

Exploring Wait Tolerance in Effective Batching for Video-on-Demand Scheduling

Hadas Shachnai*

Philip S. Yu†

Abstract

In a Video-on-demand (VOD) computer system, batching requests for the same video to share a common data stream can lead to significant improvement in throughput. Using the *wait tolerance* characteristic that is commonly observed in viewers behavior, we introduce a new paradigm for scheduling in VOD systems. We propose and analyze two classes of scheduling schemes: the Max_Batch and Min_Idle schemes, that provide two alternative ways for using a given stream capacity for effective batching. In making a video selection, the proposed schemes take into consideration the next stream completion time as well as the viewer wait tolerance. We compared the proposed schemes with the two previously studied schemes: (1) first-come-first-served (FCFS) that schedules the video with the longest waiting request and (2) the maximum queue length (MQL) scheme that selects the video with the maximum number of waiting requests. We show through simulations that the proposed schemes outperform substantially FCFS and MQL in reducing the viewer turn-away probability, while maintaining a small average response time. In terms of system resources, we show that by exploiting the viewers wait tolerance, the proposed schemes can reduce significantly the server capacity required for achieving a given level of throughput and turn-away probability as compared to the FCFS and MQL. Furthermore, our study shows that an aggressive use of the viewer wait tolerance for batching may not yield the best strategy, and that other factors, such as the resulting response time, fairness, and loss of viewers, should be taken into account.

*Contact author: The Department of Computer Science, Technion, Haifa 32000, Israel. FAX: 972-4-829-4353. e-mail:hadas@csa.cs.technion.ac.il. Part of this work was done while the author was at IBM T.J. Watson Research Center, Yorktown Heights, NY.

†IBM T.J. Watson, P.O. Box 704, Yorktown Heights, NY 10598. e-mail:psyu@watson.ibm.com

1 Introduction

In a Video-on-demand [6, 5, 8] (VOD) system, subscribers can choose both the video they wish to view and the time they wish to view it. This contrasts with services in which users can choose only from a small set of selections and watch them at pre-specified times. The server for a VOD service is essentially a homogeneous system, in which a large number of video streams of the same format, e.g., MPEG compressed, are stored and retrieved.

For a commercial environment, VOD servers may potentially need to handle thousands of requests simultaneously. There are often popular videos which are requested by many of the viewers. Batching multiple viewing requests to the same video stream can greatly increase the number of viewers supportable by a video server. In an ATM type network with multicast facility, the same video data can be sent to multiple viewers simultaneously without additional cost. Due to the isochronous requirement, each server is assumed to have a limit on the number of simultaneously supportable video streams. This is referred to as the stream capacity of the VOD server. Batching allows to support a number of viewers that is larger than the stream capacity.

In order to obtain batching we may need to delay deliberately viewer requests, so that later arrivals can share the same video stream with earlier requests. However, viewers may turn away when the *batching interval* becomes large. In general, viewers can tolerate and accept a small amount of delay (say, 5 minutes). This phenomenon is referred to as *wait tolerance*¹. As the delay grows larger, viewers are likely to leave the system without being served. On the other hand, if the batching interval is too short, more of the system stream capacity is consumed and viewers may have to wait longer in the queue for available streams. This can also result in loss of viewers.

The conventional batching approach in VOD systems uses only the states of the queues and viewers arrival times in selecting videos for showing. One scheme is the *first-come-first-served (FCFS)*, under which all video requests join a single request queue, and when a stream capacity becomes available, the request at the head of the queue is served. To support batching, all other requests queuing for the same video are also served by this stream. An alternative is to maintain a separate request queue for each video and select for the next showing the video with the longest queue. This is referred to as the *maximal queue length (MQL)* scheme.

As mentioned above, waiting requests may leave the queue when their waiting times exceed the tolerance of the requesting viewers. This loss of viewers is undesirable. Different scheduling schemes lead to various degrees of batching and losses of viewers. In selecting the next candidate for showing,

¹Conceivably, the distribution of the viewer wait tolerance can be obtained empirically by any service provider.

the FCFS scheme does not take into account the batching factor of the videos (i.e., the number of pending requests for each video) while the MQL ignores the wait time already incurred by the pending requests. Comparison study of the performance of these schemes can be found in [2], where it is shown that FCFS generally performs better than MQL.

In this paper we introduce a new scheduling paradigm that makes use of the viewer wait tolerance for effective batching. The central issues to be addressed here are the extent to which batching should be pursued, and the other factors that need to be considered. We propose and analyze two classes of scheduling schemes referred to as the Max_Batch and the Min_Idle schemes. Both classes of schemes produce delayed schedules that aim to increase the batching effect for popular videos. While the Max_Batch schemes maximize viewers batching and allow available streams to remain idle in the presence of pending requests, the Min_Idle schemes apply batching selectively only to more popular videos and avoid stream capacity from idling by serving unbatched requests of less popular videos (We elaborate on that in Section 3).

