

Improving Perceived Fairness and QoE for Adaptive Video Streams

Bjørn J. Villa

Department of Telematics
Norwegian Institute of Science and Technology
Trondheim, Norway
bjorn.villa@item.ntnu.no

Poul E. Heegaard

Department of Telematics
Norwegian Institute of Science and Technology
Trondheim, Norway
poul.heegaard@item.ntnu.no

Abstract - This paper presents an enhancement to a category of Adaptive Video Streaming solutions aimed at improving both Quality of Service (QoS) and Quality of Experience (QoE). The specific solution used as baseline for the work is the Smooth Streaming framework from Microsoft. The presented enhancement relates to the rate adaption scheme used, and suggests applying a stochastic variable for the rate adjustment intervals rather than the fixed approach. The main novelty of the paper is the simultaneous study of both network oriented fairness in the QoS domain and perception based fairness from the QoE domain, when introducing the suggested mechanism. The method used for this study is by means of simulations and numerical optimization. Perception based fairness is suggested as an objective QoE metric which, requires no reference to original content. The results show that the suggested enhancement has great potential in improving QoE, while maintaining QoS.

Keywords - Adaptive Video Streaming; Fairness; QoE.

I. INTRODUCTION

Solutions for Adaptive Video Streaming are part of the more general concept of ABR (Adaptive Bit Rate) streaming which, covers any content type. The implementation of ABR streaming for video varies between different vendors, and among the more successful one today is the Microsoft Smooth Streaming (SilverLight) framework [1]. In general, the different implementations use many undisclosed and proprietary functions, awaiting results from ongoing standardization.

The basic behavior of adaptive video streaming solutions is that the client continuously performs a measurement and estimation of available resources in order to decide which, quality level to request. The relevant resource from the network side is the available capacity along the path between the server and client. Based on this, at certain intervals the client decides to either go up or down in quality level or remain at the current level. The levels are predefined and communicated to the client by the server at session startup. The changes in quality levels are normally done in an incremental approach, rather than by larger jumps in rate level. The rationale behind this is the objective to provide a smooth watching experience for the user. However, it may also be related to the CPU monitoring done by the client, as this is a key resource required. It may be the case that even if the network can provide you with a much higher rate level, the CPU on the device being used would not be able to

process it. During the initial phase of an adaptive streaming session the potential requests of change in rate level are more frequent than later on when operating in a more steady-state phase. To some extent this is a rather aggressive behavior from a single client which, may have undesirable inter-stream impacts. At the same time, in order to give the user a good first impression and make him want to continue using the service it is desirable to reach a high quality as soon as possible.

Among the strongest drivers for commercial use of ABR based services on the Internet are Over-The-Top content providers. These are providers which, rely on the best effort Internet service as transport towards their customers. Therefore, technologies aiming at making services survive almost any network state are of great interest. In addition to focus on the network based QoS dimensions of services and involved networks, there is also a growing interest in the QoE dimension [2]. The latter should be considered as not only a richer definition of quality, but also more focused towards who decides whether something is good or bad, i.e., the end user. The evolution of successful services on Internet indicates that the focus on QoE for Over-The-Top providers is a good strategy.

A. Problem Statement

The concept of Adaptive Video Streaming is without a doubt very promising. However, as more and more services are adopting this concept the success brings new challenges. The first challenge with effects visible to the end users is how well these services behave when they compete for a shared resource, such as the broadband access to a household. With a strong dominance of video based service on the Internet this issue is important to address. As each client operates independently of each other, it has no understanding of the traffic it competes with. Different clients consider each other as just background traffic. This leads to unpredictable and potentially oscillating behavior of each session, especially in a home environment this type of interference is likely to have a very negative impact on each user QoE.

B. Research Approach

The method investigated in this paper to address the problem at hand is to apply specific changes in the algorithm used by each ABR client controlling the adaptive behavior. The specific change suggested is related to the

rate adjustment interval used [1]. The effect of changing the duration of the rate adjustment interval from a fixed value T to some stochastic variable is presented and analyzed.

