

VoD QAM Resource Allocation Algorithms

Jiong Gong¹, David Reed¹, Terry Shaw¹, Daniel Vivanco¹ and Jim Martin²

¹ Cable Television Laboratories, Inc.
858 Coal Creek Circle
Louisville, CO 80027

j.gong@cablelabs.com, d.reed@cablelabs.com, t.shaw@cablelabs.com, d.vivanco@cablelabs.com

²Department of Computer Science
Clemson University, USA
jim.martin@cs.clemson.edu

Abstract. This paper proposes a new Quadrature Amplitude Modulation (QAM) resource allocation algorithm for Video on Demand (VoD) when there is a mixture of standard definition (SD) and high definition (HD) video streams. We have developed a simulation model to compare this algorithm with two popular algorithms: the least-loaded algorithm and the most-loaded algorithm. We show that our algorithm, which we call the non-mixing algorithm, performs significantly better than the two existing algorithms by accommodating more streams thereby lowering the blocking probabilities under a range of assumptions of peak concurrent usage rate and percentage of HD streams. Using computer simulation we found that the non-mixing algorithm leads to an average of 4.39% higher allowed peak usage rate than the least-loaded and most-loaded algorithms.

Keywords— VoD, HFC networks, Broadband access, Capacity planning, Congestion control, Traffic management & control, Traffic modeling & characterization, Resource allocation, Network modeling & simulation.

1 Introduction

Video on Demand (VoD) systems over broadband access networks are likely to see a significant change in usage patterns over the next few years. The percentage of high definition (HD) VoD stream requests is likely to increase significantly from zero to approximately 10%, and peak usage is likely to increase significantly from the current average of 5% to approximately 30% as subscription-based VoD (sVoD) and digital video recorder (DVR) applications become more mainstream¹ [6]. Cable operators will require a detailed understanding of the impact of these changes in the provisioning process.

When a cable subscriber purchases a VoD selection, the video stream is assigned to a QAM modulator over a specified 6 MHz RF channel. The encoding rate of the stream along with the specific QAM configuration determines the aggregate number of streams that can be assigned to the channel. For example over a 256 QAM modulated channel, if all content is in SD format (i.e., MPEG2) and is encoded at a constant bit rate of 3.75 Mbps, 10 streams can be assigned to the same channel and thus all of the channel bandwidth is used. A VoD system, referred to as a **service group**, consists of content servers, a delivery network, a number of QAM modulators and a set of subscribers. During the purchase of a VoD selection, the resource allocation algorithm must assign a new stream to one of the modulators in the service group. If the channel capacity is an integral multiple of the bandwidth consumed by an SD flow, the QAM resource allocation algorithm is trivial. However in future VoD systems, there will be a mix of standard definition (SD) streams and HD streams.

¹ From a commercial North American Cable Operator's market forecast. One cable operator is lately seeing close to 10% peak usage rate after the introduction of sVoD service.

A common encoding rate for HD streams is 12.5 Mbps. Assuming a channel capacity of 37.5 Mbps based on 256 QAM modulation, three HD streams would completely fill a channel. The difficulty comes when a combination of SD and HD streams are assigned to a channel. In this case, some amount of the channel bandwidth will be unused. The worst case percentage of stranded bandwidth (B_s) is $B_s = \left(\frac{r_h - r_s}{Q} \right)$, where r_s and r_h denote the streaming bit rate for SD and HD streams respectively, and Q is the channel capacity [2]. In the worst case, each QAM modulator has just under r_h bandwidth stranded. This could occur if a series of HD stream requests arrive that almost fills the QAM (i.e., to the point where one more HD request would completely fill the QAM), but then an SD stream request arrives and gets allocated. For the 256 QAM scenario described above, up to 23.3% of the channel bandwidth could be stranded.

The current prevailing QAM allocation methods include two algorithms; one that allocates incoming streams starting from the lightest-loaded QAM modulator and one that starts from the busiest-loaded QAM modulator. In the rest of the paper, we refer to the former as the “least-loaded” algorithm and the latter as the “most-loaded” algorithm. It is generally believed that the most-loaded algorithm performs better than the least-loaded algorithm when there is the presence of HD VoD streams, a fact that is confirmed in our analysis. We propose and evaluate a new QAM resource scheduling algorithm called the “non-mixing” algorithm. In this paper, we present the results of a simulation-based analysis that suggest that the non-mixing algorithm can allow peak usage rates 4.39% higher than most-loaded algorithm. A further contribution of this paper is the results of a VoD usage modeling effort which was necessary to exercise our simulation model in a realistic manner.

