

Performance Modeling of Video-on-Demand Systems in Broadband Networks

Eric Wing Ming Wong, *Senior Member, IEEE*, and Sammy Chi Hung Chan, *Member, IEEE*

Abstract—A video-on-demand (VoD) system allows a viewer to choose a video of his choice such as movies, electronic encyclopedia, or educational videos, which is delivered to him real time via a network. This paper is a performance study of a distributed video server system in a fully connected backbone network. Under the assumptions of uniform loading and symmetrical network, analytical models are proposed for two server selection strategies (Random and Least Loaded) and for two reservation schemes (the strict reservation and the residual reservation). The models can be used to determine the call blocking probability and the requirements of network bandwidth given the capacity of video servers. In addition, our models provide a useful tool for VoD operators to design and dimension their systems.

Index Terms—Central video server, distributed video server system, fixed-point iteration method, local video server, video on demand.

I. INTRODUCTION

THE TELEPHONE network of the past is fast evolving into an integrated services broadband network where data and video traffic are carried in increasing proportions. With the advances in coding, storage, transmission and networking technologies, video on demand (VoD) is likely to become one of the most important traffic types in the future broadband network. VoD service can offer instant access to a large selection of video sources, such as movies, electronic encyclopedia or educational videos. VoD trials in different areas of the world [1], [2] have demonstrated that VoD is technically feasible and Laser Disk-like functions, such as pause, fast-forward, chapter search, slow play, and repeat, can easily be implemented. The classification of VoD services can be based on the degree of interaction allowed to the users. There are pay-per-view (PPV), quasi video-on-demand (Q-VoD), near video-on-demand (N-VoD), and true video-on-demand (T-VoD). Recently, the world-first commercial T-VoD system^{1,2} was launched in Hong Kong.

Our focus in this paper is on performance evaluation of distributed T-VoD. In a backbone network serving hundreds of thousands of VoD subscribers, it is imperative that many video servers are to be used and placed close to the subscribers for balancing the load in the network. In fact, a recent survey on VoD

Manuscript received April 9, 1998; revised March 6, 2001. This paper was recommended by Associate Editor S.-F. Chang.

The authors are with the Department of Electronic Engineering, City University of Hong Kong, Tat Chee Avenue, Kowloon, Hong Kong (e-mail: ewong@ee.cityu.edu.hk; schan@ee.cityu.edu.hk).

Publisher Item Identifier S 1051-8215(01)05280-6.

¹[Online]. Available: <http://www.netvigator.com/IMS>.

²[Online]. Available: <http://www.itvhk.com>.

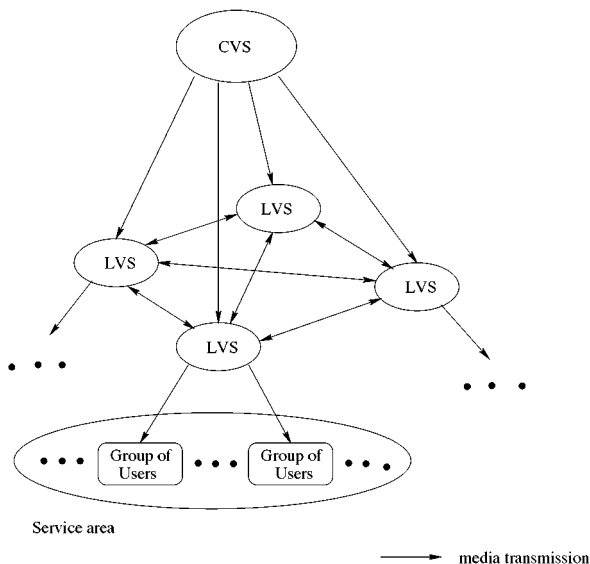


Fig. 1. A fully connected LVS network for VoD systems.

systems [3] showed that the distributed design costs no more than the centralized design, while it is capable of saving considerable network bandwidth with better service quality. Designing a distributed server system is inherently complicated. One has to decide which set of video programs and how many copies of each needed to be placed in each server location and how these should be updated with the changes in locations of other servers and traffic composition, occurrence of system fault, etc.

Fig. 1 shows an example of a VoD system in a fully connected network where there are two types of video servers: the local video server (LVS) and the central video server (CVS). A CVS stores all the video programs in high-capacity optical disks or magnetic tapes. LVSs with on-line mass storage store the popular video programs which are downloaded periodically (say daily or weekly) from CVS. The existence of LVSs releases excessive access to CVS in the following way. When the less popular videos are requested, the CVS will load the requested videos into the on-line storage and then the videos are transmitted to the users through the network. If the CVS is busy, such requests are blocked. When the subscribers request the popular videos which constitute the majority of the demand, the LVS associated with the request is tried. If it is busy, other LVSs are tried. Such a popular video request is blocked only when none of the LVSs is available. To which alternate LVS such kind of video request is assigned has a strong impact on the system performance. As an example, in Hong Kong, there are four LVSs located at Hong Kong Island, Kowloon, and New Territories. This paper focuses on the popular video accessing and studies

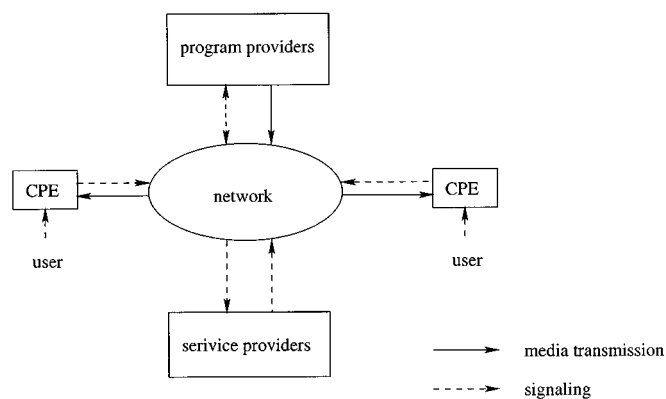


Fig. 2. Generic VoD architecture.

a number of corresponding LVS selection strategies for performance optimization.

II. VoD SYSTEM ARCHITECTURE

The generic architecture of a VoD system is shown in Fig. 2. There are four important elements: 1) the network; 2) the program provider; 3) the service provider; and 4) customer premises equipment (CPE). A brief overview of these network elements is as follows.

Network: The network provides the interconnection of the network elements in the VoD system. Beside program transfer, the network also includes other important functions, e.g., signaling and network management. The network could be of any kind as long as to be able to provide enough bandwidth for such kind of video application. A VoD implementation using ATM technology appeared in [4].

Program Provider: There can be one or more than one program provider, which provide a wide range of video programs.

Service Provider: The user generates a particular video request to the service provider, who will obtain the necessary material from program providers and deliver it to the user on the facilities of the network. Thus, the service provider acts as an agent of the user and can access various types of program providers. It is possible that the network, program providers, and service providers all belong to the same organization, but in general, they will be distinct. In fact, anyone with marketable materials can offer his services to the user through a service provider. There are two main types of storage [4]:

Video Server: This is random access type of storage, such as hard disk, solid state memory, etc., which provides a real-time playback of video program upon the user request. The video could be stored in encoded format (e.g., MPEG) so as to reduce the storage cost. In this paper, the LVS is of this type.

Video Library: It is an archiving system. Upon a user request, the program in the video library will be batch-loaded to the video servers such that the user could view the video program in real time. It takes time for batch-loading a program, so the video library should only store less popular program in order to reduce the average waiting time to the users. In this paper, CVS is of this type.

CPE: The user has a display such as TV screen, or a set-top box (STB), to control the display, and a device such as a remote

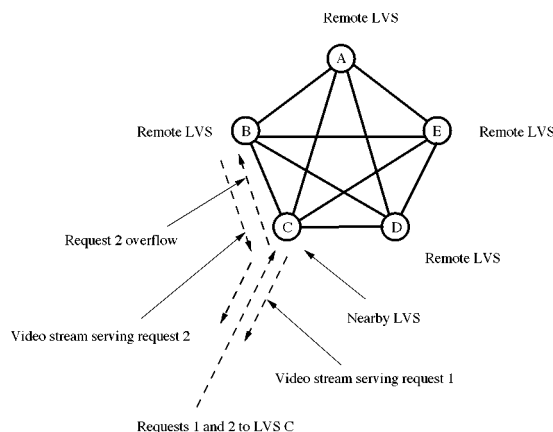


Fig. 3. The simplifying VoD system.

control or a keyboard to interact with the system. User interactive operations may include stop, speed-up, slow-down, jump forward, etc.

