was 0.2 seconds. Table 2 reports the resulting objective function values, and the tables thereafter show the solutions corresponding to each maximum latency value. The numbers in **bold** represent changes from the base model values.

According to the results, the objective values for the scenarios in which maximum latency is less than 0.2 seconds are larger. As expected, this indicates that, with more strict latency requirements, a more expensive configuration is needed. By looking at the percentage change in the objective value versus the percentage change in the maximum latency value, we can note that the model is more sensitive when the latency is reduced than when it is increased.

Note that when $\alpha_i = 0.1$ and $j = 1$, no sites were selected to receive content engines.

Table 2: Objective function values when altering maximum latency α_i .

α_i	% $\Delta\alpha_i$	Objective Value	% Δ Objective Value
100	-50%	\$1,306,284	+5.4%
150	-25%	\$1,272,272	+2.7%
250	+25%	\$1,215,364	-1.9%
300	+50%	\$1,207,636	-2.5%

Table 3: Configuration when $\alpha_i = 0.1$ for $j = 0$.

Circuit size	Locations without content engines
1	6,7,9,15,17-20, 23-25 ,28,29, 31,32,33,36,41,44,50,52,55-57,65-67 , 70, 71,72,74,75,76,77,78,80,82,87,88,91-92,94,95,97,98,101,102,103,105 , 107, 108,110,113,114,116,117,118,120,121,123,125-126,131,134,136,140
2	4,8,14, 21,22,26,34,35,38,42,43,46,48,51,53,54,60, 64,85,86,93 , 97,106,109,112,119,127,129,132,138
3	13,30, 37,58,83,96,100,111,122,124,128,130,133,137, 139,141
4	1-3,5,10,11,12,16,27,31,40,45,47,49,59, 61-63,68,69,73,79,81 , 84,89,90,99,104,115,135
5	39,45

Table 4: Configuration when $\alpha_i = 0.15$.

Circuit size	Locations without content engines
1	6,7,9,15,17-21, 23-25 ,28,29, 31,32,33,36 ,41,44,48, 50,52,55-57,64,65-67 , 70, 71,72 ,74,75, 76,77,78,80,82,87,88 ,91-93,94,97, 98,101 ,102,103,107, 110,113, 114,116,117,118,119,120,121,123,125-126,131,132,134,136,140
2	4,8,14,22, 26,34,35,37,38,42,43,46,51,53,54,60,85,86,97,105,106,109,112 , 124, 127,128-130,138
3	3,12,13,16,30,40,58,63,83,84,96,100,108,111,115, 122,133,137,139,141
4	1- 2,5,10,11,27,31,39,45,47,49,59 , 61-62,68,69, 73,79,81,89,90 , 99,104,135
5	45
	Locations with content engines
1	25,50,71,76,87,95,98

Table 5: Configuration when $\alpha_i = 0.25$.

Circuit size	Locations without content engines
1	6-9, 14,15,17-21,22,26,28,29,32,33,36,38 ,41, 43,42,44,48,51 . 52,54-57,60,64. 70,74-78, 80,82,85,88,91-94,97,102,103,106,110,112,113,116-121. 123,125-127,131,132,136, 138,140
2	3,4,11,13,34,35,37,40,46,53,58,86,105,108,109,111,115,122,124 , 128-130,133, 141
3	1-2,5,12,16,27,30, 31,39,47,49,59,61-63,66,68,69,72,73,79,81,83,84 , 96,99,100, 104,135, 137,139
4	10,45,89,90
5	65
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101, 107,114,118,134

Table 6: Configuration when $\alpha_i = 0.3$.

Circuit size	Locations without content engines
1	6-9,14,15,17-21,22,26,28,29,32,33,34,35,36,38,41,43,42,44,48,51, 52, 54-57,60,64,70,74-78, 80,82,85,88,91-94,97,102,103,106,110,112, 113,116-121,123,125-127,131,132,136,138,140
2	3,4,11-13,37,40,46,53,58,84,86,96,99,105,108,109,111,115,122,124, 128-130,133,141
3	1-2,5,16,27,30,31,39,47,49,59,61-63,66,68,69,72,73,79,81,83,100,104, 135, 137,139
4	10,45,65,89,90
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,107,114,118,134