This paper is organized as follows. Related work is presented in section 2. The VoD usage model and the proposed non-mixing algorithm are presented on section 3 and 4, respectively. In sections 5 and 6 we present our analysis methodology and simulation-based results, respectively. Finally section 7 presents the conclusion of the analysis and identifies future work items.

2 Related Work

A large amount of prior research has addressed the scalability of large-scale VoD systems. Techniques have been identified that reduce the resources that are required per session. Batching requires users to wait in a group for the same content for a predetermined amount of time and then serves them in a batch using a single multicast channel [4] [5] [1]. Periodic broadcasting schedules the transmission of content over multiple channels in periodic intervals allowing arriving users to join the next cycle [3] [12] [7]. Patching attempts to merge users who are on separate channels to an existing multicast channel [11] [13]. Piggybacking merges users on separate channels by slightly changing playback rates of users in an effort to have everyone get to the same point in the stream at which time the separate channels would be exchanged for a single multicast channel [8] [14]. While these ideas are likely to be relevant in future cable VoD systems, most current deployments are relatively small in scale. Provisioning the optimal number of QAM modulators in a VoD service set is generally based on the rule of thumb that says about 5% of the total subscriber population will use VoD during peak periods. There has been industry discussion on QAM allocation algorithms [10] [9]. However, to the best of our knowledge, there has not been an academic evaluation of QAM resource allocation algorithms.

The QAM allocation problem is essentially a **bin packing** problem. The classic bin packing algorithm packs a list of items $L = (a_1, a_2, \dots, a_n), a_i \in (0,1]$ for all i , into the minimum number of bins each with a capacity of 1. The least loaded QAM allocation algorithm is a form of best fit packing and the most loaded allocation algorithm is a form of worst fit packing [2]. In brief, a best fit packing algorithm selects the bin that has the most free space and the worst fit algorithm selects the bin that has the least free space. The standard metric that is used to evaluate bin packing algorithms is a measure of the number of bins that are required to pack various input lists. The ratio of the number of bins required by the algorithm under study to the number of bins required by an optimal algorithm (i.e., an off line algorithm) is known as the R value.

It has been shown that both the best fit and worst fit algorithms have an R value of 2 [2]. In the QAM allocation problem domain, the number of bins is fixed. Items in bins may leave after an amount of time (i.e., when the subscriber finishes watching the movie the stream is removed from the QAM). Rather than use the R metric, we are interested in the probability that a stream's request is denied due to insufficient capacity. We use the blocking rate to characterize allocation algorithm performance.

3 VoD Model

We have developed a model of VoD usage based on empirical data. The data used in this study were collected from 200 service groups of a large cable operator in North America. The average size of each service group is approximately 500 set-top boxes. Figure 1 illustrates diurnal average usage patterns over the course of one week for all requests. The results show higher usage rate values for Thursday, Friday and Saturday evening from 10:00pm until midnight. Note that the maximum 2% VoD usage rate shown in Figure 1 was the average over all 200 service groups analyzed, while some service groups exhibited peak usage rates close to 5%.

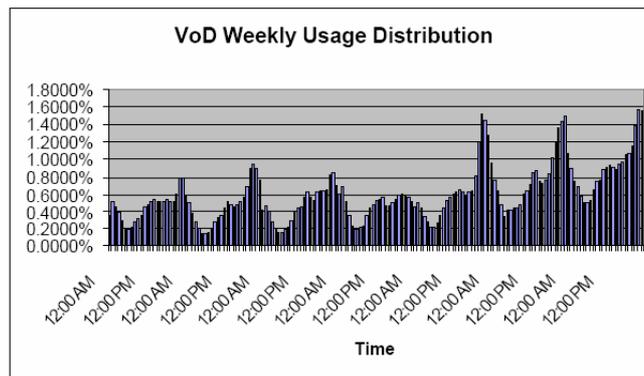


Fig. 1. Weekly VoD Usage (Sunday 12:00am through Sat 12:00am)