We focus on the main measure of system throughput: We show the computational complexity involved in finding the *optimal* strategy with respect to this measure, even for the case where we have a *complete* knowledge of system behavior². This motivated our search for efficient heuristic solutions in Section 3. A simulation program was developed to evaluate the performance of the proposed schemes. We show in Section 4 that the proposed schemes can outperform substantially FCFS and MQL in reducing the turn-away probability. A summary is given in Section 5.

We briefly comment on some related works that study alternative means to support stream sharing. In [4, 7], a bridging approach that uses buffering to support stream sharing is considered. By keeping the last few minutes of a leading video request in buffer, trailing requests can be served out from the buffer. In [9], buffers are used to support pause-resume under batching in a look-ahead scheduling scheme. In [3], stream sharing is achieved through a piggybacking approach of adjusting the playout rate of the videos in progress. By slowing down the leading video request and/or speeding up the trailing video request, the two requests can be merged into one. While batching avoids the additional resource requirement as in the bridging approach or the complexity in the piggyback approach, it does cause additional waiting of the viewers. Nevertheless, batching can be used in conjunction with either the bridging approach or the piggybacking approach to strike a balance among wait time, resource requirement and implementation complexity.

²Due to space limitations we relegate the discussion to the Appendix.

2 Preliminaries

Assume a VOD system where the initiation of video requests corresponds to arrival process \mathcal{A} with mean time $1/\lambda$ between successive requests. Upon arrival a viewer chooses with probability p_k to join the queue of video k , $1 \leq k \leq M$. The viewer waits in the queue for a certain interval of time after which he may decide to leave without being served. We define the viewer's *wait threshold* as the scheduler estimate of the amount of wait time a viewer is willing to tolerate before he considers leaving the system. In terms of queuing theory the viewers are classified as customers of the *reneging* type, with a wait time given by a random variable (*rv*) R .

The system can support N_{max} simultaneously playing streams, where each stream shows one of the M videos. Whenever a stream becomes available (upon viewing completion), the scheduler inspects the queues and chooses the next video to be shown using that stream. The average length of a video is given by L .

In our analysis of scheduling schemes for the above model, we use the following performance measures:

- *Turn-away probability*: This is the probability that an arriving viewer leaves the system without being served (due to a wait time which exceeded the viewer's wait tolerance). Minimizing the turn-away probability maximizes the system throughput, therefore in the present context we often use the term *throughput measure* when referring to minimization of the turn-away probability.

measures:

- *Average response time*: For any viewer that is served by the system (i.e. a non-reneging viewer), the time that elapses from the viewer's arrival until the start of his show is the viewer's response time. The reneged viewers are ignored, since the response time is not defined for a reneged viewers.
- *Fairness*: Let $q_t(i)$ denote the turn-away probability associated with the i -th video $1 \leq i \leq M$, and let q_t be the average turn-away probability of the system. We say that a system is *totally fair* if $q_t(i) = q_t$, $\forall 1 \leq i \leq M$. As in [2], we measure the degree of unfairness in the system as the average deviation, σ_f , from q_t , using a partition of the series of requests, sorted by video ids, to $|S|$ subsets of roughly equal number of requests.

Thus,

$$\sigma_f = \frac{\sum_{i=1}^{|S|} (q_t(i) - q_t)^2}{|S|} \quad (1)$$

For a given scheme, the system's *degree of fairness* increases as σ_f gets closer to 0.

While all of the above measures are important for achieving good quality of service, we refer in our analysis to the system throughput as the *main* measure, with the average response time and the fairness ratio the *secondary* measures. This is due to the fact, that the throughput measure is essential in determining the efficiency of a given scheme. Indeed, some of the schemes may provide small response time with high degree of fairness and yet perform very poorly: A typical example is the *Last Come First Serve* strategy, which minimizes the average response time, while causing a large fraction of the arriving viewers to leave the system without being served (i.e. the turn-away probability is close to 1). Thus, the last two measures provide only secondary information in the overall performance evaluation of a given scheme.

3 The Wait Tolerance Batching Approach

In this section we introduce heuristic batching schemes that effectively use the viewers wait tolerance. Intuitively, our schemes use some limited knowledge of the reneging pattern of the viewers for delaying the start of a new showing based on the future completion schedule, and thus increase the batching effectiveness in the system. We keep this delay short enough, so as to minimize the expected loss of viewers. Therefore, the overall throughput of the system increases.