The ABR solution used as reference point for the work is the one from Microsoft (Smooth Streaming). However, the key principles would still apply to other solutions based on similar principles.

C. Paper Outline

The structure of this paper is as follows. Section II provides an overview of methodology and metrics; Section III describes the simulation model; Section IV presents simulation results; Section V gives an analysis of the results; Section VI provides the conclusions and an outline of future work.

II. RELATED WORK

It has been shown in [3] that competing adaptive streams can cause unpredictable performance for each stream, both in terms of oscillations and ability to achieve fairness in terms of bandwidth sharing. The experimental results presented give clear indication on that competing ABR clients cause degraded and unpredictable performance. Apart from this paper, the topic at hand does not seem to have been addressed by the academic research community to the extent it deserves.

In another paper [4], the authors have investigated how well adaptive streaming performs when being subject to variable available bandwidth in general. Their findings were that the adaptive streams are performing quite well in this type of scenario except for some transient behavior. These findings do not contradict the findings in [3] as the type of background traffic used do not have the adaptive behavior itself, but is rather controlled by the basic TCP mechanisms.

Rate-control algorithms for TCP streaming in general and selected bandwidth estimation algorithms are described in [5]. This work is relevant to any TCP based application delivering a video stream.

In some of our own previous work we have described and analyzed how competing adaptive streams can be controlled using a knowledge based bandwidth broker in the home gateway [6] [7].

III. METHODOLOGY AND METRICS

In this section, we introduce the relevant performance metrics and together with motivation for the chosen focus. Thereafter, some candidate methods on how to improve the performance metrics are given, and finally, the specific method subject for study is presented.

A. Flow Based Performance Metrics

For transport flows it is common [8] to focus on the following metrics in order to assess their performance: inter-flow fairness, stability and convergence time. This in addition to the general QoS metrics: bandwidth, packet loss, delay and jitter. The same metrics can be applied to adaptive

video streams as they by definition also are flows with similar concerns. The analysis of these metrics can be done from a strict network oriented perspective (QoS), but to some extent also bridged over to a user perception domain (QoE). When focusing on the inter-flow fairness metric this is traditionally analyzed [9] using, e.g., the Jain's fairness index [10], the product measure [11] or Epsilon-fairness [12] for flows with equal resource requirements. For flows with different resource requirement, the Max-Min fairness [13], proportional fairness [14] or minimum potential delay fairness [15] approaches are commonly seen. Real life adaptive video streams would typically belong to the last category.

Max-Min fairness: The objective of max-min fairness is to maximize the smallest throughput rate among the flows. When this is met, the next-smallest throughput rate must be as large as possible, and so on. Max-min fairness can also be explained by considering it as a progressive filling algorithm, where all flows start at zero and grow at the same pace until the link is full. With this approach the max-min fairness gives priority to the smallest flows. The least demanding flows always have the best chance of getting access to all the resources it needs.

Proportional fairness: The original definition of proportional fairness comes from economic disciplines [14] for the purpose of charging. The original definition is used in the relevant RFC [9] but it does not come across as very constructive for the purpose of analyzing fairness in single resource (e.g., bandwidth) sharing among flows. In this context more recent definitions and interpretations are more suitable [16]. The principle of this would be that a resource allocation is considered proportional fair if it is made to the flow which, has the highest ratio between potential maximum resource consumption and its average resource consumption so far. A further simplification would be to use the current resource usage (if greater than 0) instead of the average in the ratio calculation. The same ratio numbers for each flow could then be used to give a view on the current system fairness by comparing them. If they are all equal the system could be stated as proportionally fair.

Minimum potential delay fairness: The idea behind minimum potential delay fairness is based on the assumption that the involved flows are generated by applications transferring files of certain sizes. A relevant bandwidth sharing objective would be to minimize the time needed to complete those transfers. However, this does not apply to an adaptive streaming scenario and is therefore not discussed any further.