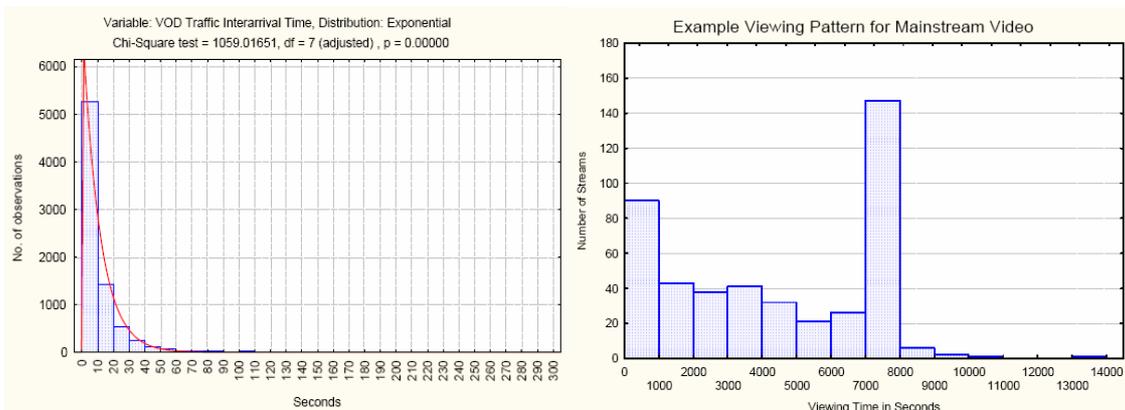


Fig. 2. Histogram and Fitted Exponential Distribution of VoD Request Interarrival Times

Fig. 3. Distribution of Stream Length for Mainstream Video Titles

We modeled the interarrival times of VoD request streams and their duration. Our results indicate that interarrival times follow an exponential distribution, although each of the main genres of content, including mainstream movies, adult content and video browsing have different fittings. Video browsing refers to short-lived streams mainly generated by a sVoD user who browses the available VoD channels available in his/her subscription package. Figure 2 shows the fit of interarrival VoD request times associated with 7800 instances of stream arrivals observed in the data set. The x-axis represents 10-second windows over a period of 24 hours. The solid line is a fitted exponential distribution curve with a λ of 0.091.

We have also modeled the viewing time distribution of mainstream movies. Our results suggest a mixed probability distribution. As illustrated in Figure 3, there was a significant mode at a viewing time of 2 hours which is the average length of mainstream movies. The data suggests a large number of early exits, which can be associated with video browsing generated by sVoD users.

4 Non-Mixing Algorithm

In this section we describe a new QAM allocation algorithm called the non-mixing algorithm. We start by describing a mathematical framework to model the problem. Suppose a collection of n QAM modulators is deployed to serve a VoD service group. Let q_i , $i = 1, 2, \dots, n$, denote the used capacity of each QAM modulator i . Total capacity, Q , which is usually 37.5 Mbps for a 256 QAM, is assumed to be the same for all QAM modulators. Therefore, the remaining capacity that can be used for new stream requests on that QAM modulator is then $Q - q_i$. Let r_s and r_h denote the streaming bit rate, respectively, for SD and HD streams. The two types of streams may arrive at a collection of QAM resources according to two distinct random processes, such as the Poisson process, but exit the system based on the same holding time distribution. We call the current state of any given QAM (q_i) at a particular time as an allocation. We define an allocation as inefficient, if,

$$Q - q_i < r_h, \forall i, \text{ and } \sum Q - q_i \geq r_h \quad (1)$$

In other words, none of the QAM modulators individually has the capacity, even though the sum of all available resources on each QAM modulator is able to support one or more HD stream requests. A better scheduling algorithm would generate fewer cases of inefficient allocations. Note that while each type of stream is assumed to be in itself modulus in its own bit rate, they jointly are not when they are mixed together in a QAM modulator. As a result, inefficiency tends to arise when different stream types are mixed together. Both most-loaded and least-loaded algorithms lead to mixed allocations at the QAM modulators. In the following lines the non-mixing algorithm is going to be presented. Let's first start by defining 4 possible states for any QAM on the system at any given time, depending on its current allocation;

- No streams have been allocated.
- A mixture of SD and HD streams are occupying it.
- Only SD streams are occupying it.
- Only HD streams are occupying it.

Mathematically, we denote these four types accordingly by defining a state function as:

$$S_i(q_i) = \begin{cases} 1, & \text{if } q_i = 0 \\ 2, & \text{if } q_i = x_i r_s + y_i r_h, x_i \neq 0, y_i \neq 0 \\ 3, & \text{if } q_i = x_i r_s, x_i \neq 0 \\ 4, & \text{if } q_i = y_i r_h, y_i \neq 0 \end{cases} \quad (2)$$

where x_i and y_i are positive integers representing the number of SD and HD streams, respectively, occupying QAM modulator i . In the above four states, we call a QAM modulator in state 1 an empty QAM

modulator. We call a QAM modulator in state 2, that is $S_i(q_i) = 2$, a mixing QAM modulator. QAM modulators in state 3 and 4 are called non-mixing SD and HD QAM modulators, respectively. The algorithm selects a QAM using the following prioritized rules:

- Select a non-mixing QAM modulator of the same stream type.
- Select an empty QAM modulator.
- Select a mixing QAM modulator.
- The last resort is to create another mixing QAM modulator by selecting an existing QAM that currently has only SD or only HD streams.