III. SIMPLIFYING ASSUMPTIONS, SERVER SELECTION STRATEGIES, AND RESERVATION SCHEMES

As stated before, we, in this paper, focus only on the access of popular movies since it makes greater impact on the system performance. Therefore, we only consider a fully connected network with M LVSs like the one in Hong Kong (i.e., ignoring the existence of CVS for the study). We further assume the network is symmetric and uniformly loaded. Fig. 3 shows such a simplified network. Each (popular) movie is copied to every LVS. In the VoD system, each LVS consists of an array of disks, each of which has a certain number of read-write heads. The I/O transfer rate of each head is generally much higher than the required delivery rate of one movie, so each head can serve more than one movie simultaneously. For simplicity, but without loss of generality, we assume that each head can at most serve one user at a time. There are N heads in an LVS. In other words, each LVS can support a maximum of N users simultaneously.

There are a number of server selection strategies for selecting an LVS to serve a user request. For all strategies considered, a user request is assigned to its nearby LVS located at the local switch first. If this LVS is busy, a remote LVS is tried. Fig. 3 gives an example showing the nearby and remote LVS for requests 1 and 2. Accessing remote LVSs consumes additional bandwidth from the network and may lead to congestion. Fig. 3 gives an example showing the use of remote LVS (for request 2) does cause additional bandwidth when compared with the use of nearby LVS (for request 1). In addition, the *remote* traffic may cause more local traffic overflow (from LVS B, in this example) and make the situation worse. Therefore, it is important to select an appropriate remote LVS such that the (additional) bandwidth cost can be reduced/minimized. The specific remote LVS to be chosen depends on the server selection strategies, two of which are studied in this paper.

1) *Random Strategy:* An LVS is chosen among the available remote LVSs at random. If all the LVSs are busy, the request is rejected.

2) *Least-Loaded Strategy*: The remote LVS with the lightest load is chosen. If all the LVSs are busy, the request is rejected.

Moreover, in order to minimize the system cost, it may be necessary to restrict the volume of remote traffic. This can be achieved by reserving some server capacity to local requests only. For this purpose, two reservation schemes are studied in this paper:

1) *Strict Reservation Scheme*: r server heads are reserved solely for local video requests. In this scheme, each LVS keeps a record of the number of overflow calls occupying its servers. When this number reaches $N - r$, further remote overflow calls will be blocked, even if there are free servers available.

2) *Residual Reservation Scheme*: The last r unoccupied server heads are reserved solely for local video requests. In other words, a video server will not serve any remote video requests if the number of its available server heads is equal to or less than r .

IV. ANALYSIS OF VoD SYSTEMS

Analytic models have been very attractive, not only because of the relatively small computational requirements of the numerical solutions as compared with simulation (we did find our model is much more efficient than simulation, as shown in numerical results section), but also because the modeling problems have been intellectually stimulating and because of problem with simulation other than computational expense [5]. In this paper, we will develop a queueing model for the distributed VoD system and focus on the derivation of the blocking probability and the bandwidth cost (on the fully connected core network) given the capacity of video servers. Other performance measures, such as the end-to-end service response delay, the amount of video buffers required, the network capacity required, etc., are beyond the scope of this paper.

Each LVS can be modeled as a multiserver queue and the set of LVSs in a network can be modeled as a queueing network. Under the Poisson arrival and exponential service time assumptions, the network can be represented as a multidimensional Markov chain with the dimension equal to the total number of LVSs. Although this chain can be solved theoretically, solutions beyond the three-node cases are computationally infeasible due to the curse of dimension.

With that, the development of efficient computational procedures for good approximate solutions is of interest. The fixed-point iteration method [6] for the analysis of dynamic routing in circuit switched networks appears to be applicable for the analysis of the VoD system. But a direct application leads to unacceptable errors when checked by computer simulation. This is due to the assumption that all LVSs are statistically independent. A detailed study of the VoD system reveals that some essential internal structures of the VoD system need to be captured in the model and this leads to a set of considerably more complicated system equations. Our solution approach, however, still makes use of the fixed-point iteration method. Various examples with diverse system parameters show that the new analytical model can give quite accurate results. Therefore, the contribution of the paper is to present a new approximate analytic model with computational efficiency and reasonable

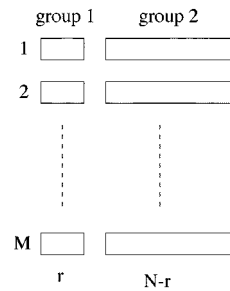


Fig. 4. Queueing model for the residual reservation scheme.

accuracy. Using the model, we could design and dimension the VOD systems. Examples through case studies are shown in Section VI. Some interesting insight and design rules were also found.

Let λ_D be the rate of local offered traffic to an LVS and λ_O be the rate of overflow traffic from an LVS to other (remote) LVSs. Let both of these two traffic streams be Poisson processes and let the service time (or movie showing time) be exponentially distributed with mean equal to one unit time (this assumption is justified by considering the effect of user interaction and the varying lengths of different video programs). Note that, like Erlang B formula for call-blocked-clear telephone networks, the video duration does not need to be exponentially distributed. This exponential distribution assumption used in this paper is just for the sake of simplicity. Here, we will study two LVS selection strategies: Random and Least Loaded, both with two (strict and residual) reservation schemes. The definition of major symbols are summarized in Appendix (Table II).

A. Random Strategy

1) *Residual Reservation*: First, we will find the blocking probability of the whole network. The network can be modeled as two groups of queues, as shown in Fig. 4.

Each group consists of M multiserver queues. Each queue in group 1 has $N - r$ servers, while each queue in group 2 has r servers, which means a total of N servers in an LVS. Note that a server (different from a video server), which is a common terminology in queueing theory, is equivalent to a server head in an LVS. Consider a particular video server (say video server 1). It consists of queue 1 from group 1 and queue 1 from group 2. These two queues are called companion queues. The servers in the two companion queues are used in the following way: the servers in the group 1 queue will always serve the local requests first, while the servers in the group 2 queue will serve the local requests only when the servers in the group 1 queue are all busy. Note that if a request departs from the group 1 queue, one request, if any, in the group 2 queue will be transferred to the group 1 queue such that the above condition holds again (i.e., the group 1 queue will be always filled up first). This transfer is possible since group 1 and group 2 are only *logical* concepts. According to the residual reservation scheme, an LVS consists of a group 1 queue and a group 2 queue where the group 1 queue can serve both local traffic and remote traffic (i.e., overflow traffic from other LVSs), while the group 2 queue is reserved for local traffic only. The input traffic to group 1 represents the aggregated fresh

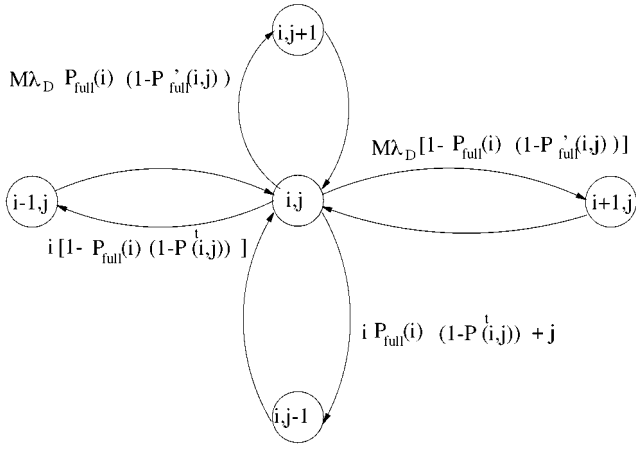


Fig. 5. Transition rates at state (i, j) for the residual reservation scheme.

offered traffic to the network and its rate is equal to $M\lambda_D$. Regarding to the queueing model in Fig. 4, the operation of residual reservation scheme is as follows. When a call (request) arrives, it is randomly directed to one of the queues in group 1 (say queue 1). If the queue is full, it will go into its companion queue (i.e., queue 1) in group 2. If that queue is full as well, it will be redirected back (overflowed) to one of other available queues (i.e., among queues 2 to M) in group 1. If all those queues are full, the call will be blocked.

Let (i, j) represent the state of the network in which there are totally i and j calls in groups 1 and 2, respectively. Fig. 5 shows the transition rates into and out of state (i, j) .