We classify the M videos as *hot* or *cold* videos, based on their request frequencies. Assuming the videos are numbered in decreasing request frequencies, a classification threshold $\tau \in \{1, 2, \dots, M\}$ determines the boundary between these two categories. The selection of τ and the sensitivity of the proposed schemes to its value will be discussed later. We considered the following two classes of batching schemes:

1. **Max_Batch Schemes** – A video is eligible for scheduling only if some of its requests have wait times reaching or exceeding a pre-specified *batching threshold*. This in fact imposes a *minimum wait time requirement*. If no video satisfies this eligibility criterion, an available stream remains unused and the scheduling point is delayed to a later time, in which a request

<p>$\mathcal{H} = \{1, \dots, M\}$.</p> <p>At each scheduling point t_0 (that can either be a stream completion, a new request arrival or some delayed scheduling instant)</p> <p>Let $L_{\mathcal{H}}$ denote the subset of videos in \mathcal{H} with waiting requests that reached or exceeded the batch threshold by t_0.</p> <p>if $L_{\mathcal{H}} \neq \emptyset$ a video is scheduled using one of the criteria:</p> <p>(a) Maximal queue length (the <i>BMQ scheme</i>).</p> <p>(b) Maximal expected loss till the next completion at t_1, i.e. $\max_{k \in L_{\mathcal{H}}} l_k(t_1 - t_0)$ (the <i>BML scheme</i>).</p> <p>else</p> <p>delay the schedule until a request reaches its batch threshold</p>
--

Figure 1: The BMQ and BML schemes

reaches the batching threshold. The objective is clearly to maximize batching, which may result in stream idling. The batching threshold is chosen based on the wait tolerance.

2. **Min_Idle Schemes** – These schemes try to pursue batching on hot videos without the minimum wait time requirement (used in the Max_Batch approach). Cold videos are always eligible for scheduling. The eligibility criterion is thus substantially relaxed. The objective is to provide small response time and to decrease loss of viewers among those requesting cold videos. This is done by minimizing stream idling.

The Min_Idle schemes maintain two sets of video queues, denoted by \mathcal{H} and \mathcal{C} , while the Max_Batch schemes only maintain one set (\mathcal{H}) of queues. Under the Min_Idle schemes, the partition of videos to the sets \mathcal{H} and \mathcal{C} is defined dynamically based on the current state of the system and the classification threshold τ .

3.1 The Max_Batch schemes

The Max_Batch schemes are presented in Figure 1. When a stream becomes available at t_0 , we pick the next video among the videos with requests exceeding a pre-specified *batch threshold*, denoted by ν . (This eligible subset of videos is denoted by $L_{\mathcal{H}}$.) We consider two stream selection strategies: The first strategy chooses the video with the maximum queue length to increase the batching effect and is thus called the *Max_Batch MQL* (or *BMQ*) scheme. Since the video with the maximal queue may not have the largest number of long waiting requests which are most likely to leave, the other

At each scheduling point t_0 (that can be either a stream completion or an arrival of a request)

Define \mathcal{H} and \mathcal{C} :

A video k is in \mathcal{H} , if

(a) $k \leq \tau$.

(b) $n_k > 1$.

(c) $n_k = 1$ and $t_0 - a_1 > \nu$;

/* A single viewer that exceeded batch threshold */

otherwise the video is in \mathcal{C} .

Let $L_{\mathcal{H}}$ denote the subset of videos in \mathcal{H} with a positive

expected loss in the interval $[t_0, t_1)$, i.e. $L_{\mathcal{H}} = \{k \in \mathcal{H} \mid l_k(t_1 - t_0) > 0\}$.

if $L_{\mathcal{H}} \neq \emptyset$ a video is scheduled using one of the

criteria as defined for the Max_Batch case

(*IMQ* and *IML* respectively)

else

schedule a video in \mathcal{C} (if any) using FCFS.

Figure 2: The IMQ and IML schemes

strategy chooses the video that will incur the maximal expected loss if delayed until the next stream completion instant t_1 (i.e., the next time a stream becomes available). This is referred to as the *Max_Batch with minimal loss* (or *BML*) scheme. With no explicit knowledge of the viewer behavior in calculating the expected loss, the scheduler simply assumes that any viewer reaching the batch threshold will depart. (Therefore a video will become eligible for scheduling as soon as some viewer reaches the batch threshold.) Under this assumption, the expected number of departures during the interval $[t_0, t_1)$ (denoted by $l_k(t_1 - t_0)$ in step (b) of Figure 1) in the k -th video queue can be easily calculated based on the arrival time of each request in the queue. In our experimental studies, different viewer behaviors were considered, but the scheduler always used this simplified assumption. If at the stream completion time there is no viewer exceeding the batch threshold, the stream remains idle. The schedule is delayed until one of the viewers reaches the batch threshold.

The batch threshold parameter needs to be determined. While a natural choice is to assign to it the size of the wait threshold (i.e. $\nu = \omega$), our experimental studies (as described in the next section) showed that as the viewer wait threshold increases, choosing the maximal possible delay for batching (that does not incur loss of viewers) provides only a small increase in system throughput, and results in significant increase in the average response time. Thus, the Max_Batch schemes perform better by

choosing $0 < \nu(\psi) \leq \omega$ for a desired *upper bound* ψ on the average response time. By optimizing on the value of ν we can optimize the throughput of the Max_Batch schemes while keeping the average response time small.