If there are multiple QAM modulators available within the same state class, priority is given to those QAM modulators that have a larger likelihood of becoming a non-mixing QAM modulator or an empty QAM modulator once some streams start to drop. This implies the following rules:

- If multiple non-mixing QAM modulators are available to a stream request of the same stream type, priority should be given to the busiest non-mixing QAM modulator because other mixing QAM modulators have a higher likelihood of being non-mixing or empty.
- If multiple mixing QAM modulators are available to a SD or HD stream request, priority is given to the busiest mixing QAM modulator, because other mixing QAM modulators have a higher likelihood of being non-mixing or empty.
- If multiple non-mixing QAM modulators are available to a stream request of a different type, that is if a stream request will have to create a new mixing QAM modulator, priority is given to the least busy QAM modulator, because it has the highest likelihood of becoming non-mixing again.

Refer to Appendix A for further details of the algorithm.

5 Analysis Methodology

5.1 Simulation Model

We developed a simulation model with which we can evaluate the performance of a set of QAM allocation algorithms and also be used as a capacity planning tool for cable operators. The model simulates a pool of 256 QAM modulators in a VoD service group. Session requests are either SD or HD streams. SD and HD stream requests have been modeled as independent Poisson processes with interarrival times exponentially distributed. The aggregate stream request is the combination of the SD and HD streams, which also follows a Poisson process with interarrival times exponentially distributed [14]. Equation 3 shows the relationship used to calculate the aggregate VoD request interarrival rate, λ , based on the number of users in a service group and the aggregated concurrent usage rate during the peak hour.

$$\lambda = (\text{Number_user}) * (\text{Peak_usage_rate} / 3600) \quad (3)$$

Since the peak-usage rate is defined as the maximum number of stream requests during the peak one hour time period, this parameter was converted from hours into seconds. Equations 4 and 5 represent the SD and HD mean interarrival rates, respectively.

$$\lambda_{SD} = (\text{Percentage_SD_streams}) * \lambda \quad (4)$$

$$\lambda_{HD} = (\text{Percentage_HD_streams}) * \lambda \quad (5)$$

Arrival requests have been already classified in section 3 in three genres; mainstream movies, adult content and video browsing, and each of them is characterized by their own unique average duration time. Stream durations for each of these genres have been modeled as independent random variables distributed

exponentially. This conclusion has been found from the empirical data shown in Figure 3. Note that this figure shows the aggregate stream duration distribution, thus this is the aggregation of three exponential distributions with different average duration times. Equation 6 shows the aggregate stream duration, μ , based on the weighted average based on the proportions of the genres that make up the streams, where m represents the number of stream types (in this case $m=3$).

$$\mu = \sum_{j=1}^m (\text{Percentage_movie_type}_j) * (\text{Average_duration_movie_type}_j) \quad (6)$$

The proposed VoD simulation engine presented in this paper replicates a real-world stream processing experience as follows;

- Accepted streams are released from the QAM modulator when their duration expires.
- Incoming stream requests are compared with the available QAM capacity.
 - If the available capacity is insufficient to handle the request, it is denied and the number of sessions rejected count is incremented by one.
 - If the request is accepted it is placed in an empty channel of one of the available QAM modulators. The channel selection is determined by the stream allocation algorithm that has been configured.

Three allocation algorithms are implemented in the simulation model: least-loaded, most-loaded, and non-mixing. In the most-loaded algorithm, the available QAM capacity remaining within a service group is placed in an array and sorted from the lowest to the highest. The QAM modulator that has the smallest remaining capacity represents the most-loaded or busiest QAM modulator. The incoming stream request is assigned to the most-loaded QAM modulator with sufficient capacity to handle it. In the least-loaded algorithm the reverse occurs, arriving requests are assigned to the QAM modulator that has the largest remaining capacity enough to handle the request. In the non-mixing algorithm, the available QAM capacities are grouped in virtual clusters. The incoming stream is assigned to a QAM channel according to the rules described in section 4.

5.2 Assumptions

The model we developed can accommodate a great number of scenarios depending on the streaming bit rates, the size of the service group, the precise mixture of SD and HD streams and other factors. Table 1 and 2 show the system level assumptions and the stream characteristic assumptions, respectively. The values presented in these tables were obtained from current deployments and real usage VoD patterns data.

System Level Assumptions	
Modulation Technique	256-QAM
SD Bit Rate	3.75 Mbps
HD Bit Rate	12.5 Mbps
Channel Capacity	37.5 Mbps
Number of users on Service Group	500

Table 1. System Level Assumptions.

Stream Characteristics Assumptions			
	Percentage of movie type in SD streams	Percentage of movie type in HD streams	Average Duration
Mainstream Movies	40%	57%	2 hours
Adult Movies	30%	-	20 minutes
Browsing stream	30%	43%	15 minutes

Table 2. Stream Characteristics Assumptions.