At state (i, j) , an arrival will cause a transition to state $(i, j+1)$ if it is directed to a full queue in group 1 and the companion queue in group 2 is not full; otherwise, the arrival will cause a transition to state $(i+1, j)$. Let $P_{\text{full}}(i)$ be the probability that an arrival is directed to a full queue when there are i calls in group 1. To find $P_{\text{full}}(i)$, consider the problem of inserting i distinct balls into M urns where each urn has the capacity to store N balls. The number of distinct occupancies is (see [7] for details)

$$C_N(i, M) = \sum_{j=0}^M (-1)^j \binom{M}{j} \binom{i + M - j(N+1) - 1}{M-1}. \quad (1)$$

$C_N(i, M)$ can be calculated by the following recursion:

$$C_N(i, m) = \begin{cases} 1, & m = 1; 0 \leq i \leq N \\ 0, & m = 1; i > N \\ \sum_{j=0}^{\min(i, N)} C_N(i-j, m-1), & 2 \leq m \leq M. \end{cases}$$

$P_{\text{full}}(i)$ is simply given by

$$P_{\text{full}}(i) = \begin{cases} 0, & i < N-r \\ \frac{C_{N-r}(i - (N-r), M-1)}{C_{N-r}(i, M)}, & N-r \leq i \leq M(N-r). \end{cases}$$

Let $P'_{\text{full}}(i, j)$ be the probability that a call entering group 2 finds a full queue when the network is state (i, j) . When there

are i calls in group 1, there are at most $\hat{M} = \lfloor (i/(N-r)) \rfloor$ queues that can be filled up. Therefore, only \hat{M} companion queues in group 2 can be occupied since, as mentioned before, group 1 queue must be filled up first. In other words, with i calls in group 1, j is limited to be $\hat{M}r$. So, we have

$$P'_{\text{full}}(i, j) = \begin{cases} 0, & j < N-r \\ \frac{C_r(j-r, \hat{M}-1)}{C_r(j, \hat{M})}, & r \leq j \leq \hat{M}r. \end{cases}$$

Therefore, the total transition rate from state (i, j) to state $(i, j+1)$ is $M\lambda_D P_{\text{full}}(i)(1 - P'_{\text{full}}(i, j))$. On the other hand, the total transition rate from state (i, j) to state $(i+1, j)$ is hence $M\lambda_D [1 - P_{\text{full}}(i)(1 - P'_{\text{full}}(i, j))]$.

At state (i, j) , two types of departures would lead to a transition to state $(i, j-1)$. The first type is simply a departure of a call in the group 2 queue. For the second type, if a queue in group 1 is full and its companion queue in group 2 is not empty, a call departure in the group 1 queue effectively causes a call departure in the group 2 queue since, as mentioned before, the companion group 1 queue has to be always filled up first. The probability $P^t(i, j)$ that, at state (i, j) , a queue in group 2 has no (local) calls at all is given by

$$P^t(i, j) = \frac{C_r(j, \hat{M}-1)}{C_r(j, \hat{M})}. \quad (2)$$

Therefore, the total transition rate from state (i, j) to state $(i, j-1)$ is $i P_{\text{full}}(i)(1 - p^t(j)) + j$. At state (i, j) , a transition to state $(i-1, j)$ occurs when the corresponding group 1 queue is not full (and hence no group 2 call will try to convert to a group 1 call), or when the corresponding group 1 queue is full but the companion queue in group 2 is empty. This happens with rate $i[1 - P_{\text{full}}(i)(1 - p^t(i, j))]$.

Fig. 6 depicts the state space of the network. Let $P(i, j)$ be the probability that the network is in state (i, j) . A new call will be blocked only if all of the queues in group 1 are fully occupied plus the corresponding queue in group 2 being full, since a call will try the nearby LVS and all the remote LVSs before being blocked. So, the call-blocking probability is given by

$$P_B = \sum_{j=r}^{Mr} P'_{\text{full}}(M(N-r), j) P(M(N-r), j). \quad (3)$$

Now, we consider the cost due to transmission of video from remote LVSs. To do that, we try to focus our attention on an LVS and model its behavior. When the nearby LVS is full, the random strategy directs the call (request) to one of the admissible remote LVSs at random. Since the network is symmetric and the loading is uniform, the local call blocking probability is equal to P_B calculated above. Let us denote an LVS to be in state i when it has i servers occupied, and the corresponding state probability be P_i . Since each remote LVS carries the remote traffic (requests) from $M-1$ possible LVSs, but each of them has also $M-1$ remote LVSs to direct its remote traffic to, the total remote traffic A_i to an LVS at state i is

$$A_i = \begin{cases} A, & i = 0, 1, \dots, N-1-r \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

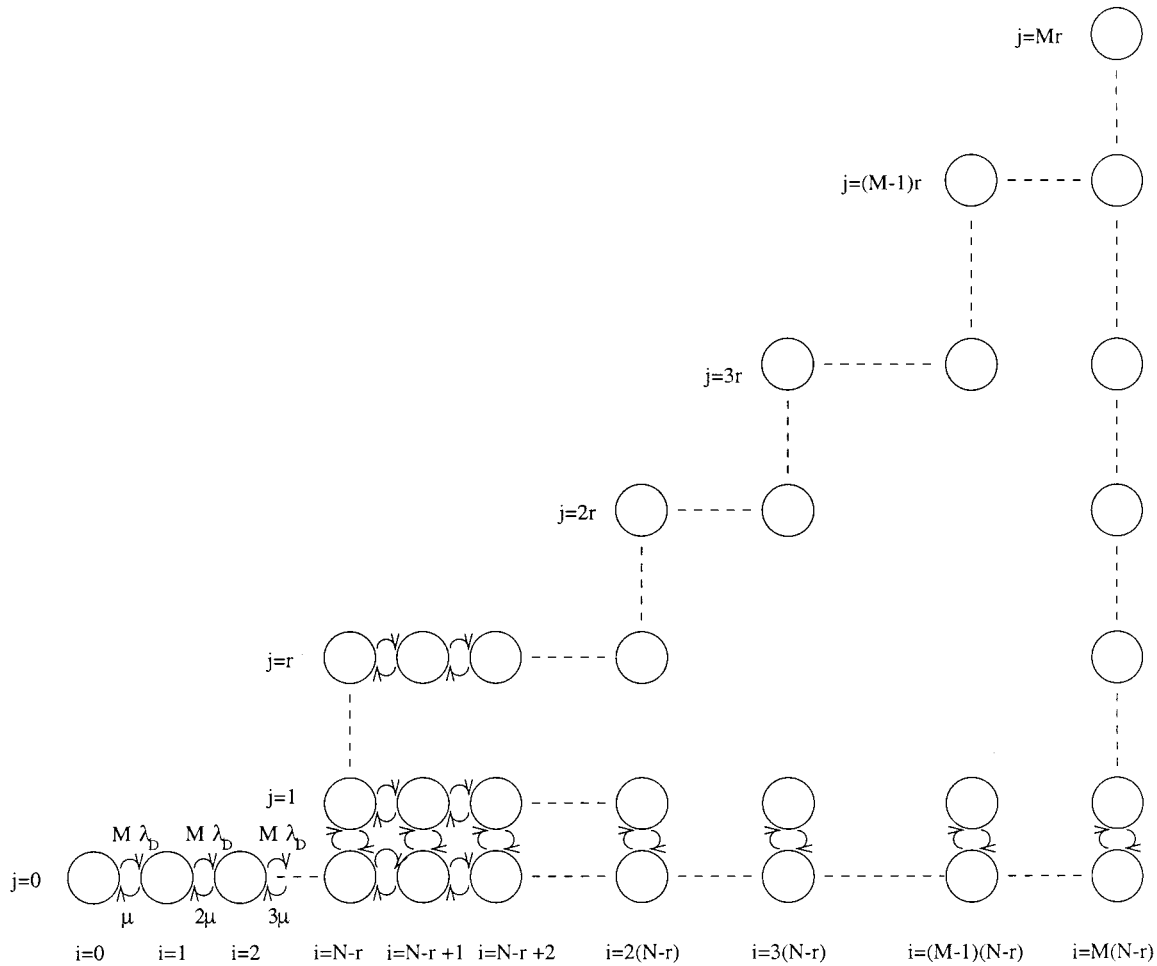


Fig. 6. State transition diagram of the network with the residual reservation scheme.

where A is a constant and A_i is related to λ_D by the following or equation of conservation of flow:

$$\sum_{i=0}^{N-r-1} A_i P_i = \lambda_D P_N \left[1 - \left(\sum_{j=N-r}^N P_j \right)^{M-1} \right]. \quad (5)$$

When an LVS is in state i , the call arrival rate λ_i and the call departure rate μ_i are

$$\begin{aligned} \lambda_i &= \alpha(\lambda_D + A_i), & i &= 0, 1, \dots, N-1 \\ \mu_i &= i, & i &= 1, 2, \dots, N. \end{aligned} \quad (6)$$

Here, we assume that the effective total arrival rate in each state is the actual total arrival rate scaled by a constant α such that the carried traffic is matched with that calculated from the previous network model

$$\alpha \left[\lambda_D(1 - P_N) + A \sum_{i=0}^{N-r-1} P_i \right] = \lambda_D(1 - P_B). \quad (7)$$

From the balance equation, we have

$$P_i = \frac{\mu_{i+1}}{\lambda_i} P_{i+1} \quad (8)$$

$$P_i = \frac{N!}{i! \prod_{k=i}^{N-1} \lambda_k} P_N. \quad (9)$$

Substituting it into the normalization equation, we obtain

$$\sum_{i=0}^{N-1} \frac{N!}{i! \prod_{k=i}^{N-1} \lambda_k} P_N + P_N = 1. \quad (10)$$

Therefore, P_N , and hence P_i 's can be solved by (9) and (10).