B. Perception Based Performance Metrics

There is a wide range of metrics which, influence how satisfied an end user is with a service such as e.g., video streaming. Many of these are not related to network aspects, and therefore difficult to influence by means in this domain. However, one of the perceived performance metrics which, could be correlated with network aspect is the notion of

perceived fairness. It is then of great interest to try and find methods of influencing this in a positive manner.

Looking at fairness from an end user perception, research from the social science and psychology domain [17] states that this is closely related to what is called 'Social Justice'. In this context a queuing system or any other resource allocation mechanism would be considered as a 'Social System'. It has further been found that users react negatively to any system behavior which, gives better service to other user, unless justification is provided. Such system behavior is considered un-fair, i.e., in violation with the social justice of the system as the end users considers it as discrimination.

The end user notion of system discrimination has been suggested by [18] as an important measure of perceived service quality, and more specifically the perceived fairness is stated to be closely related to the discrimination frequency. It should be noted that analyzing this type of end user perceived discriminations has a challenge in terms of handling the false positive and false negative cases.

Applying the concept of discrimination to competing adaptive streams, it would be related to situations where end user expectations are not met during steady state periods and also negative changes in service delivery during more transient periods. In other words, whatever makes the end user think that he is being discriminated due to other users in the system, will lead to reduced perceived service quality.

In order to use this type of perceived end user discrimination as a measure for how well the algorithm which, controls the adaptive streams are performing, a clear definition regarding what end users are considering as discrimination is required. This could, e.g., be periods with session rate below some threshold, any change in session rate to a lower level or the session rate change frequency.

C. Methods for Improving Performance

There are several things that one could try to incorporate into the adaptive algorithms controlling the ABR service [1] in order to make them perform better in a multi-stream scenario.

The selected performance metrics to be studied are from the network side proportional fairness, and from the end user side the perceived fairness metric as earlier described. Whether it is possible to improve both these fairness metrics at the same time will be an important part of the results.

Randomization of time intervals: The fixed rate adjustment intervals (T) used by each adaptive stream while in steady-state may be a contributing factor to inaccurate estimations of available bandwidth and thereby oscillating behavior. An alternative to fixed intervals would be to randomize them by using a per-session stochastic parameter (within certain reasonable bounds). By doing so the available bandwidth estimation methods may become more accurate.

Back-off periods: Whenever a service is reducing its rate level due to observed congestion it may try to increase again

after the same amount of time (T). In addition to the previous described randomization of this interval, one could also consider introducing a back-off period. This would imply that after a service has reduced its rate level, it enters a back-off period of a certain duration during which, no increase is allowed.

Threshold based behavior: Rather than using the same intervals of potential rate changes all the time, one could introduce a threshold for when it operates more or less aggressive. This threshold could be the mean available rate level for a specific session, or even a smoothed average value for the actual achieved level. This concept is applied with success in more recent TCP versions for the purpose of optimizing performance.

The method chosen for the simulations is according to the first approach described, i.e., a randomization of the intervals between each potential rate change as originally suggested in [3]. As baseline for the simulations, the fixed interval with $T=2s$ has been used. Then as stochastic alternatives, both a uniform distribution and a negexp distribution have been implemented. The uniform distribution gives values of T between [1.6, 2.4]s, while the negexp alternative gives values of T according to the distribution function with $\lambda=0.5$ and expected value $(1/\lambda) = 2s$.

IV. SIMULATION MODEL

As the adaptive streaming solutions of today are highly proprietary, the details concerning their implementation are not disclosed. Due to this, there will always be some degree of uncertainty concerning their internal functions.

A. Assumptions

One of the key functions of an ABR client is the method used for determining whether to go up or down in rate level during times of varying available bandwidth. From studying live traffic it does not seem as if the clients use additional network probing beyond the actual information obtained through download of video segments. Further on, in the likely absence of a per stream traffic shaper at the server side (for scalability and performance reasons), it will give a traffic pattern for each stream which, typically contains a sequence of busy and idle periods. The measured busy period rate is then higher than the actual stream rate level. Also, it is likely that there will be sub-periods within the busy periods where per packet rate is close to the total available bandwidth. As such, the client can probably obtain a rather accurate indication of maximum available bandwidth by just looking at minimum observed inter-arrival time of packets of known size belonging to the same stream.