6 Results

The performance of the stream allocation algorithms was measured by calculating the average blocking probability first. The analysis varies the peak usage rate, SD and HD stream composition percentages and the QAM pool size in a service group.

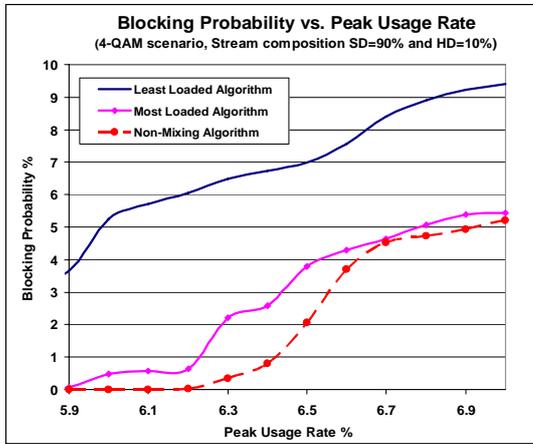


Fig.4.a.

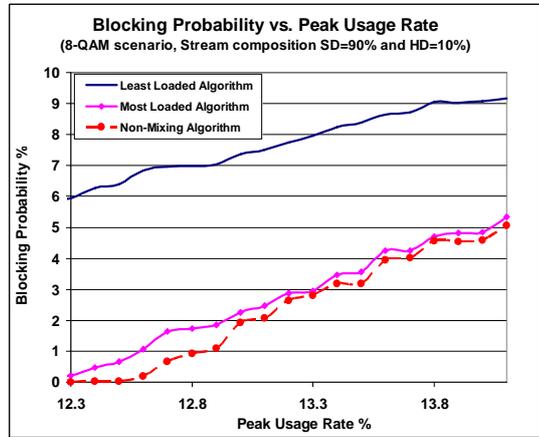


Fig.4.b.

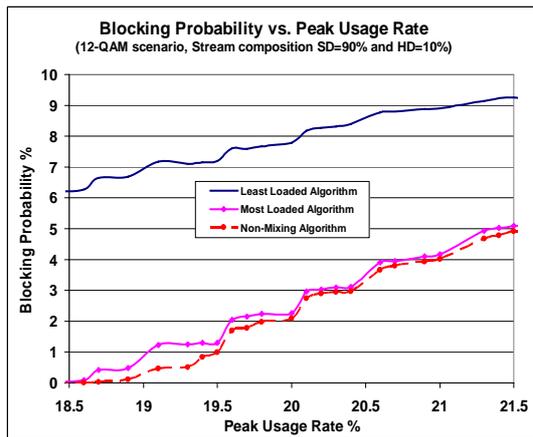


Fig.4.c.

Fig. 4. Blocking Probability vs. Peak-Usage Rate for Least-Loaded, Most-Loaded and Non-Mixing QAM Allocation Algorithms for 90% SD and 10% HD Streams, (a) 4 QAM scenario, (b) 8 QAM scenario, (c) 12 QAM scenario.

Figures 4.a, 4.b and 4.c show the blocking probability for the three mentioned algorithms against a range of peak-usage rates for systems with 4, 8 and 12 QAM modulators, respectively, for the case where the traffic consists of 90% SD streams and 10% HD streams. Figures 5.a, 5.b and 5.c show similar results for the case where the traffic consists of 70% SD streams and 30% HD streams. From these results, it can be seen that non-mixing allocation algorithm leads to a lower blocking probability than the other two algorithms at all usage levels. Filling a QAM modulator with only one type of stream can guarantee maximum capacity utilization given the modular nature of the streaming bit rates. On the other hand, a mixing QAM modulator is likely to have stranded bandwidth that is not sufficient to accommodate an incoming HD stream. Figures 4 and 5 also indicate the poor ability of the least-loaded algorithm to efficiently allocate streams on congested VoD systems. Figure 4.a illustrates that in a VoD system consisting of 10% HD content with 4 QAM modulators, and under 6% peak usage level, the most-loaded and the non-mixing algorithm lead to a blocking probability close to 0%, while the least-loaded algorithm results in a blocking probability close to 4%.

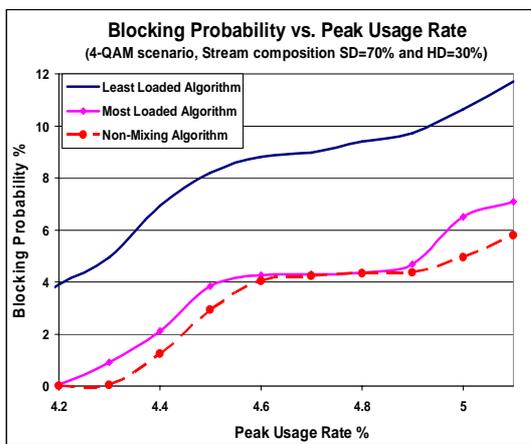


Fig.5.a.