Let \mathcal{P} denote the set of P_i . For a given \mathcal{P} , A can be calculated by using (5). For a given P_B , \mathcal{P} can be calculated by (6)–(10). Hence, these equations can be written in the fixed-point model form $A = f_1(\mathcal{P})$ and $\mathcal{P} = g_1(A)$ [6]. A , α and \mathcal{P} can be computed by the Successive Over-Relaxation method.

The transmission cost per server in a unit time C_t is

$$C_t = \alpha \lambda_D P_N \left[1 - \left(\sum_{i=0}^r P_{N-i} \right)^{M-1} \right] C \quad (11)$$

where C is the cost of transmitting a call (video) with unit length, $\alpha \lambda_D P_N$ is the remote traffic rate from an LVS, and

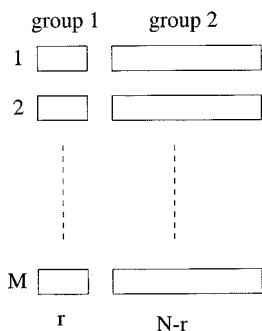


Fig. 7. Queueing model for the strict reservation scheme.

$[1 - (\sum_{i=0}^r P_{N-i})^{M-1}]$ is the probability that at least one remote LVS is not full [and hence can handle remote (video) requests].

Note that in the special case when $r = 0$, P_B is simply given by $E(M\lambda_D, MN)$, where $E(A, k)$ is the Erlang B formula with offered traffic A and k servers. Using this P_B , the transmission cost can be calculated by the above method with r set to 0.

2) *Strict Reservation*: Another reservation scheme is to reserve r servers solely for local calls. So, each LVS keeps a record of the number of overflow calls (from other queues in group 2) occupying its servers. When this number reaches $N - r$, further remote overflow calls will be blocked, even if there are free servers available. A network deploying this reservation scheme can be modeled as two groups of queues as in Fig. 7.

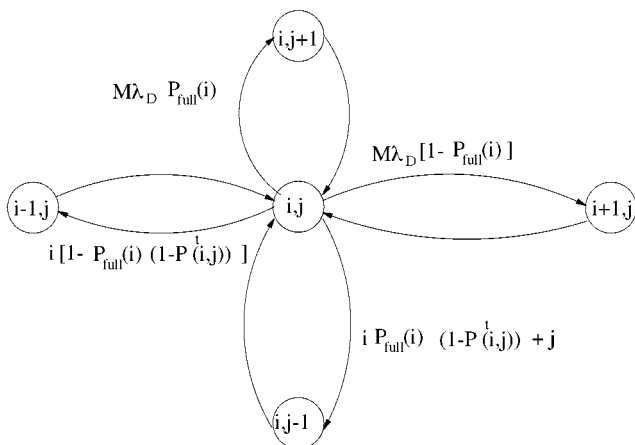
Each group consists of M multiserver queues. Each queue in group 1 has r servers, while each queue in group 2 has $N - r$ servers. It means that an LVS consists of a group 1 queue and a group 2 queue, where the group 1 queue is reserved for local traffic only, while the group 2 queue can serve both local and remote (overflow) traffic. The input traffic to group 1 represents the aggregated fresh offered traffic to the network and its rate is equal to $M\lambda_D$. When a call arrives, it is randomly directed to one of the queues in group 1. If the queue is full, it will go into its companion queue in group 2. If that queue is full as well, it will be randomly redirected to one of other available queues in group 2. If all those queues are full, the call will be blocked.

Let (i, j) represent the state of the network in which there are totally i and j calls in groups 1 and 2, respectively. Fig. 8 shows the transition rates into and out of state (i, j) .

At state (i, j) , an arrival will cause a transition to state $(i, j + 1)$ if it is directed to a full queue in group 1. The probability $P_{\text{full}}(i)$ that an arrival is directed to a full queue when there are i calls in group 1 is given by

$$P_{\text{full}}(i) = \begin{cases} 0, & 0 \leq i < r \\ \frac{C_r(i-r, M-1)}{C_r(i, M)}, & r \leq i \leq Mr. \end{cases}$$

Now, we consider the transition from state (i, j) to state $(i, j - 1)$. Two types of departure will lead to this kind of state transition. The first is simply a call departure in a group 2 queue. The other occurs when a queue in group 1 is full and its companion queue in group 2 has at least one local call; a call departure in the group 1 queue effectively causes a call departure in the group 2 queue, since we can always find a


 Fig. 8. Transition rates at state (i, j) for the strict reservation scheme.

local call to replace the departed one. Let us calculate this probability. Among the j calls in group 2, the number of local calls \hat{j} can be estimated by

$$\hat{j} = \left[\frac{\sum_{i>r} (i-r) P_{ij}}{\sum_{i>r} (i-r) P_{ij} + \sum_{i,j} j P_{ij}} \right] j,$$

where P_{ij} is the probability that in an LVS, there are i local calls and j overflow calls. The calculation of P_{ij} will be discussed in more detail later.

The probability $P^t(i, j)$ that, at state (i, j) , a queue in group 2 has no local calls at all is given by

$$P^t(i, j) = \frac{C_{N-r}(\hat{j}, M-1)}{C_{N-r}(\hat{j}, M)}.$$

Therefore, the total transition rate from state (i, j) to state $(i, j - 1)$ is $i P_{\text{full}}(i) (1 - p^t(j)) + j$. Fig. 9 depicts the state space of the network.

For a call to be blocked, all the queues in group 2 must be fully occupied. Therefore, the call blocking probability is given by

$$P_B = \sum_{i=r}^{Mr} P_{\text{full}}(i) P(i, M(N-r)). \quad (12)$$

Now, we consider the cost due to transmission of a video from remote LVSs under this reservation scheme. Again, we switch our focus to an LVS. Let (i, j) represent the state of an LVS in which there are i direct calls and j overflow (remote) calls. Note that when $0 \leq i \leq r$, j can take any value between 0 and $N - r$, but when $i > r$, j can only be between 0 and $N - i$. The state space is as shown in Fig. 10.

Let Ω_1 and Ω_2 be the sets of states in which direct calls and overflow calls can be accepted, respectively. They are given by

$$\begin{aligned} \Omega_1 &= \{(i, j) : 0 \leq i < r, \quad 0 \leq j \leq N - r\} \\ &\quad \cup \{(i, j) : r \leq i < N, \quad 0 \leq j < N - i\} \\ \Omega_2 &= \{(i, j) : 0 \leq i < r, \quad 0 \leq j < N - r\} \\ &\quad \cup \{(i, j) : r \leq i < N, \quad 0 \leq j < N - i\}. \end{aligned}$$