3.2 The Min_Idle schemes

The Min_Idle schemes are presented in Figure 2. Here the videos are partitioned to two sets. A video is considered to be in \mathcal{H} , if it satisfies one of the following criteria: (1) it is a hot video, (2) its queue has more than one request, and (3) the only request in its queue has wait time that exceeded the batch threshold (i.e. the difference between arrival time a_1 and scheduling time t_0 exceeds the batch threshold). Otherwise, a video belongs to the set \mathcal{C} . Videos in \mathcal{H} will be delayed to achieve batching, but will be given also preferential treatment in scheduling when the batch threshold is approaching. Criterion (2) above provides means to recognize hot videos dynamically, and reduces the sensitivity to the classification threshold τ as will be shown in the next section. Criterion (3) brings the long waiting cold video requests into the \mathcal{H} set to avoid cold video requests being completely blocked out by the hot video showings.

When a video stream becomes available we determine the eligible subset ($L_{\mathcal{H}}$) of videos which have requests with wait times that will exceed the batch threshold upon the next stream completion. If $L_{\mathcal{H}}$ is not empty, one of the videos in $L_{\mathcal{H}}$ will be chosen. The selection criterion can again be either the maximum queue length (referred to as the *IMQ* scheme) or the maximum expected loss (referred to as the *IML* scheme) as in the Max_Batch schemes. If $L_{\mathcal{H}}$ is empty, a cold video in the set \mathcal{C} will be chosen. Since videos with more than one outstanding requests are in \mathcal{H} , each video in \mathcal{C} can have a queue of length at most one. We schedule these queues by FCFS.

The rationale for maintaining the sets \mathcal{H} and \mathcal{C} in the Min_Idle schemes is that different scheduling criteria can be applied to hot and cold videos. Hot videos (or any video with a queue length of two or more) will be scheduled only if further delay will result in viewers exceeding the batch threshold. This will improve batching. If no hot video viewer is going to exceed the batch threshold before the next stream completion, a video in \mathcal{C} is scheduled. Since it is hard to achieve batching effect for cold videos, these video requests will be scheduled as soon as there is additional stream bandwidth available (when no hot video needs to be scheduled immediately). This is in contrast to the Max_Batch schemes, where scheduling of every video is delayed as much as possible. To avoid cold video requests being completely blocked out by the hot video showings, Min_Idle schemes also transfer the long waiting cold video requests into set \mathcal{H} .

Except for the BMQ scheme, all the other schemes (BML, IML and IMQ) take into consideration the next stream completion time in addition to the viewer wait tolerance in making the scheduling decision. Information on the next stream completion time can be used to help decide whether a viewer can be further postponed. We note that batching or any other stream sharing schemes generally work most effectively when VCR-like functions are not supported. In this case, the completion schedule can be easily calculated since the playout time of each video is known a priori. Indeed, a main reason for the complexity of using batching when VCR-operations are allowed, is the fact that once a viewer ‘pauses’, he is disconnected from his batching group. This viewer may not find an idle stream when he presses ‘resume’. This may result in a significant degradation in the system’s quality of service. When VCR-like functions need to be supported, the completion schedule has to be dynamically maintained. Since we are only looking for the next stream completion time (within the range of the batch threshold), the prediction of completion time is made nearly at the end of a video showing, not at the beginning. Thus the probability of error and the effect of uncertainty are greatly reduced.

4 Simulation Study

4.1 Model description

We assume that video requests are generated by a Poisson process and the interarrival time is exponentially distributed with mean $1/\lambda$. Upon arrival to the system, a viewer chooses to wait for a showing of video k with probability p_k , $1 \leq k \leq M$, with the p_k ’s given by the *Zipf* distribution:

$p_k = \frac{c}{k}$, where $c = 1/(\sum_{k=1}^M \frac{1}{k})$ is a normalization constant. A viewer may leave the system without

being served. Two scenarios of viewer renegeing behavior are considered. In the first scenario, there is a minimum amount of time, ω , that a viewer is willing to wait before renegeing, while in the second one a viewer may leave at *any* time after his arrival to the system, in other words, while in the first scenario $\omega > 0$, in the second scenario $\omega \equiv 0$.

- **Scenario 1 (Deterministic case):** After waiting an amount of time given by the wait threshold ω , the viewer’s remaining time until departure is given by the rv R , which is determined by an Exponential distribution. That is, $R \sim \exp(\mu)$, where $1/\mu$ is the mean time that elapses after ω until the viewer leaves the system. Thus the viewer’s overall wait time is given by $W = \omega + R$.

- **Scenario 2 (Probabilistic case):** The viewers departure time W follows a Normal distribution with mean μ and variance σ^2 , i.e. $W \sim N(\mu, \sigma^2)$. A *confidence parameter* ζ is used by the scheduler for determining a wait threshold ω for the viewers. The value of ω is defined as

$$\omega = \min \{r : \text{Prob}(W > r) \leq 1.0 - \zeta\} , \quad (2)$$

i.e., the value of ω is chosen as the minimal amount of time after which a viewer is likely to leave the system immediately with probability higher than ζ .