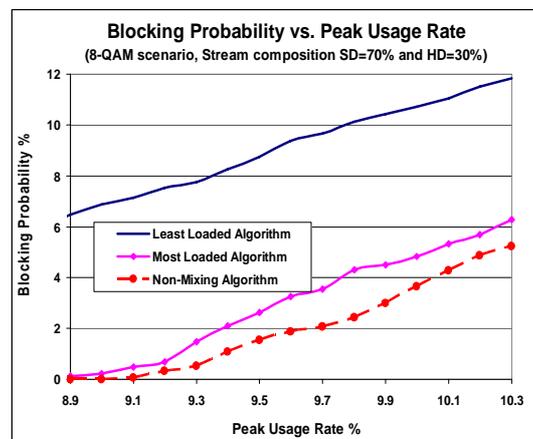


Fig.5.b.

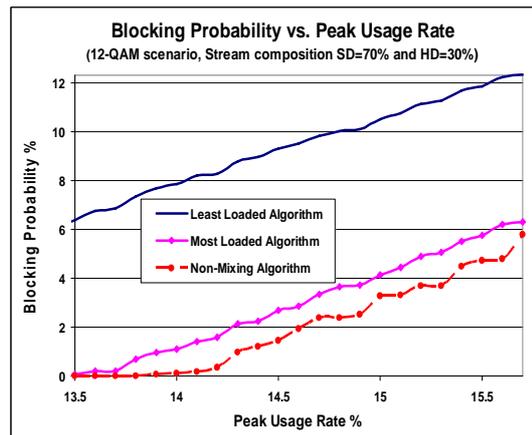


Fig.5.c.

Fig. 5. Blocking Probability vs. Peak-Usage Rate for Least-Loaded, Most-Loaded and Non-Mixing QAM Allocation Algorithms for 70% SD and 30% HD Streams, (a) 4 QAM scenario, (b) 8 QAM scenario, (c) 12 QAM scenario.

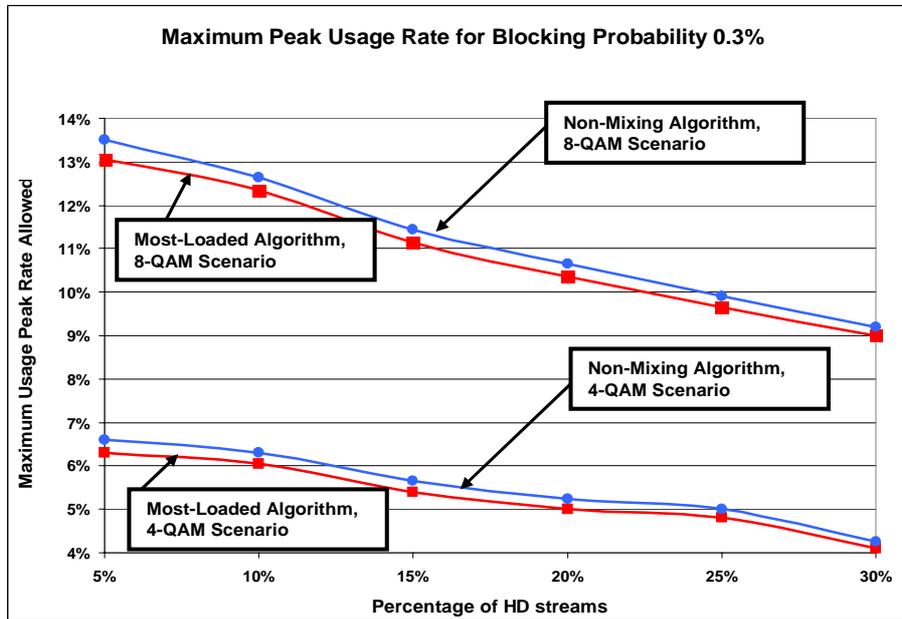


Fig. 6. Maximum Peak-Usage Rate Allowed vs. Percentage of HD streams using No-Mixing and Most-Loaded Algorithm for 4 and 8-QAM Scenarios

Figure 6 shows the maximum peak-usage rate that can be supported to meet a 0.3% blocking probability objective in a 4 QAM and 8 QAM VoD systems as a function of the percentage of HD streams. The percentage of capacity improvement of the non-mixing algorithm over the most-loaded algorithms ranges between 3.66% to 5% for the 4 QAM scenario, and between 2.22% to 3.45% for the 8 QAM scenario. For the 4 QAM and 8 QAM scenario an average of 4.39% and 2.71%, respectively, higher allowed peak usage rate can be perceived. As the traffic load increases, it becomes more difficult for the non-mixing algorithm to keep QAM modulators non-mixed. In this case, the non-mixing algorithm has a tendency to behave like the most-loaded algorithm.