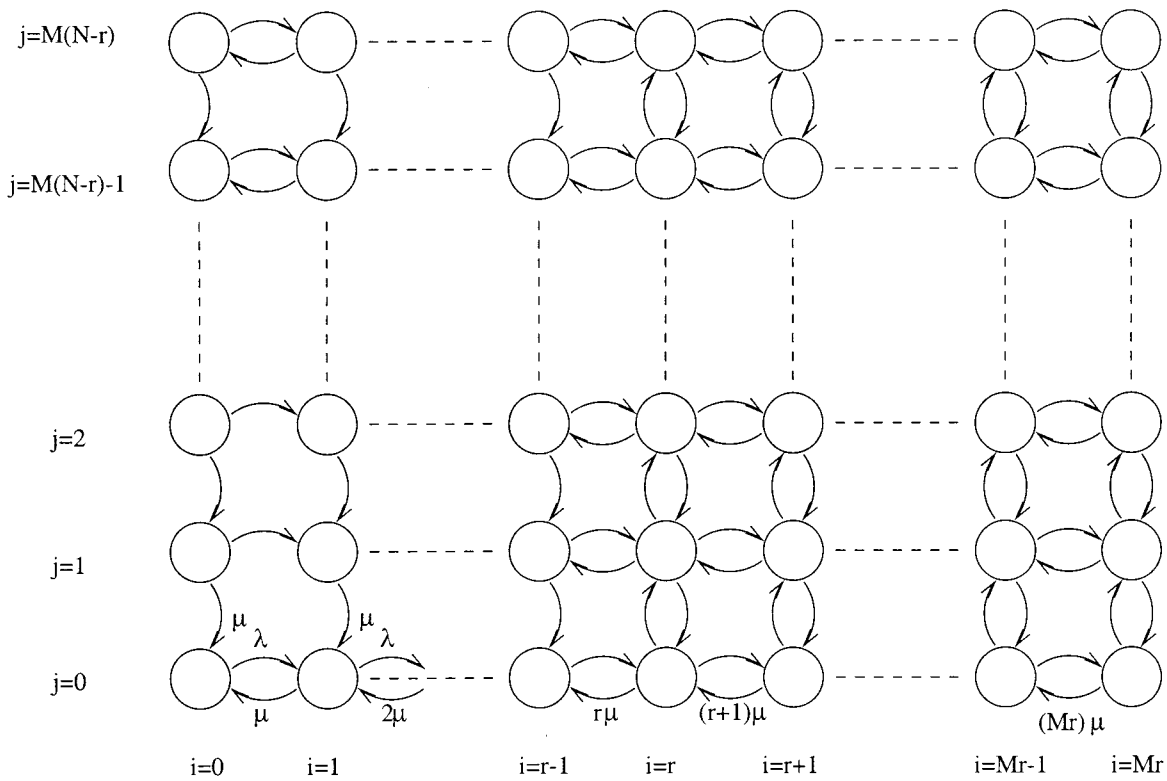


Fig. 9. State transition diagram of a network under the strict reservation scheme.

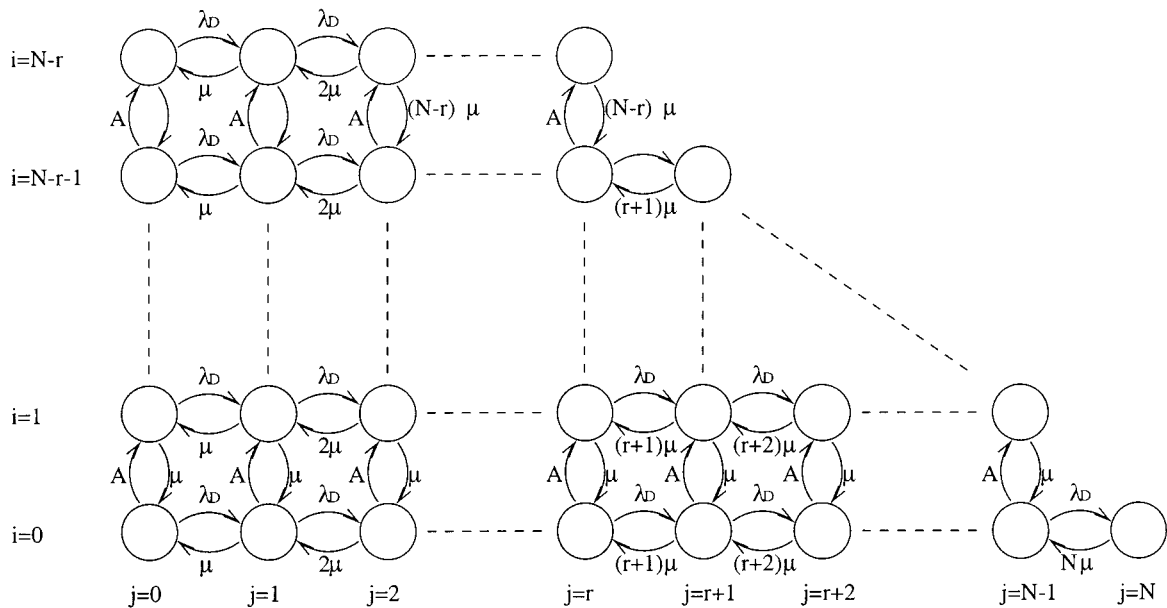


Fig. 10. State transition diagram of a single server with Random strategy and the strict reservation scheme.

As in the case of residual reservation, the total remote traffic $A_{i,j}$ to an LVS at state (i, j) is

where A_{ij} is related to λ_D by the following equation of conservation of flow:

$$A_{ij} = \begin{cases} A, & (i, j) \in \Omega_2 \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

$$\sum_{(i,j) \in \Omega_2} A_{ij} P_{ij} = \lambda_D \left(1 - \sum_{(i,j) \in \Omega_1} P_{ij} \right) \cdot \left[1 - \left(1 - \sum_{(i,j) \in \Omega_2} P_{ij} \right)^{M-1} \right]. \quad (14)$$

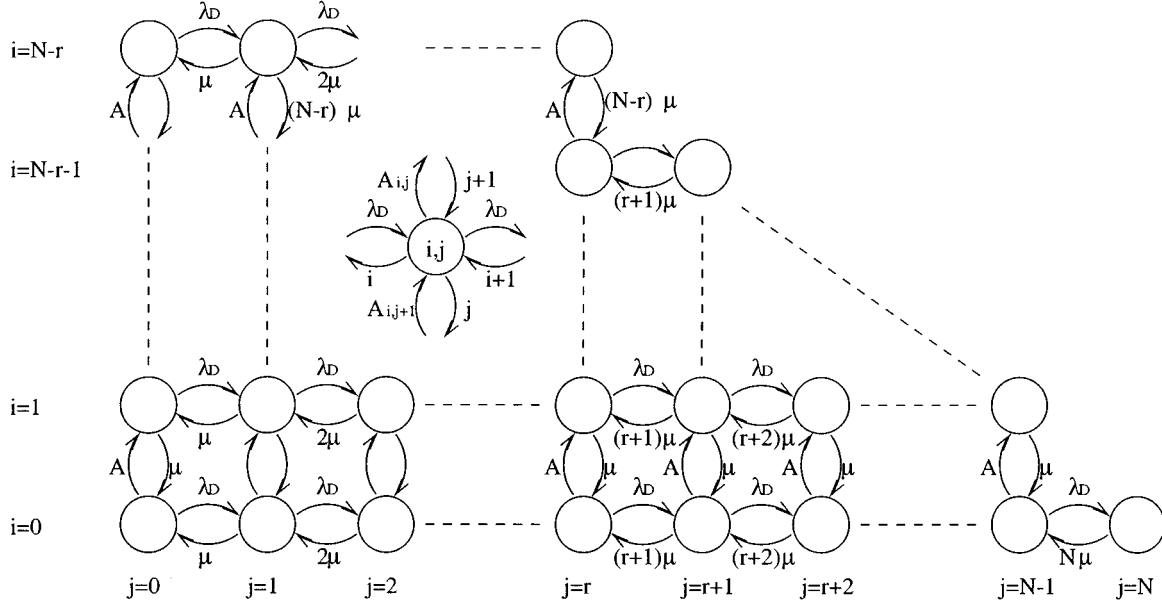


Fig. 11. State transition diagram of an LVS with Least Loaded strategy and the strict reservation scheme.

The transition rates from (i, j) to $(i + 1, j)$ is $\alpha\lambda_D$ for $(i, j) \in \Omega_1$, and that from (i, j) to $(i, j + 1)$ is αA_{ij} . Again, we scale up the actual traffic rates to match the carried traffic obtained from P_B

$$\alpha \left(\lambda_D \sum_{(i,j) \in \Omega_1} P_{ij} + \sum_{(i,j) \in \Omega_2} A_{ij} P_{ij} \right) = \lambda_D (1 - P_B). \quad (15)$$

Using (14) and (15) and the global balance equations of each state, the set of P_{ij} can be obtained by solving a system of linear equations. Then, using the same method as described in the case of residual reservation, for a given P_B obtained from the network model, the set of P_{ij} , A , and α can be solved. Once a set of P_{ij} is calculated, it can be used to calculate a new set of $P(i, j)$, and then a new P_B . This process is repeated until both P_B and all P_{ij} converge.

Therefore, C_t is given by

$$C_t = \alpha \lambda_D \left(1 - \sum_{(i,j) \in \Omega_1} P_{ij} \right) \cdot \left[1 - \left(1 - \sum_{(i,j) \in \Omega_2} P_{ij} \right)^{M-1} \right] C. \quad (16)$$

B. Least-Loaded Strategy

1) *Residual Reservation*: First, P_B is found using the same method as in Random strategy. Then, we proceed to calculate the transmission cost.