However, not all streams will have interleaved busy periods so there is a good chance for each stream to overestimate the potential for additional bandwidth. There is a wide range of bandwidth estimation methods and a few of these are described in [19], but again - as the details of the

adaptive streaming solutions are not disclosed we will not discuss this part any further. Independent of which, method being used, there will be some degree of uncertainty which, contributes to variable performance. Further on, we assume the following to be true for the ABR sessions to be studied

- No stream coordination at server side
- No involvement from mechanisms in the network between the client and server
- All clients operate independently and do not communicate
- All clients are well behaved in the sense that they follow the same scheme
- At each defined stream rate level there are no variations due to i.e., picture dynamics

B. Session Type and Schedule

The ABR sessions used in the simulator are based on profiles observed in commercial services. The quality levels defined are {0, 250, 750, 1500, 2500, 3500, 5000} Kbps. All sessions are of the same type. The sessions are initiated by 10 different users and start time scheduling are done according to stochastic distributed parameters t_a – Uniform [0, 2000] ms and t_b – Uniform [0, 60] s. This gives that all sessions start during the first 60 seconds (t_b), but shifted by some milliseconds (t_a) in order to avoid synchronization of the rate adjustment intervals.

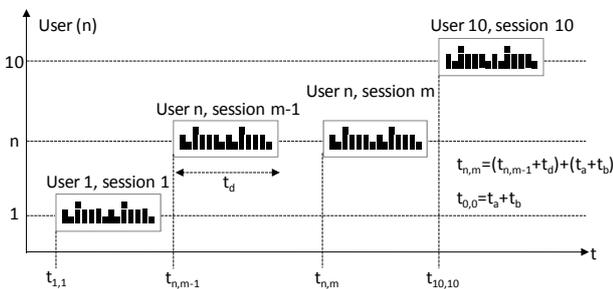


Figure 1. Session scheduling per user

During one simulation run, each user executes a total of 10 sessions sequentially. Time for starting the next session (m) for specific user (n) is noted $t_{n,m}$ (cf. Figure 1). The duration of each session t_d is deterministic and set to 40 minutes. A total of 10 simulation runs using different seeds are executed, corresponding to an aggregated session time of approximately 66 hours per user.

C. Rate Adaption Algorithm

The model for rate adaption per session is based on periodic estimation of available bandwidth $A_s(t)$ and calculation of a smoothed average $SA_s(t)$.

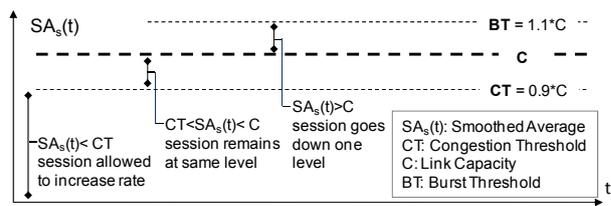


Figure 2. Thresholds for smoothed average

This smoothed average (cf. Figure 2) is compared to a congestion threshold (CT), the link capacity (C) and a burst threshold (BT) in order to trigger a rate adjustment.

Whenever the sum of requested rates from sessions is above the burst threshold (BT), the next session which, calculates $SA_s(t)$ will be forced down, independent of the value of $SA_s(t)$. This function is implemented in the simulator in order to incorporate the somewhat unpredictable behavior during times of heavy congestion.

The calculation of smoothed average $SA_s(t)$ is based on [3], and is expressed in (1). The parameter δ gives the weighting of the estimated available bandwidth for the two periods included in the calculation.

$$SA_s(t) = \delta A_s(t_{i-1}) + (1 - \delta)A_s(t_i) \quad (1)$$

The available bandwidth estimation function used in the simulations is based on the assumption that sessions running at high rates are able to make more accurate estimations than those running at lower rates. An abstraction of the function itself is made by a number of n bandwidth samples $C_{i,j}$ (cf. Figure 3)

A specific session is then given access to a number of these samples according to its current rate level, and then it will use this as basis for its estimation. A high rate gives a high number of samples available, and then, also, a higher degree of accuracy.