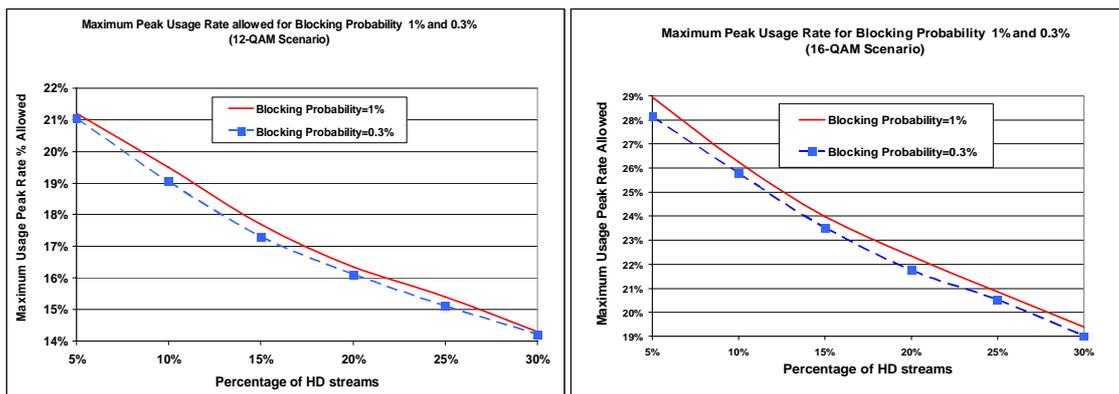


Fig.7.a.

Fig.7.b.

Fig. 7. Maximum Peak-Usage Rate Allowed vs. Percentage of HD streams using No-Mixing Algorithm, (a) 12-QAM Scenario, (b) 16-QAM Scenario.

Figures 7.a and 7.b show that the maximum peak-usage for 12 and 16 QAM systems respectively when subject to blocking rate objectives of 0.3% and 1%. Most providers would consider a blocking rate of 0.3% acceptable and a 1% rate marginally acceptable. These figures suggest that the maximum peak-usage rate that can be supported to achieve blocking probability objectives decays as the percentage of HD streams increase. To demonstrate how the results from Figure 7 might be used for provisioning, assume that a hypothetical VoD system will experience a peak-usage rate of 20%. For this load, none of the 4 QAM or 8 QAM systems for the 500 home service group can support this volume of traffic, regardless of the stream composition (see Figures 4 and 5). Figure 7.a suggests that a 12 QAM system could handle this load as long as the traffic mix contains less than 7.5% HD streams. Figure 7.b suggests that a 16 QAM system could handle this load for traffic that includes up to 27% HD streams.

7 Conclusions

Our results highlight the effect that QAM allocation algorithms can have on the efficiency of a VoD system. The two commonly deployed algorithms, least-loaded and most-loaded, are designed for current generation VoD systems that offer only SD streams. Future systems will involve a mix of SD and HD streams. We have shown that a least-loaded algorithm can result in more than a five fold increase in blocking probability compared to a most-loaded algorithm when subject to varying levels of SD and HD stream requests. Our proposed algorithm, the non-mixing algorithm, is able to demonstrate better performance under all cases of usage level and under all cases of HD percentage assumptions.

Many VoD systems are deployed using 4 QAM modulators in a service group. Our analysis shows that changing to the non-mixing algorithm can support up to 6.2% peak-hour concurrent usage that contains 10% HD streams, which is difficult to be accommodated with most-loaded or least-loaded algorithms (see figure 4.a). With more HD content, the non-mixing algorithm can generate an average of 4.39% higher allowed peak usage rate over most-loaded algorithm. 6.2% seems to be a reasonable peak-hour concurrent usage assumption in the near term for many VoD systems in North America that are currently experiencing peak-hour concurrent usage below 5%.

The benefits of the non-mixing algorithm over the most-loaded algorithm depend primarily on its ability to avoid, to the extent possible, mixing SD and HD streams. This is driven by several factors, including SD and HD traffic composition, SD and HD streaming bit rates, traffic load, and many others. The benefits do not appear to significantly depend on the number of QAM modulators in the system. However, the number of QAM modulators that are needed to meet a blocking probability objective is highly dependent on the percentage of HD streams in the traffic mix. Future work includes the evaluation of the allocation algorithms in systems that have VoD, switched digital broadcast and high speed data (DOCSIS) [16] traffic using the same set of QAM resources. We also plan to develop models and tools that can be used for capacity planning in the next generation cable systems.