4.3.2 Accounting for Percentage of Web Traffic

In the base model, the packet arrival rates associated with the use of a content engine were based solely upon http traffic. These value were obtained by surveying the pilot content engine web caching statistics for an approximate percentage of savings on http traffic. However, the packet arrival rates associated without the use of a content engine were based upon all types of network traffic arriving at the site. Since it is not possible to measure what percentage of the total traffic is http, we estimated the value and then varied it to see how it affects the solution. In the experiments, we adjust the packet arrival rate with a content engine ($\hat{\lambda}_i$) by varying the percentage of http savings in order to account for the fact that not all traffic at the site is http. In other words, we calculated the packet arrival rate with a content engine as follows:

$$\hat{\lambda}_i = \lambda_i(1 - \gamma(\text{current percentage}))$$

Note that in the base model, γ is equal to 1. Three alternate values of γ were evaluated: 0.25, 0.5, and 0.75. The objective value for each value of γ turned out to be \$1.24 million. Results for each value of γ are given in Tables 7, 8, and 9. The interesting result from this experiment was that the solutions had the same objective values as the base model, and the same configurations of bandwidth values and content engine placements. In other words, each run of the experiment resulted in the same solution that is found in Table 1.

It can then be assumed that the packet arrival rates with the content engine are still

relevant even though they do not account for the all types of packets arriving at the site. Attempting to include this fact in the model by varying the $\hat{\lambda}_i$ values does not make much of a difference in the overall solution.

Table 7: Configuration when $\gamma = 0.25$.

Circuit size	Locations without content engines
1	6-9,15,17-21,26,28,29,32,33,36,41,42,44,48,52,54-57,60,64,70,74-78,80,82,85,88,91-94,97,102,103,106,107,110,112,113,116-121,123,125-127,131,132,134,136,140
2	4,14,22,34,35,37,38,43,46,51,53,86,105,109,122,124,128-130,133,138
3	1-3,5,12,13,16,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,135,137,139,141
4	10,11,31,39,45,69,89,90,104,
5	65
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,114,118

Table 8: Configuration when $\gamma = 0.5$.

Circuit size	Locations without content engines
1	6-9,15,17-21,26,28,29,32,33,36,41,42,44,48,52,54-57,60,64,70,74-78,80,82,85,88,91-94,97,102,103,106,107,110,112,113,116-121,123,125-127,131,132,134,136,140
2	4,14,22,34,35,37,38,43,46,51,53,86,105,109,122,124,128-130,133,138
3	1-3,5,12,13,16,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,135,137,139,141
4	10,11,31,39,45,69,89,90,104,
5	65
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,114,118

Table 9: Configuration when $\gamma = 0.75$.

Circuit size	Locations without content engines
1	6-9,15,17-21,26,28,29,32,33,36,41,42,44,48,52,54-57,60,64,70,74-78,80,82,85,88,91-94,97,102,103,106,107,110,112,113,116-121,123,125-127,131,132,134,136,140
2	4,14,22,34,35,37,38,43,46,51,53,86,105,109,122,124,128-130,133,138
3	1-3,5,12,13,16,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,135,137,139,141
4	10,11,31,39,45,69,89,90,104,
5	65
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,114,118

4.3.3 Adjusting the Circuit Costs

In the third experiment, we varied the circuit costs to evaluate their effect on the solution. We used the following procedure to determine the adjustment parameters. The original cost values for each bandwidth type (v_k) were plotted against the bandwidth values (β_k). Then, a regression line was fit to the data points and an equation for the regression line was found. The equation was in the form

$$v_k = a \ln(\beta_k) - 325$$

where a is the cost parameter that was manipulated. The value for a is 140 in the base model. The other values for a were 100, 120, 160, and 180.

The objective function values for this experiment are reported in Table 10. The solutions corresponding to each value of a are also shown in the Tables 11, 12, 13, and 14. Again the values in **bold** are the differences from the base model.

When the cost is low (i.e., $a = 100$), the model is more inclined to choose circuits of varying sizes and fewer content engines. However, as the cost increases, the model tends to choose more content engines and more circuits of a small bandwidth value. This indicates that the cost is too expensive for a higher circuit and that the better configuration is to use smaller circuits in conjunction with content engines.

In addition, the percentage change in the objective value decreases rapidly as a decreases. Examining the percentage change in objective values from the base model also indicates

that total cost decreases at a greater rate than it increases when the circuit costs are altered.

This experiment suggests that if the cost of a circuit were to continue to increase, then the company would be more likely to use content engines than they are today with respect to the base model. Therefore, content engines could be an economical solution in the long term.