Obviously, taking ζ close to 1 causes larger loss of viewers, while a small value for ζ decreases the batching effect in the system, since the viewers' wait tolerance allows further delay in their schedule.

4.2 Performance Results

We first describe the set of base parameter values used in the simulations: Assuming a system of $M = 100$ videos with the average length of a showing $L = 120$ minutes³, we used a server capacity of 800 streams and a video request rate of 50 arrivals per minute.

The classification threshold for distinguishing between videos in \mathcal{C} and in \mathcal{H} was chosen to be $\tau = 20$. For Scenario 1, the wait threshold is assumed to be $\omega = 5.0$ minutes and the average remaining time till departure was 3 minutes (that is, $R \sim \text{exp}(1/3)$). We set the batch threshold equal to the wait threshold, i.e. $\nu = \omega$ (unless otherwise specified). For Scenario 2, with the renegeing pattern that corresponds to the Normal distribution, we use the value $\mu = 5.0$ minutes as mean for W with the variance $\sigma^2 = 2.5$ minutes.

In Figures 3–6 we compare the performance of the proposed schemes to the MQL and FCFS schemes. In these experiments the differences in performance among the four proposed schemes (IML, IMQ, BML and BMQ) were subtle compared to the difference between any of these schemes and the MQL or FCFS, thus we use the numerical results obtained for the IML to represent the performance of all four schemes. In the sequel we focus on the four proposed schemes and demonstrate some of the differences between our four schemes (see Figures 7–11).

In Figure 3 we show the turn-away probability vs. arrival rate to the VOD system. We note, that IML leads to lower turn-away probability, e.g., for the rate of 50 arrivals per minute the IML decreases by 7% the turn-away probability obtained by the FCFS and by 16% the turn-away probability measured under MQL. As the load increases, the turn-away probability under FCFS and MQL

³To simplify the performance study, we do not assume any VCR-like functions being supported.

approaches 29% and 34% respectively, while under IML it reduces to at most 17%. Figure 4 shows that IML also provides smaller average response time compared to FCFS, except for low arrival rate; The MQL yields a typical response time of at most 1 minute at the cost of causing more than 25% of the viewers to renege before being served (see Figure 3). In Figure 5 we examine the fairness issue. The unfairness of IML is close to that of FCFS and is much better than the unfairness measured under the MQL.

Next we use Scenario 2 to study the effect of uncertainty on the performance of our schemes. We consider the case where $W \sim N(5.0, 2.5)$. The value of the confidence parameter (ζ) used by the scheduler is 0.75. Figure 6 shows the turn-away probabilities of IML, MQL and FCFS under different server capacities. We conclude that even with uncertainty on the wait threshold, IML outperforms MQL and FCFS. In fact, changing the confidence parameter to any other value between 0.35 and 1 would increase the turn-away probability by at most 5%, as shown in Figure 7.

In Figures 8–11 we present some distinctions between the proposed schemes. Figures 8 and 9 compare between the MQL-Based criterion (used by the BMQ and IMQ schemes) and the Minimal-Loss criterion (used by the BML and IML schemes):

Figure 8 shows the advantage of the IMQ in providing inherently smaller average response time. The difference is most significant for low arrival rates, in which case, under the Max_Batch schemes every request needs to wait out the batch threshold. When the load increases, the response time of the Max_Batch schemes somewhat improves, as an arriving viewer is likely to find his queue non-empty, and therefore larger amount of viewers will be scheduled before exhausting their batch threshold. Under the Min_Idle schemes, as the load increases the response time first goes up, due to stream contention, and then reduces with the increase in the batching effect. This reduces the average wait time in the system for high arrival rates. Note, that the preference given by the MQL criterion to *large* batches of viewers increases significantly the wait times of viewers in the cold video queues. However, since these viewers often leave the system without being served their wait times have no effect on the system response time. On the other hand, the Minimal-Loss-based criterion tends to prefer *long waiting* viewers. That implies scheduling more frequently small batches of viewers waiting in the cold video queues. Choosing to schedule at any time a small batch of viewers results in further delay in the schedules of larger batches, and a corresponding increase in the overall system response time. This explains the fact, that the IMQ leads to smaller response time than IML and the other schemes under any load.