Reference:

1. C. Aggarwal, J. Wolf, P. Yu, "On Optimal Batching Policies for Video-on-Demand Server", ACM International Conference on Multimedia Systems, pp. 253-258, June 1996.
2. E. Coffman, M. Garey, D. Johnson, "Approximation Algorithms for Bin Packing: A Survey", Approximation Algorithms for NP-hard Problems, pp 46-89, PWS Publishing Company, 1995.
3. T. Chiueh, C. Lu, "A Periodic Broadcasting Approach to Video-on-Demand Service", Proc. SPIE, vol 2615, pp. 162-169, 1996.
4. A. Dan, D. Sitaram, P. Shahabuddin, "Scheduling Policies for an On-demand Video Server with Batching", ACM International Conference on Multimedia, pp. 15-23 1994.
5. A. Dan, D. Sitaram, P. Shahabuddin, "Dynamic Batching Policies for an On-demand Video Server", ACM Multimedia Systems, vol 4, pp. 112-121, 1996.
6. J. Flint, "Marketers Should Learn to Stop Worrying and Love the PVR", The Wall Street Journal, Oct 2005.
7. L. Gao, J. Kurose, D. Towsley, "Efficient Schemes for Broadcasting Popular Videos", NOSSDAV 98, July 1998.

8. L. Golubchik, C. Lui, R. Muntz, "Adaptive Piggybacking: A Novel Technique for Data Sharing in Video-on-Demand Storage Servers", ACM Multimedia Systems, vol. 4, no#0, pp 14-55, 1996.
9. J. Gong, Y. Syed, "Optimal QAM Assignment in the Presence of Mixed SD and HD Stream", NCTA NationalShow 2005.
10. G. Hardin, "Session Resource Management: How to Slice the Pie Allocating Bandwidth for Standard and High-Def VOD", Communications Technology Magazine, May 2005, available at : http://www.ct-magazine.com/archives/ct/0505/0505_sessionresource.htm
11. K. Hua, Y. Cai, S. Sheu, "Patching: A Multicast Technique for True Video-on-demand", IEEE Multimedia, vol. 4, pp. 51-62, 1997.
12. L. Juhn, L. Tseng, "Harmonic Broadcasting for Video-on-demand Service", IEEE Transactions on Broadcasting, vol 43 pp.268-271, Sept 1997.
13. W. Liao, V. Li, "The Split and Merge Protocol for Interactive Video-on-Demand", IEEE Multimedia, vol. 4, pp.51-62, 1997.
14. S. Lau, J. Lui, L. Golubchik "Merging Video Streams in a Multimedia Storage Server: Complexity and Heuristics", Multimedia Systems, vol. 6, no. 1, pp29-42, 1998.
15. S. Ross, "Introduction to Probability Models", Academic Press, 2003.
16. DOCSIS® Specifications, Cable Television Laboratories, Inc.(<http://www.cablemodem.com/primer/>)

Appendix A; Non-mixing Algorithm

In this appendix, we show the details of the non-mixing algorithm, taking a SD stream request as an example. The mathematical notations are defined in Section 4.

1. Identify a set of I , s.t. $Q - q_i \geq r_s$ for $\forall i, i \in I$
 - 1.1 If I is empty, reject the stream request;
2. Identify a subset of J , $J \subseteq I$, s.t. $S_j(q_j) = 3, j \in J$;
 - 2.1 If J is empty, go to the next step;
 - 2.2 If J has multiple elements, select $j^* = \arg \underset{j \in J}{\text{Min}} Q - q_j$;
 - 2.3 If there are multiple j^* , select randomly among j^* ;
3. Identify a subset of J , $J \subseteq I$, s.t. $S_j(q_j) = 1, j \in J$;
 - 3.1 If J is empty, go to the next step;
 - 3.2 If J has multiple elements, select j^* randomly;
4. Identify a subset of J , $J \subseteq I$, s.t. $S_j(q_j) = 2, j \in J$;
 - 4.1 If J is empty, go to the next step;
 - 4.2 If J has multiple elements, select $j^* = \arg \underset{j \in J}{\text{Min}} Q - q_j$;
 - 4.3 If there are multiple j^* , select randomly among j^* ;
5. Identify a subset of J , $J \subseteq I$, s.t. $S_j(q_j) = 4, j \in J$;
 - 5.1 If J has multiple elements, select $j^* = \arg \underset{j \in J}{\text{Max}} Q - q_j$;
 - 5.2 If there are multiple j^* , select randomly among j^* ;