Table 10: Sensitivity of objective function value to circuit costs.

a	$\% \Delta a$	Objective value	$\% \Delta$ Objective value
100	-28.6%	\$982,313.56	-20.7%
120	-14.3%	\$1,116,635.08	-9.9%
160	+14.3%	\$1,317,963.90	+6.4%
180	+28.6%	\$1,365,921.21	+10.2%

Table 11: Configuration when $a = 100$.

Circuit size	Locations without content engines
1	7-9,19-21,26,41,42,48,54,56,60,64,70,74,76-78,85,88,92,93,97,103,106,112,113,118,119,123,125-127,131,132,134,136,140
2	4,6,14,15,22,28,29,32,33,34,35,36,37,38,43,44,46,51,52,53,55,57,75,80,82,86,91,102,105,107,109,110,116,117,120,121,122,124,128-130,133,138
3	1-3,5,12,13,16,17,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,135,137,139,141
4	10,11,18,31,39,45,69,89,90,94,104,
5	65
	Locations with content engines
1	25

Table 12: Configuration when $a = 120$.

Circuit size	Locations without content engines
1	6-9,15,19-21, 23,24,26,28,29,32,33,36,41,42,44,48,50,52,54-57,60,64,70,71,74,75,76,77,78,80,82,85,87,88,91-94,95,97,98,101,102,103,106,107,110,112,113,114,116,118,119-121,123,125-127,131,132,134,136,140
2	4,14,22,34,35,37,38,43,46,51,53,86,105,109,122,124,128-130,133,138
3	1-3,5,12,13,16,17,27,30,40,47,49,58,59,61-63,66,68,72,73,79,81,83,84,96,99,100,108,111,115,117,135,137,139,141
4	10,11,18,31,39,45,69,89,90,104,
5	65
	Locations with content engines
1	25,67

Table 13: Configuration when $a = 160$.

Circuit size	Locations without content engines
1	6-9,10,11-13,15,17-21, 22,26,28,29,30,31,32,33,34,36,38,39,41,42,43,44,45,46,48,49,51, 52, 54-57,58,60,64,65,66,69,70,72,74-78,80,82,83,85,88,89,90,91-94,96,97,100,102,103, 104,105,106,107,108,110,112,113,116-121,123,125-127,129,131,132,134,136,137,138,139,140
2	4,14,35,37,46,53,86,109,122,124,128,130,133
3	1-3,5,12,16,27,40,47,59,61-63,68,73,79,81,84,99,111,115,135,141
	Locations with content engines
1	23,24,25,50,67,71,76,87,95,98,101,114,118,126

Table 14: Configuration when $a = 180$.

Circuit size	Locations without content engines
1	1-3,5,6-8,10-13,15,16,17-21, 22,26,27,28,29,30,31,32,33,34,36,38-40,42,43,44,45-47,48,49,51,52, 54-57,58-61,64,65,66,68,69,70,72,73,74-78,79,80,81,82,83,84,85,88,89,90,91-94,96,97,99,100,102,103,104,105,106,107, 108,109,110,111,112,113,115,116-121,122,124,123,125-127,129,131,132,134,136,137,138,139,140
2	4,14,35,37,46,53,86,128,130,133
3	135,141
	Locations with content engines
1	9,23,24,25,41,50,60,67,71,76,87,95,98,101,114,118,126

4.3.4 Restricting the Model to Use Larger Circuit Sizes

Although the base model gives us a potentially feasible solution for implementation, it does not necessarily offer the most realistic solution for the company in this case study. The solution to the base model selects many smaller circuits, which may not be feasible for many of the larger remote sites at the company due to requirements not captured in this model. A circuit with a smaller bandwidth will not necessarily be able to handle the traffic as it does with the existing circuit today.

In order to make the model more representative of the current design, we developed a variation of our model in which we restrict the set of feasible circuit sizes. Given a current circuit size of k , the only feasible choices are $k - 1$, k , and $k + 1$. These options were selected since the company would likely only consider changing bandwidth by one size.

Including this restriction and solving the base model, we obtained an objective function value of \$1,440,088, or an increase of 16.2%. The optimal configuration scenario is given in Table 15. Although this model offers a more realistic view of the types of circuits that would be used at this company, note that it does not choose to use any content engines. Thus, if the company wants to achieve the savings offered by the content engine, they cannot restrict the choice of bandwidth size.

Table 15: Base Model Solution with Circuit Size Restriction.