When the local LVS is full, the Least Loaded strategy will direct the call to the remote LVS with the maximum available servers. When there is more than one LVS of such kind, one is chosen at random.

Consider a particular LVS A . If this LVS is full, the overflow calls of rate λ_O will be routed randomly to one of the least-loaded (remote) LVSs. Let there be a total of β such LVSs. Then the remote load from server A that falls on a particular least-loaded LVS is λ_O/β . Let Z_k be the probability that an LVS has k or more occupied servers. Then, we have

$$Z_k = \sum_{n \geq k} P_n. \quad (17)$$

Given that remote LVS B of LVS A has k occupied servers, the probability $f(\beta | k)$ that $\beta - 1$ other remote LVSs also have k occupied servers each and each of the remaining $M - 1 - \beta$ remote LVSs has more than k occupied servers is given by

$$f(\beta | k) = \binom{M-2}{\beta-1} (Z_k - Z_{k+1})^{\beta-1} Z_{k+1}^{M-1-\beta} \quad (18)$$

where $Z_k - Z_{k+1}$ is the probability that an LVS has exactly k occupied servers, i.e., P_k . Therefore, given that LVS B has k occupied servers, the amount of traffic y_k that gets routed from LVS A to remote LVS B is

$$\begin{aligned} y_k &= \sum_{\beta=1}^{M-1} \frac{\lambda_O}{\beta} f(\beta | k) \\ &= \frac{\lambda_O}{M-1} \frac{Z_k^{M-1} - Z_{k+1}^{M-1}}{Z_k - Z_{k+1}} \\ &= \frac{\lambda_O}{M-1} \frac{Z_k^{M-1} - Z_{k+1}^{M-1}}{P_k}. \end{aligned} \quad (19)$$

Since LVS B carries the remote traffic from $M - 1$ possible LVSs, the total remote traffic A_i on LVS B is

$$A_i = \begin{cases} (M-1)y_i, & i = 0, 1, \dots, N-1-r \\ 0, & \text{otherwise} \end{cases} \quad (20)$$

where $\lambda_0 = \lambda_D P_N$.

Note that A_i is related to the carried load as follows:

$$\alpha \left(\lambda_D \sum_{i=0}^{N-1} P_i + \sum_{i=0}^{N-r-1} A_i P_i \right) = \lambda_D (1 - P_B). \quad (21)$$

When LVS B has i servers occupied, the call arrival rate λ_i and the call departure rate μ_i are

$$\begin{aligned} \lambda_i &= \lambda_D + A_i, & i &= 0, 1, \dots, N-1 \\ \mu_i &= i, & i &= 1, 2, \dots, N. \end{aligned} \quad (22)$$

Therefore, for state n , the global balance equation is given by

$$(\lambda_n + \mu_n)P_n = \lambda_{n-1}P_{n-1} + \mu_{n+1}P_{n+1} \quad (23)$$

with the understanding that $P_n = 0$ for $n < 0$ or $n > N$.

Again, using the P_B obtained from the network's model, A_i , P_i , and α can be solved using the Successive Over-Relaxation method.

C_t is given by

$$C_t = \alpha \lambda_D P_N \left(1 - \left(\sum_{i=0}^r P_{N-i} \right)^{M-1} \right) C. \quad (24)$$

In the special case when $r = 0$, P_B is again given by $E(M\lambda_D, MN)$, and using this P_B , the transmission cost can be calculated by the above method with r set to 0.

2) *Strict Reservation*: Firstly, P_B is obtained by the same method as in Random strategy.

Now, for the transmission cost, consider an LVS, the state transition diagram is as shown in Fig. 11.

The state transition diagram is very similar to that of Random strategy, except that the overflow traffic is state dependent. In state (i, j) , let $i + j = k$, and the rate of overflow calls A_{ij} is

$$\begin{aligned} A_{ij} &= (M-1) \sum_{\beta=1}^{M-1} \frac{\lambda_o}{\beta-1} \binom{M-2}{\beta-1} \left(\sum_{\substack{m+n=k \\ n < N-r}} P_{mn} \right)^{\beta-1} \\ &\quad \cdot \left(\sum_{m+n > k} P_{mn} + \sum_{0 \leq m \leq r} P_{mN-r} \right)^{M-1-\beta} \end{aligned} \quad (25)$$

with

$$\alpha \left(\lambda_D \sum_{(i,j) \in \Omega_1} P_{ij} + \sum_{(i,j) \in \Omega_2} A_{ij} P_{ij} \right) = \lambda_D (1 - P_B). \quad (26)$$

Using the same method as in Random strategy, we can solve P_B , P_{ij} , and α . The transmission cost is simply given by (16).

V. NUMERICAL RESULTS

In this section, some numerical results from our models and simulations are presented, given $M = 5$, $N = 50$ and $\lambda_D = 47.5$. The simulations are carried out using OPNET version 6.0³. First, we consider the residual reservation schemes. Fig. 12 plots the blocking probability against r . The blocking probability predicted by our model is quite close to the simulation results.

Fig. 13 plots the normalized transmission cost against r . Both simulation and analytical results show that when r is small, the transmission cost of Least Loaded strategy is less than Random strategy. This is because Least Loaded strategy always directs

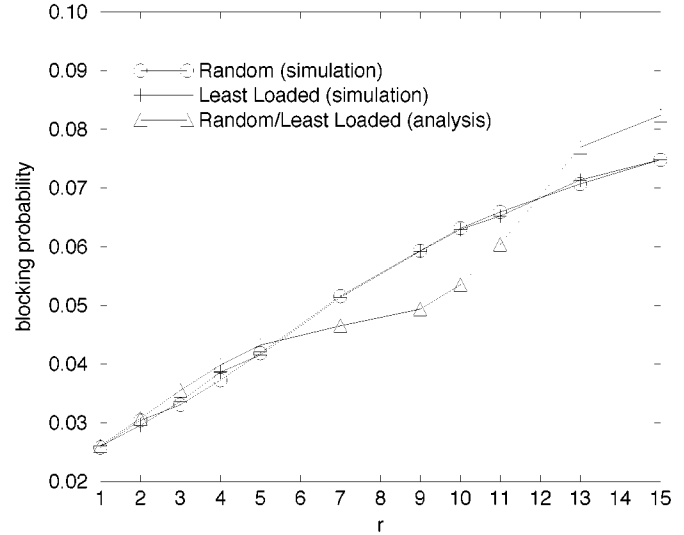


Fig. 12. Blocking probability against r for the residual reservation schemes.

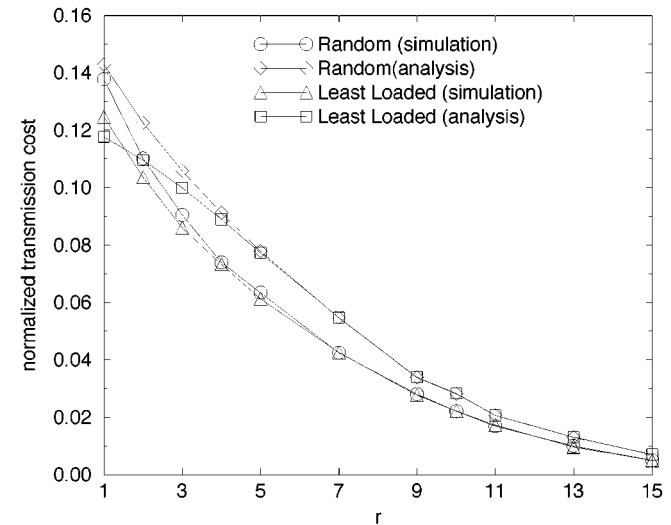


Fig. 13. Normalized transmission cost against r for the residual reservation schemes.

overflow traffic to the least-loaded LVS, and hence reduces the probability of (further) overflow. However, as r increases, the difference becomes less and finally, when r becomes large, the performance of the two strategies converges. This is because when r becomes large, each LVS is dominated by local calls and not many overflow calls would be accepted in both schemes.

Now, we consider the strict reservation schemes. Fig. 14 plots the blocking probability against r . Comparing with simulation results, it can be seen that our model quite accurately predicts the blocking probability for a wide range of r .

Fig. 15 plots the normalized transmission cost against r for the strict reservation schemes. Both simulation and analytical results show that the transmission cost of Least Loaded strategy is always less than Random strategy. Compared with Fig. 13, we can see that the residual reservation is a more effective scheme to smoothly limit/control the transmission cost (i.e., network bandwidth utilization).

³http://www.mil3.com.