Figure 9 provides a distinction between the MQL-based criterion and the Minimal-Loss-based

criterion with respect to fairness. We note that as the load increases, the MQL-based schemes yield higher unfairness. Again, this shows a side effect of over-emphasis on batching. Furthermore, BML (respectively, BMQ) shows slightly smaller unfairness than IML (respectively, IMQ) at the expense of higher requirement of minimal server capacity, as shown in Figure 11. This is due to the fact that by imposing a minimum wait time requirement, the Max_Batch schemes in fact limit the number of streams that can be simultaneously used by each video. (The upper limit is set by the length of a video divided by the batch threshold.) This prevents the hot videos from using up all the streams and improves fairness. The Min_Idle schemes do not impose such a wait time restriction but allow for dynamic adjustment of the cold videos into the set \mathcal{H} . The turn-away probabilities are shown in Figure 10. At high arrival rate, Minimal-Loss-based criterion leads to slightly smaller turn away probability, especially for the Min_Idle schemes. Figure 11 indicates that for a wait threshold of 5.0 minutes the Min_Idle schemes require less server capacity for guaranteeing a turn-away probability of at most 5%. However, as noted above, the difference between the schemes is small.

5 Conclusion

We have presented a paradigm for multimedia scheduling schemes that explores viewers wait tolerance for increasing the system throughput. Based on this approach we have introduced four schemes, which provide substantial improvement compared to previously studied FCFS and MQL schemes, especially for scenarios of heavy loads, for which the efficient use of server capacity is crucial in providing good quality of service.

The four alternative schemes (IML, IMQ, BML and BMQ) proposed differ not only in their eligibility criteria for scheduling but also in stream selection rules. The eligibility criterion is introduced to facilitate batching. We found that an aggressive pursue of batching may not lead to the best performance. Other factors like response time, expected viewer losses and stream idling also need to be taken into consideration.

Under the Max_Batch (BML and BMQ) approach, videos are not eligible for scheduling until they have viewers with wait time exceeding the batch threshold. The Min_Idle (IML and IMQ) approach relaxes the eligibility criterion: Batching is only pursued for the hot videos without imposing the requirement to wait out the whole batch threshold. Cold videos are always eligible for scheduling under the Min_Idle schemes. These schemes require less server capacity for guaranteeing a given turn-away probability.

Two streams selection criteria were considered. We studied the MQL criterion (BMQ and IMQ)

versus the Minimal-Loss criterion (BML and IML) for stream selection. It is found that for any load, the MQL based schemes, that try to optimize batching, yield smaller response time (at the expense of fairness) while keeping the throughput close to the throughput provided by the Minimal-Loss based schemes.

References

- [1] E. F. Beckenbach, Editor, Applied Combinatorial Mathematics, J. Wiley and sons, 1964.
- [2] A. Dan, D. Sitaram and P. Shahabuddin, "Scheduling Policies for an On-Demand Video Server with Batching", *Proc. ACM Multimedia '94*, SF, CA, Oct. 1994, pp. 391-398.
- [3] L. Golubchik, J.C.S. Lui and R. Muntz, "Reducing I/O Demand in Video-on-Demand Storage Servers", *Proc. ACM SIGMETRICS '95*, Ottawa, Ontario, Canada, May 1995.
- [4] M. Kamath, D. Towsley and K. Ramamritham, "Buffer Management for Continuous Media Sharing in Multimedia Database Systems", *Technical Report 94-11*, Univ. of Massachusetts, Feb. 1994.
- [5] P. V. Rangan and H. M. Vin, "Designing File Systems for Digital Video and Audio", *Proc. of 12th ACM Symposium on Operating Systems*, 1991.
- [6] W. Sincoskie, "System Architecture for Large Scale Video on Demand", *Computer Networks ISDN System*, Vol. 22, 1991, pp. 155-162.
- [7] D. Rotem, and J.L. Zhao, "Buffer Management for Video Database Systems", *Proc. Intl. Conf. on Data Eng.*, Taipei, Taiwan, March 1995, pp. 439-448.
- [8] H. M. Vin and P. V. Rangan, "Designing a Multi-User HDTV Storage Server", *UCSD Technical Report CS92-225*, Jan. 1992.
- [9] P. S. Yu and J.L. Wolf and H. Shachnai, "Design and Analysis of A Look-ahead Scheduling Scheme to Support Pause-Resume for Video-on-Demand Applications", *Multimedia Systems*, Vol. 3, No. 4, Sept. 1995, pp. 137-149.

Appendix: The Optimal Strategy

In this section we formulate the problem of maximizing the throughput in a VOD system as a stochastic optimization problem, and derive the functional equation defining the optimal scheme. Assuming a complete knowledge of the arrival process \mathcal{A} , the distribution of the wait tolerance rv R and the request probabilities for the videos p_1, \dots, p_M , our objective is to maximize the throughput in a given time interval $[0, T)$.

Typically, the scheduler is activated when a stream becomes available, i.e. *scheduling points* are identified with video stream completions. Given a scheduling point t_0 and the state of the previously scheduled videos, we denote by t_l , $l \geq 1$ the time of the l -th completion after t_0 . Let $l_k(t_l - t_0)$ be the expected loss of viewers for video k in the interval $[t_0, t_l)$, i.e., given the set of viewers waiting in queue k at t_0 :

$$l_k(t_l - t_0) = E(\text{Departures in } (t_l - t_0)) .$$

Using the rv R , with $a_r < t_0$ the arrival time of the r -th viewer in the queue of video k , we have

$$l_k(t_l - t_0) = \sum_{r=1}^{n_k(t_0)} \text{Prob}(R < (t_l - a_r)) . \quad (3)$$

where $n_k(t_0) \equiv n_k$ is the number of viewers in the queue of video k at t_0 .