Circuit Size	Locations without content engines
1	2,12,25,27,37,42,53,84,85,93,97,108,111,133
2	1,3,5,9,11,15,16,22,38,40,41,48,60-63,65,67,72,81,86, 117,119,124,135,136
3	4,6-8,10,13,14,19,20,23,24,26,28-33,35,36,39,44,47, 49,50,52, 54-59,66, 68-71,75-80,82,87-92,95,98-104,107, 110,112-116,118,120-123,125,126,128,130-132,134,137,140
4	21,45,83,96,109,139,141
5	17,18,34,43,46,51,64,74,94,105,106,127,129,138

5 Conclusions

The research in this case study presents one possible solution for reducing latency in a large corporate network. This model could be extended within the company in the case study to include the European and Asian regions as well.

The content engine is a possible solution for improving latency, but it is just one of many solutions. This thesis has shown that the bandwidth at remote sites plays a key role in network efficiency. Changing the circuit size has been the most likely way to keep latency at a reasonable level. Circuits can be expensive and methods such as using content engines could become more valid as circuit costs rise and networks continue to develop. If the use of content engines were to become more prevalent, investigating the possibility that remote sites could share the content engines would be important. If the network topology were extended to a meshed topology (rather than just a star topology), the remote sites could transmit information to one another without having to go back to the hub. This change could offer yet another way to reduce latency.

References

- [1] Bartlett, John, Peter Sevcik, and Sean Moore. Economies of QOS on WAN Access Lines. *Business Communications Review*, Oct. 2004.
- [2] Berger, Rosemary T. and S. Raghavan. Long-Distance Access Network Design. *Management Science* **17**, 309 – 325, 2004.
- [3] Berman, Oded and Dmitry Krass. Facility Location Problems with Stochastic Demands and Congestion. *Facility Location: Applications and Theory.*, Z. Drezner and H.W. Hamacher, eds. Springer-Verlag, Heidelberg, 2002.
- [4] Fabbi, Mark. *Latency: The Silent Killer of Application Performance*. Gartner Research. 14 Sept. 2004.
- [5] ILOG CPLEX 9.0. ILOG CPLEX Division, Incline Village, NV, 2005.

Appendix

This appendix presents the original cost and arrival data used in the experiments. Tables 16 and 17 present data concerning bandwidth and original problem parameters while the remaining Tables 18 through 21 present the original location-specific data.

Table 16: Original Data for Each Bandwidth Type.

k	v_k	β_k
1	2640	56000
2	4572	128000
3	5544	256000
4	6516	5120000
5	6852	768000
6	8784	1544000

Table 17: Original Data Parameters.

c_1	s	α_I	T
5500	200000000	0.2	44736000

Table 18: Original Data for Locations 1 through 35.

i	d_i	λ_I	λ_i	$E[L]$
1	2055567	1.49	1.78	3198
2	864845	2.62	3.57	2684
3	4669174	1.14	1.39	2186
4	805444	8.91	12.09	512
5	2654565	0.01	0.02	4369
6	2654565	16.20	19.14	512
7	2439894	22.35	30.91	512
8	4537224	2.44	3.17	512
9	25910	2.07	2.60	512
10	25910	18.11	22.53	1028
11	25910	2.19	2.98	3557
12	392889	0.53	0.73	2769
13	922717	7.18	8.88	938
14	1143503	1.50	1.99	836
15	1365126	1.98	2.67	1497
16	1365126	1.47	1.82	2899
17	1365126	23.48	30.37	512
18	4380795	38.32	45.14	512
19	4341431	41.35	52.63	512
20	798186	3.18	4.08	336
21	1989181	2.34	3.27	434
22	3398452	2.77	3.94	628
23	3837899	4.56	5.82	482
24	501407	3.47	4.70	512
25	1102803	1.36	1.75	512
26	1102803	0.82	1.12	512
27	865731	2.94	3.57	2710
28	642144	8.99	11.51	512
29	2707834	2.04	2.86	1296
30	2035370	10.94	13.41	857
31	2035370	1.12	1.34	4166
32	1290758	9.09	11.49	512
33	443246	6.73	8.47	512
34	2340083	5.16	7.05	512
35	2340083	5.56	6.79	512

Table 19: Original Data for Locations 36 through 70.