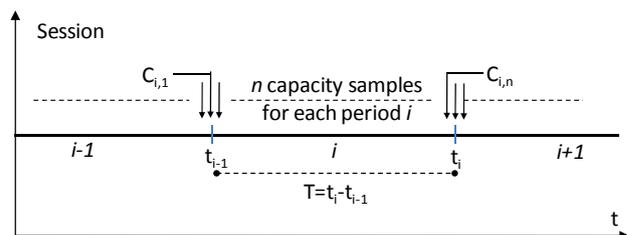


Figure 3. Capacity samples per period

The number of samples $x_{s,i}$ available to a specific session s for period i is given by its ratio between current rate $R_s(t_i)$ and max rate R_smax , multiplied by n as per (2).

$$x_{s,i} = n \frac{R_s(t_i)}{R_smax} \quad (2)$$

In the simulations, the value of n was set to 20 and R_s was according to the session definition 5000Kbps. The available bandwidth estimated $A_s(t)$ for period i is then given by the following (3).

$$A_s(t_i) = \sum_{l=1}^{x_i} C_{i,l} / x_{s,i} \quad (3)$$

By combination with the expression for $SA_s(t)$ it gives the following expression (4).

$$SA_s(t) = \delta \sum_{l=1}^{x_{i-1}} C_{i,l} / x_{s,i-1} + (1 - \delta) \sum_{l=1}^{x_i} C_{i,l} / x_{s,i} \quad (4)$$

The value of δ was set to 0.8 as per [3], thus giving most weight to the available bandwidth estimation from the previous period.

D. Simulation Tool

The simulator was built using the process oriented Simula [20] programming language and the Discrete Event Modeling On Simula (DEMOS) context class [21].

This programming language is considered as one of the first object oriented programming languages, and remains a strong tool for performing simulations.

V. RESULTS

The simulation results are presented for different congestion levels on the access link. The chosen capacities are 10, 20, 30 and 40Mbps. The lowest capacity would represent a highly congested scenario. The simulations were also run for all levels from 10 to 40 with increments of 200Kbps but for the sake of clarity these details are left out as they did not change the conclusions.

The studied fairness parameters (proportional and perceived), are compared for the 10 independent users sharing the access link. In order to present more information regarding variations in quality levels, a presentation of Coefficient of Variation (CV) is given. Values for CV below 1 is considered low-variance, while above 1 is considered high-variance.

The simulation results to be presented are based on that all users are accessing the same service, with identical session properties (i.e., quality levels). However, the simulations were also run for other service types and a mixture of services. These results are also left out, as they did not change the conclusion.

A. Proportional Fairness

Proportional fairness is measured as achieved session average rate per user, divided by session max – as per the definition earlier (cf. Figure 4, Figure 5, Figure 6). A high value is good and the maximum value is 1.