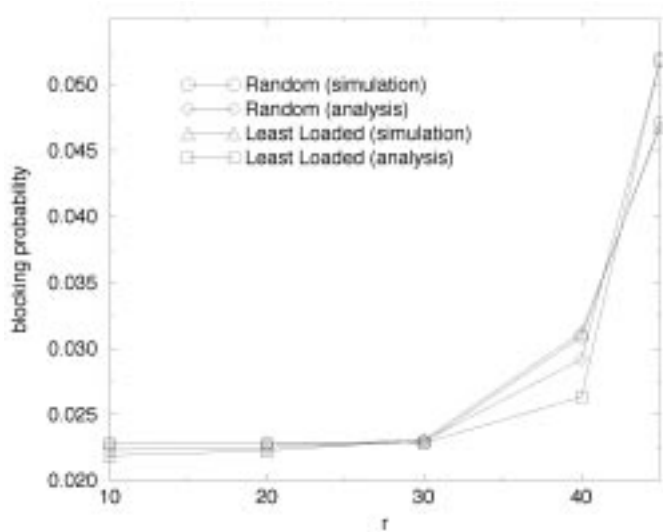


Fig. 14. Blocking probability against r for the strict reservation schemes.

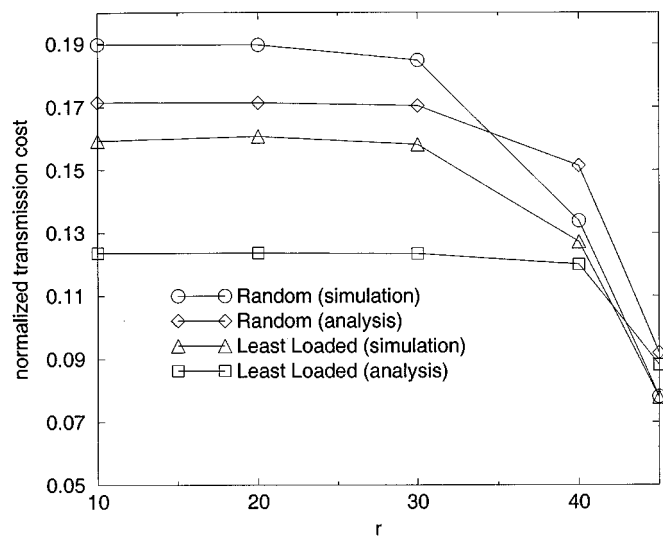


Fig. 15. Normalized transmission cost against r for the strict reservation schemes.

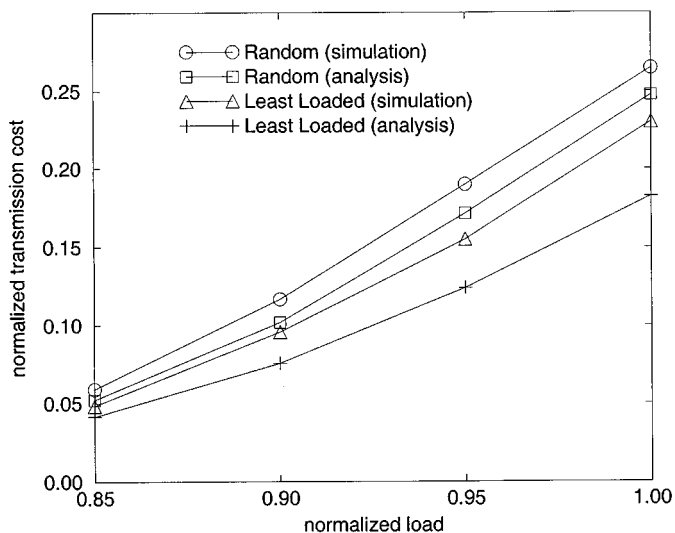


Fig. 16. Normalized transmission cost against loading.

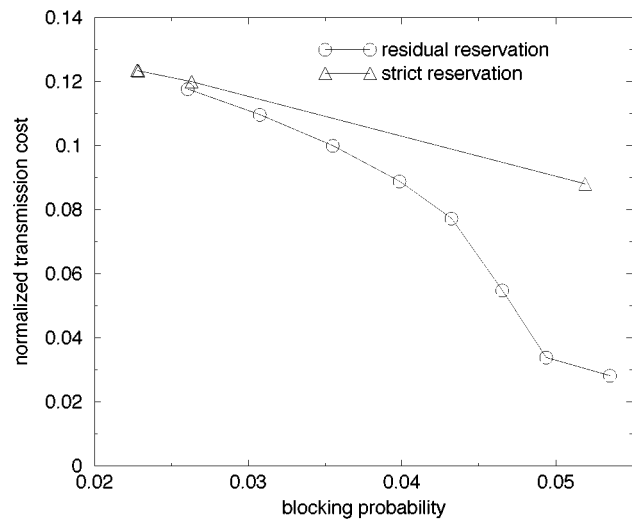


Fig. 17. Normalized transmission cost against blocking probability.

TABLE I
SPEED COMPARISON BETWEEN ANALYSIS AND SIMULATION

	analytical	simulation
$N=50, r=0$	0.06	698.64
$N=50, r=10$	22.03	557.56
$N=60, r=0$	0.02	632.32
$N=60, r=10$	31.83	589.10

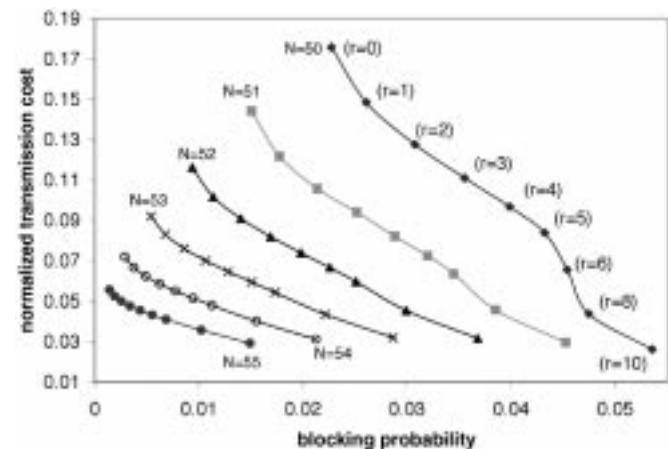


Fig. 18. Normalized transmission cost against blocking probability for different r and N .

Fig. 16 plots the normalized transmission cost for both Random and Least Loaded strategies when there is no reservation in LVs for local calls. As with reservation, the transmission cost for the Least Loaded strategy is always less than the Random strategy.

Fig. 17 plots the normalized transmission cost against blocking probability for both reservation schemes with the Least Loaded strategy. It can also be seen that, given a blocking probability, the residual reservation scheme gives less communication cost than the strict reservation scheme. From these results, we can conclude that residual reservation with Least Loaded strategy always performs better than other combinations of the reservation scheme and selection strategy. From

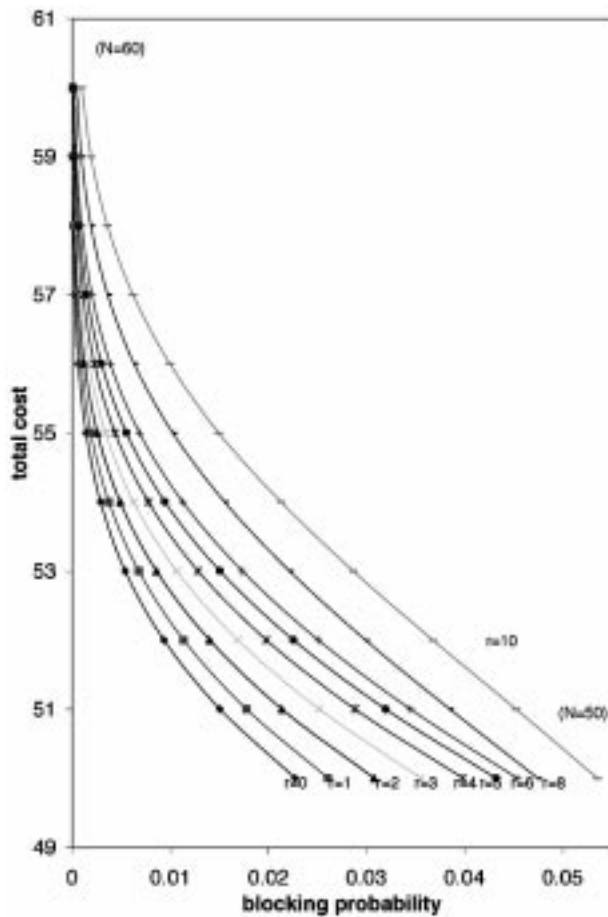


Fig. 19. Total cost against blocking probability $a = 0.01$.

now on, we focus our study on residual reservation with the Least Loaded strategy.