i	d_i	λ_j	λ_i	$E[L]$
35	2340083	5.56	6.79	512
36	1449327	13.69	17.32	551
37	1449327	3.21	3.88	940
38	2460655	1.15	1.54	908
39	2761747	1.40	1.73	3941
40	1937972	0.72	0.96	2460
41	1365126	2.99	3.58	360
42	4581062	1.04	1.38	512
43	38801	8.45	11.10	356
44	4347755	10.06	13.02	512
45	155350	2.33	3.20	5512
46	1785181	6.75	8.09	512
47	1721918	1.02	1.27	3741
48	958525	1.49	1.98	512
49	4329843	3.91	5.06	2182
50	4478981	2.11	2.48	512
51	2063887	4.87	6.10	512
52	2063887	16.33	19.37	512
53	2063887	7.24	9.01	605
54	3837320	4.42	5.39	512
55	348198	8.21	10.29	512
56	1370421	3.63	5.04	1677
57	1370421	13.47	19.22	512
58	848028	4.16	5.23	1275
59	1504624	0.50	0.67	3439
60	2663754	3.67	4.71	512
61	1988618	1.01	1.39	3071
62	1962482	1.09	1.52	3207
63	2091117	1.08	1.31	3016
64	48908	0.71	0.90	717
65	482996	4.22	5.10	5678
66	2471341	0.51	0.71	3343
67	2471341	3.95	5.07	512
68	4338534	0.54	0.72	3472
69	2360682	1.30	1.78	3654
70	839563	5.55	7.35	3247

Table 20: Original Data for Locations 71 through 105.

i	d_i	λ_I	λ_i	$E[L]$
71	95659	2.97	4.20	512
72	1130918	1.70	2.16	2806
73	2450226	1.02	1.37	3385
74	2450226	22.03	30.06	4091
75	1803141	2.52	3.52	807
76	1281150	4.04	5.30	512
77	474467	1.97	2.58	4370
78	2426408	2.16	2.99	3330
79	2426408	0.74	1.03	4008
80	2426408	7.43	9.71	512
81	1954822	1.11	1.58	3178
82	4341801	3.05	4.25	915
83	4319769	9.36	13.01	855
84	1052302	0.60	0.71	2832
85	12875	3.22	4.35	512
86	476446	10.11	12.09	512
87	154240	1.74	2.37	616
88	485668	17.20	21.47	1082
89	4184053	44.77	52.68	512
90	4184053	41.33	48.99	512
91	863172	14.69	18.32	512
92	2378707	6.08	8.18	1394
93	2427502	0.62	0.86	512
94	115873	8.37	9.93	1968
95	4268270	2.22	2.73	512
96	131741	4.88	6.19	1280
97	285997	0.62	0.76	791
98	227803	1.78	2.36	512
99	140303	0.53	0.67	2786
100	840174	8.18	9.74	1157
101	1304293	3.46	4.79	512
102	1675391	7.83	9.73	580
103	628368	0.75	1.00	3631
104	2464694	1.00	1.23	4045
105	2464694	9.85	13.74	512

Table 21: Original Data for Locations 106 through 131.

i	d_i	λ_I	λ_i	$E[L]$
106	121232	3.04	4.29	512
107	3802365	4.52	5.65	512
108	3802365	0.27	0.37	2193
109	1058739	10.15	12.92	512
110	2331441	5.91	7.48	512
111	2331714	9.78	12.74	749
112	553566	4.20	5.04	512
113	1954822	610.52	869.06	4848
114	4567672	3.97	4.90	512
115	752513	0.57	0.71	2529
116	717687	8.08	11.22	512
117	4549519	0.90	1.22	3689
118	3074233	2.90	3.55	512
119	4440615	1.21	1.46	512
120	4629648	8.09	10.44	512
121	4629648	6.70	8.77	512
122	4629648	2.88	4.00	1132
123	4629648	10.63	14.06	1081
124	4313348	0.87	1.23	1351
125	1858052	3.96	5.46	2693
126	1797525	2.83	3.38	512
127	143618	4.69	5.55	512
128	240822	8.98	11.04	689
129	2057385	8.55	11.66	512
130	4523576	10.94	13.42	512
131	3823898	13.35	17.54	512
132	3823898	1.45	1.97	512
133	4580322	1.49	2.05	1438
134	4580322	4.91	5.81	512
135	291887	1.31	1.70	3202
136	1338894	1.29	1.56	4398
137	2248366	17.83	23.03	512
138	2078323	4.55	6.27	512
139	98701	21.37	27.25	512
140	4343265	48.70	62.20	512
141	8047	6.78	8.77	993

Biography

The author, Gretchen Trump, was born in Riverside, CA on March 29, 1979 but spent the majority of her youth living in Lancaster, PA. In 2000, she obtained her B.S. degree in Mathematics from Muhlenberg College located in Allentown, PA. Since then she has worked as a Network Engineer for a Fortune 500 company while attending Lehigh University part-time as a graduate student.

**END OF
TITLE**