A straightforward approach is the *Greedy* scheme defined as follows:

- At any scheduling point t_0 , $\forall 1 \leq k \leq M$ compute the expected number of viewers reneging from queue k before the next completion, given by $l_k(t_1 - t_0)$, assuming k is not scheduled for showing at t_0 .
- Schedule the video with maximal expected loss in the interval $[t_0, t_1)$.

It is easy to verify, that Greedy is not optimal: This is due to the fact that Greedy's selections of videos do not use information on

- (i) The expected number of *arrivals* to each of the queues before the next scheduling point, and
- (ii) The expected loss of viewers, in the long run, i.e., after the *next* scheduling point, following a selection of a video at some time instance t .

Consider, e.g., the case where $M = 10$, and there are 4 videos with request probabilities $p_i = 0.24$ each, and the request probabilities for the other (cold) 6 videos are negligibly small (and sum to 0.04). Suppose, that in some scheduling point at time t , there is a long waiting request to one of the cold movies j , and there is a large number of requests in each of the hot movies queues. Since the wait times of these requests are small, the Greedy chooses to schedule j . Assume that the next scheduling point is at $t + \delta$, and after that at $t + 2\delta$. Then, if wait times of the requests in the hot movies queues are approaching at time $t + \delta$ their wait tolerance threshold, obviously, selecting a hot movie at time t would result in a smaller overall expected loss of viewers to the system.

Hence, computing the expected loss of viewers only in the interval $[t, t + \delta)$ is insufficient for the minimization of the expected loss in the interval $[0, T)$.

Yet, as shown in our simulation study, the Greedy scheme provides a significant increase in the

system throughput, compared to the basic MQL and FCFS schemes (We call this scheme below the Max_batch scheme).

We proceed with the computation of the optimal scheme. Let i be the state of the system given by the waiting times of the viewer in queues $1, \dots, M$:

$$i = (w_{1,1}, \dots, w_{1,n_1}, w_{2,1}, \dots, w_{2,n_2}, \dots, w_{M,1}, \dots, w_{M,n_M}) , \quad (4)$$

where $w_{j,k}$ is the wait time of the j -th viewer in the queue of video k . We denote by $f(i)$ the expected loss of viewers when the system starts from state i , and the optimal strategy is used.

A decision D in state i to schedule all viewers in queue k , for some $1 \leq k \leq M$ yields the expected number of renegeing viewers in the interval $[t_0, t_1]$ given by

$$r_i(D) = \sum_{m \neq k} l_m(t_1 - t_0) . \quad (5)$$

Let $p_{ij}(D)$ be the probability of transition to state j from state i given the decision D . Bellman's principle of optimality [1] asserts that if decision D is taken in state i , and at the next scheduling point the optimal strategy f is used, the total expected loss of viewers is

$$r_i(D) + \sum_j p_{ij}(D)f(j) . \quad (6)$$

Since this holds for any choice of D , clearly an optimal policy D is chosen so as to minimize (6). Hence, the optimal strategy satisfies

$$f(i) = \min_D [r_i(D) + \sum_j p_{ij}(D)f(j)] . \quad (7)$$

Using successive approximations, f can be computed directly whenever the state space is finite. For the above problem, the renegeing property can be used to truncate the infinite state space to a finite space for which the optimal strategy may be computed.

We now discuss the time and space complexity for computing the optimal strategy. Using the classical method of successive approximations, we define

$$f_0(i) = \min_D r_i(D)$$

$$f_{n+1}(i) = \min_D r_i(D) + \sum_j p_{ij}(D)f_n(j) .$$

The first function, $f_0(i)$, represents the minimum expected loss of viewers after any particular scheduling point, whereas $f_n(i)$ gives the minimum expected loss of viewers under the constraint, that there are at most $(n + 1)$ additional scheduling points till the end of the scheduling interval. Typically, the difference $f(i) - f_n(i)$ is geometrically decreasing, i.e., for some $0 < x < 1, c > 0$,

$$f(i) - f_n(i) \leq cx^n ,$$

therefore it is sufficient to successively approximate $f_n(i)$ for relatively small values of n , which may be no larger than 10. However, we note that for any state i , as defined in (4), computing $f_{n+1}(i)$

involves a multiplication of the matrix $\{p_{ij}(D)\}$ (of the size $N \times N$, where N is the number of states), and the vector $(f_n(1), \dots, f_n(N))$. Then we minimize over M possible decisions. Suppose that the length of each queue never exceeds some $0 < l < 50$, and wait times are from the set $\{1, \dots, k\}$ in minutes, then $N = O(k^{l \cdot M})$. This implies that, e.g., for $M = 50$, and $k < 10$, $N > 2^{500}$. The space complexity is similar, since we need to store the transition matrix $\{p_{ij}(D)\}$. Thus, even for moderate values of M and T , the computational complexity as well as the space required for solving the functional equation for f can be enormous, thus we proceed with a study of approximation schemes. We note, that this computational complexity does not depend on the length of the scheduling interval T , and is therefore incurred also for the case where T is very small.