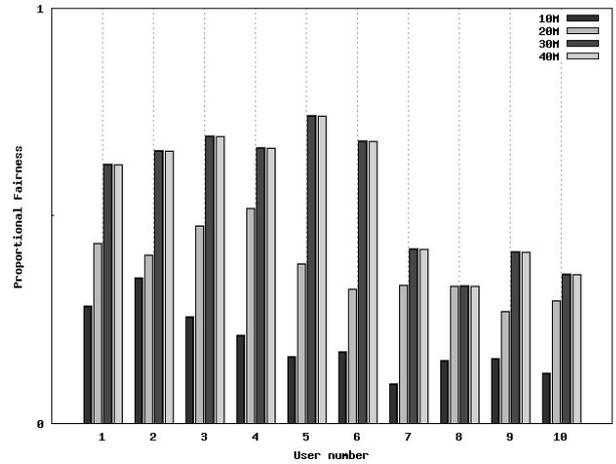


Figure 4. Proportional Fairness, fixed T=2s

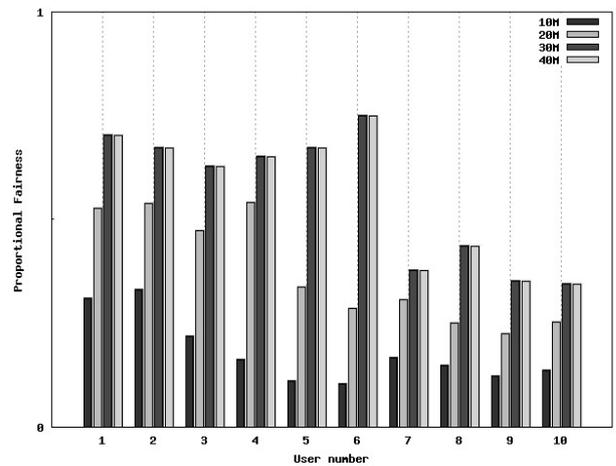


Figure 5. Proportional Fairness, Uniform T [1.6, 2.4]

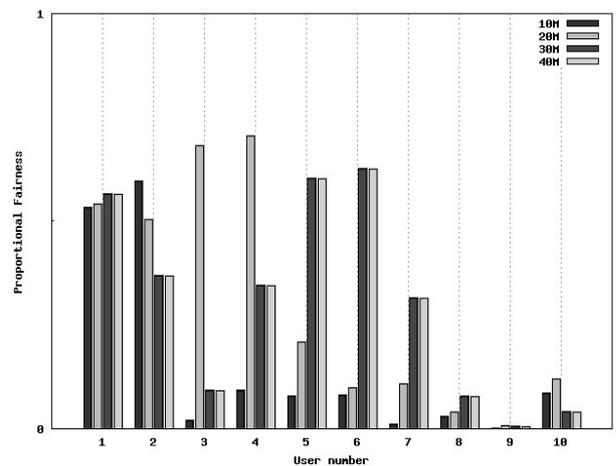


Figure 6. Proportional Fairness, negexp T [$\lambda=0.5$]

B. Perceived Fairness

The perceived fairness metric is calculated as the number of quality (rate) level reductions per minute (cf. Figure 7, Figure 8, Figure 9). Here, a low metric value is good – as it would reflect less rate reductions per minute.

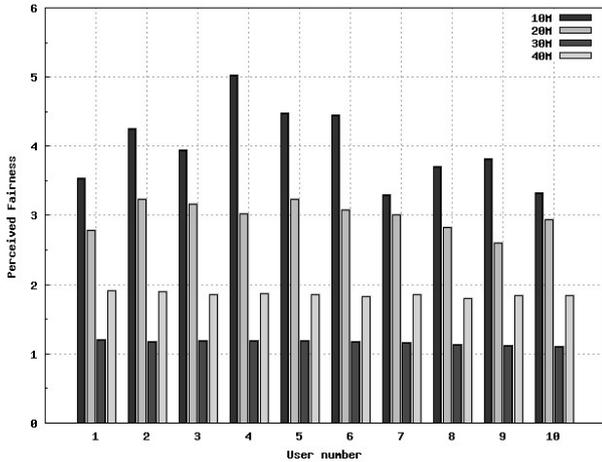


Figure 7. Perceived Fairness, fixed T=2s

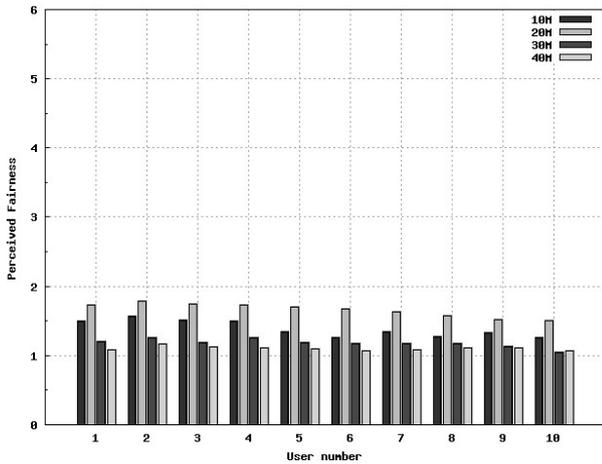


Figure 8. Perceived Fairness, Uniform T [1.6, 2.4]