Although our model does not provide a closed form solution for the performance measures, it is still computationally more efficient than simulation. Table I gives a snapshot of the running times (in seconds) of two approaches using a Sun Ultra 1 machine. Obviously, solving our models requires much less time than simulation.

VI. APPLICATIONS OF THE MODELS

Our models provides a useful tool for VoD operators to design and dimension their systems. For example, using our model, the transmission cost and blocking probability can be obtained for different reservation parameter r and disk heads N , as depicted in Fig. 18. When designing a system to meet a specified quality of service requirement (blocking probability), this set of curves can be used to choose appropriate values for r and N , so that the transmission cost can be well controlled.

Alternatively, for a newly established VoD system, our model can be used to assess the relationship between total cost and the blocking probability. Here, the total cost T is modeled as

$$T = a \times \text{transmission cost} + N \quad (27)$$

where a is a scaling factor and N is a measure of the cost of a video server, assuming that the cost of a video server increases with the number of disk heads. When a is larger than 1, it means that the transmission cost is relatively higher than the server

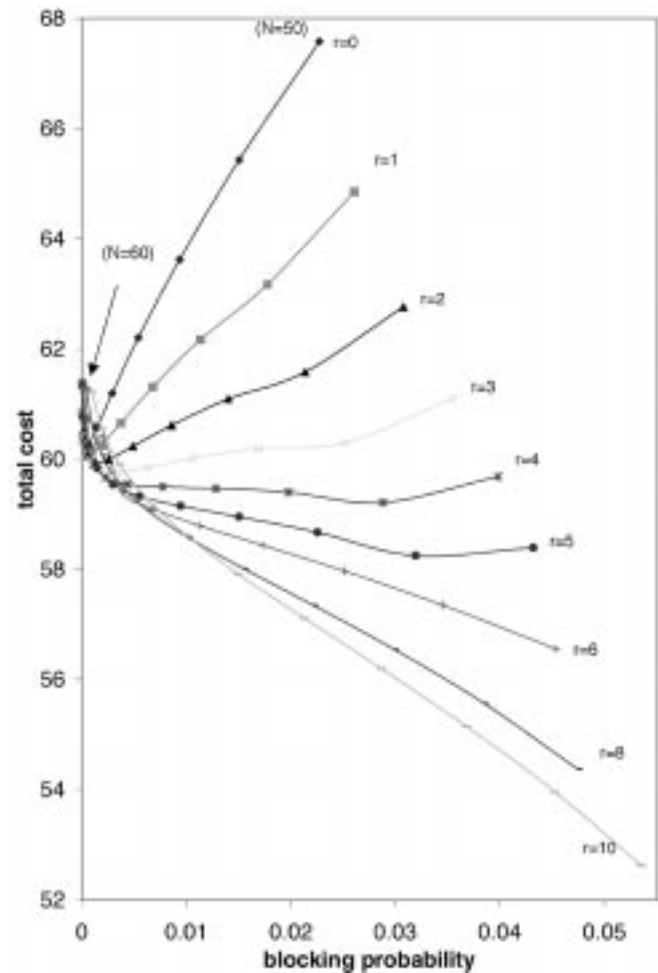


Fig. 20. Total cost against blocking probability $a = 100$.

cost. On the other hand, when a is less than 1, it means that the server cost is relatively higher than transmission cost. Figs. 19 and 20 plot the total cost against blocking probability with $a = 0.01$ and 100, respectively. Comparing these two figures, it can be seen that a can affect the relationship between total cost and blocking probability. In Fig. 19, the total cost decreases with r , while in Fig. 20, the total cost increases when r decreases. Also, in Fig. 20, for a given a blocking probability, there exists an optimal combination of N and r such that the total cost can be minimized. For example, if a blocking probability of 0.01 is required, $N = 55$ and $r = 8$ should be chosen such that total cost is minimized. In addition, through these results obtained, it can be seen that the introduction of server strategy schemes is very useful in order to reduce system cost.

VII. CONCLUSION

In this paper, a fully connected video server network architecture for VoD systems was studied. Under the assumptions of uniform loading and symmetrical network, analytical models were proposed for two server selection strategies (Random and Least Loaded) and for two reservation schemes (strict reservation and residual reservation). The models can be used to determine the call-blocking probability and the requirements of network bandwidth given the capacity of LVSs. It was shown that the models

TABLE II
DEFINITION OF MAJOR SYMBOLS

Symbols	Definition
M	the number of LVSs in the system
N	the number of server/disk heads in an LVS
r	reservation parameter
(i, j)	the state of the network in which there are totally i and j in groups 1 and 2, respectively
$P(i, j)$	the probability that the network is in state (i, j)
λ_D	the rate of local offered traffic to an LVS
λ_O	the rate of overflow traffic from an LVS
$C_N(i, M)$	the number of distinct occupancies when inserting i distinct balls into M urns where each urn has the capacity to store N balls
$P_{\text{full}}(i)$	the probability that an arrival is directed to a full queue when there are i calls in group 1
$P'_{\text{full}}(i, j)$	the probability that a call entering group 2 finds a full queue when the network is state (i, j)
$P^t(i, j)$	the probability that a queue in group 2 has no local calls at state (i, j)
P_B	the system blocking probability
C_t	the transmission cost (per server) on the fully connected core network in a unit time

are, in general, quite accurate when compared with the simulation results. Numerical results showed that: 1) the Least Loaded strategy always provides less communication cost than Random strategy as expected; 2) given the same blocking probability, the residual reservation scheme gives less communication cost than the strict reservation scheme; and 3) the residual reservation is a more effective scheme to smoothly limit/control network bandwidth utilization than the strict reservation scheme. In addition, the models provide a useful tool for VoD operators to design and dimension their systems. For example, when designing a system to meet a specified quality of service requirement (blocking probability), a set of system parameters can be obtained by using the models such that the total system cost is minimized.

APPENDIX

Please see Table II.

ACKNOWLEDGMENT

The authors are grateful for the valuable suggestions from Professor P. T. S. Yum.

REFERENCES

- [1] J. Sutherland and L. Litteral, "Residential video services," *IEEE Commun. Mag.*, pp. 36–41, July 1992.
- [2] Y. H. Chang *et al.*, "An open-systems approach to video on demand," *IEEE Commun. Mag.*, pp. 68–80, May 1994.
- [3] S. A. Barnett and G. J. Anido, "A cost comparison of distributed and centralized approaches to video-on-demand," *IEEE J. Select. Areas Commun.*, vol. 14, pp. 1173–1183, Aug. 1996.
- [4] D. Deloddere, W. Verbiest, and H. Verhille, "Interactive video on demand," *IEEE Commun. Mag.*, pp. 155–162, May 1994.
- [5] C. H. Sauer and E. A. MacNair, *Simulation of Computer Communication Systems*. Englewood Cliffs, NJ: Prentice-Hall, 1993.

- [6] K. W. Ross, *Multiservice Loss Models for Broadband Telecommunication Networks*. New York: Springer-Verlag, 1995, ch. 7.
- [7] I. Chlamtac and A. Ganz, "Design and analysis of very high-speed network architectures," *IEEE Trans. Commun.*, vol. 36, pp. 252–262, Mar. 1988.



Eric W. M. Wong (S'87–M'90–SM'00) received the B.Sc. and M.Phil. degrees in electronic engineering from the Chinese University of Hong Kong, Hong Kong, in 1988 and 1990, respectively, and the Ph.D. degree in electrical and computer engineering from the University of Massachusetts, Amherst, in 1994.

He joined the City University of Hong Kong as an Assistant Professor in the Department of Electronic Engineering in 1994. His research interests are in high-speed networks, video-on-demand, satellite communications, and dynamic routing. The most

notable of these involved the analytical modeling of the least loaded routing scheme in circuit-switched networks. The model drastically reduces the computational complexity of designing and dimensioning telephone systems. The model has also served as a core for the analysis of many other advanced routing schemes, such as the Real-Time Network Routing currently implemented in the AT&T telephone network.



Sammy C. H. Chan (M'90) received the B.E. and M.Eng.Sc. degrees in electrical engineering from the University of Melbourne, Melbourne, Australia, in 1988 and 1990, respectively, and the Ph.D. degree in communication engineering from the Royal Melbourne Institute of Technology, Australia, in 1995.

From 1989 to 1994, he was with Telecom Australia Research Laboratories, first as a Research Engineer, and between 1992 and 1994 as a Senior Research Engineer and Project Leader. In December 1994, he joined the Department of Electronic

Engineering, City University of Hong Kong, where he is currently an Associate Professor.