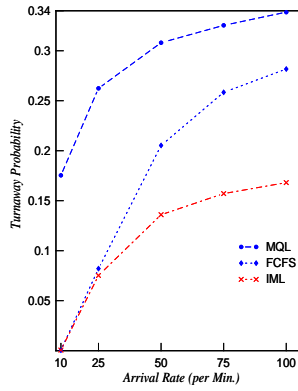


Figure 3: Turn-away Probability vs. Arrival Rate (Comparison to FCFS and MQL)

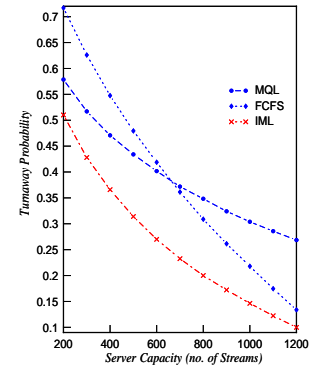


Figure 6: Turn-away Probability vs. Server Capacity (Reneging by the Normal Distribution)

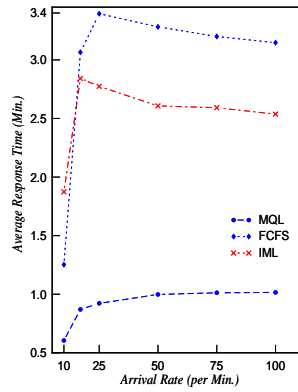


Figure 4: Average Response Time vs. Arrival Rate

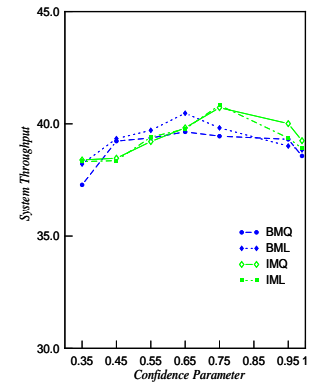


Figure 7: System Throughput for Reneging by the Normal Distribution with Varying Confidence Parameter

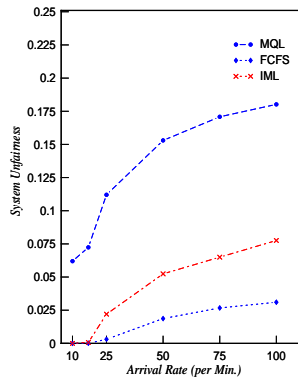


Figure 5: Unfairness vs. Arrival Rate

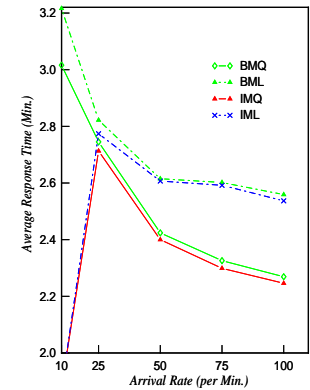


Figure 8: Average Response Time for Varying Arrival Rate

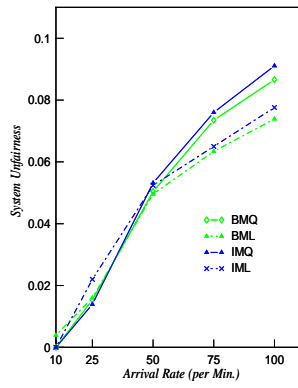


Figure 9: System Unfairness for Varying Arrival Rate

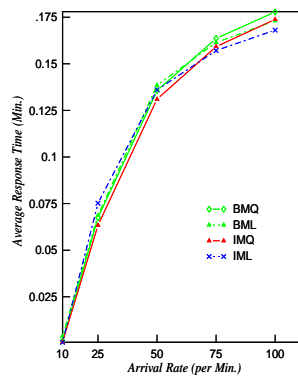


Figure 10: Turnaway Probability for Varying Arrival Rate (the Max_Batch Schemes vs. the Min_Idle Schemes)

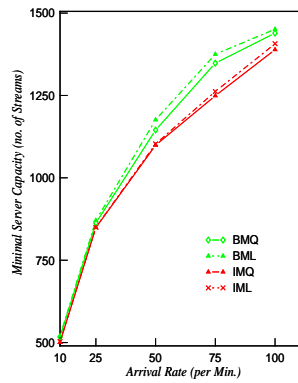


Figure 11: Minimal Server Capacity For a Turn-away of less than 5%