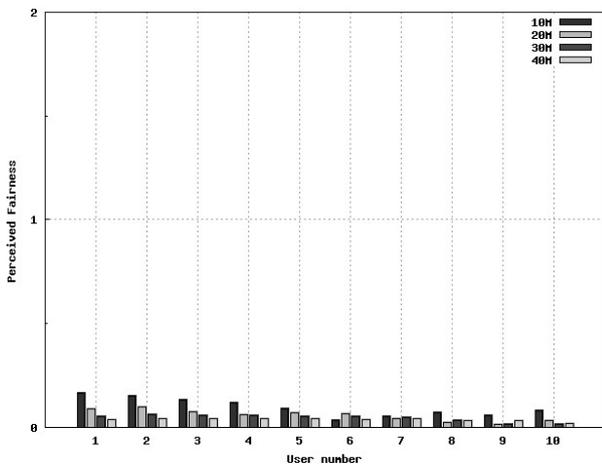


Figure 9. Perceived Fairness, negexp T [λ=0.5]

C. Coefficient of Variation (CV)

The Coefficient of Variation is calculated as Standard Deviation/Mean Value for sessions belonging to a user (cf. Figure 10, Figure 11, Figure 12). Values below 1 indicate low-variance which, is preferred.

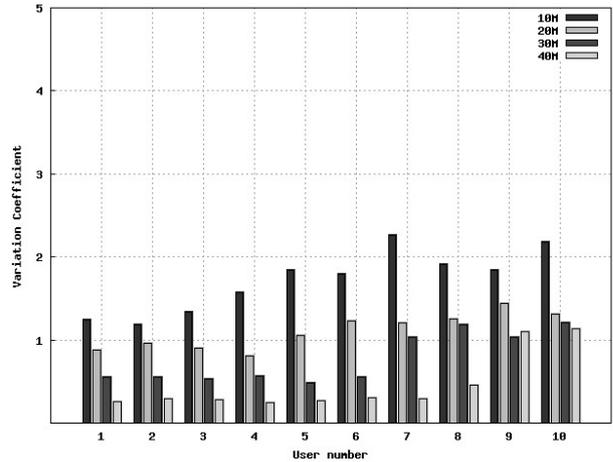


Figure 10. Coefficient of Variation, fixed T=2s

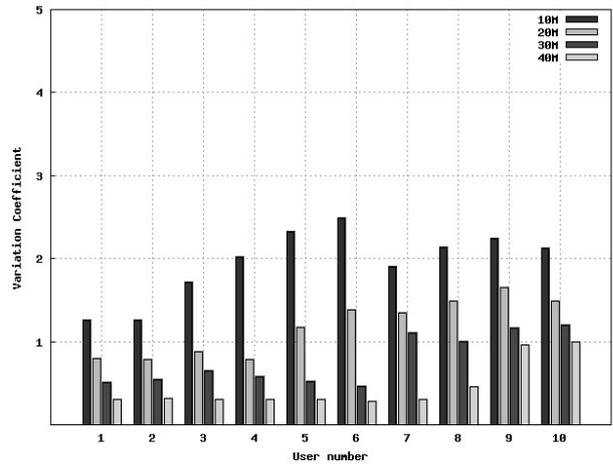


Figure 11. Coefficient of Variation, Uniform T [1.6, 2.4]

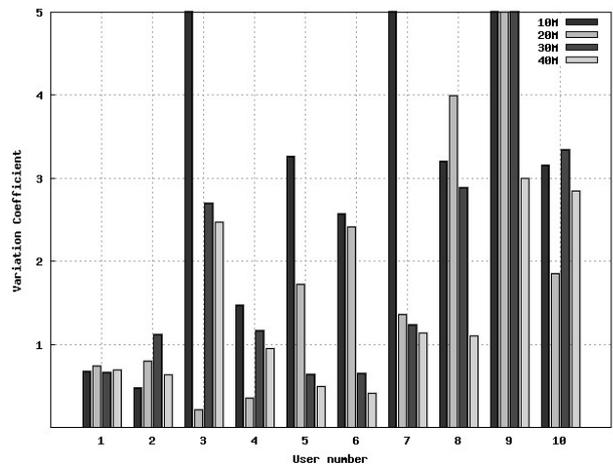


Figure 12. Coefficient of Variation, negexp T [λ=0.5]

VI. ANALYSIS

As expected, the randomization of time interval duration does have an effect on the parameters studied. However, the effect is not always positive.

Concerning proportional fairness, the introduction of a uniform T variable does not have a significant effect. The result can be viewed as neutral. On the other side, when the negexp T variable is used a clear negative effect is observed as the difference between users becomes significant.

For the perceived fairness metric, both the use of a uniform T and a negexp T have a significant positive effect. The best results are achieved for the negexp case which, gives values well below 1 for all congestion levels and users. It may be considered promising that the effect is especially strong during high times of high congestion (link capacity of 10M and 20M).

Regarding Coefficient of Variation, the results are similar to Proportional Fairness. A uniform T give no change, while a negexp T gives a negative change.

TABLE I. SUMMARY OF SIMULATION RESULTS

	<i>Proportional Fairness</i>	<i>Perceived Fairness</i>	<i>Coefficient Variation</i>
uniform T	neutral	positive	neutral
negexp T	negative	positive	negative

The somewhat intuitive explanation to why changes could be expected is that some of the negative effects of a fixed adjustment interval as illustrated in Figure 13 are reduced. In the case of fixed periods, each session would get the same periodic view on the link utilization, always missing or including some other traffic. This gives a certain degree of error in the available bandwidth estimation functions.

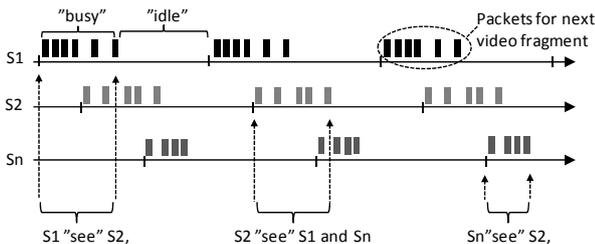


Figure 13. Problem with fixed estimation periods

Each session estimates available bandwidth only during its busy periods (ref Section III Subsection C). This means that in order to get an accurate estimation it is beneficial for it to have overlapping busy periods with as many other sessions as possible.

A. Burst Period Duration

The duration of the busy period for a specific session depends on both its current rate level and the rate

adjustment interval. The dependency of the rate level follows from the obvious relation to data volume to be transferred per time unit for a specific rate level, while the dependency of rate adjustment interval follows from the requirement to maintain the same average amount of data received over time.

At the beginning of each interval the client requests the next video fragment for a specific rate level, with duration equal to its rate adjustment interval. This is illustrated in Figure 14 where two sessions running at the same rate level, but with different rate adjustment intervals have different busy period durations.

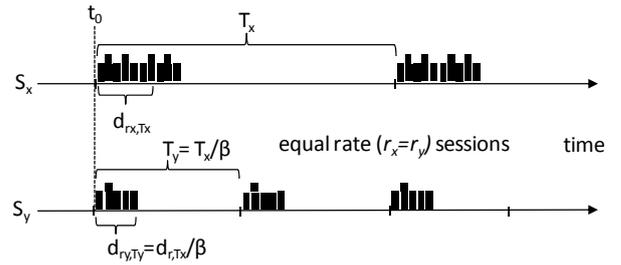


Figure 14. Equal rate sessions with different busy periods

Any two sessions (S_x , S_y) running at the same rate level, will have a relation between their burst period durations expressed by the parameter β . This parameter is given by the following expression (5).

$$\beta = \frac{d_{r_x, T_x}}{d_{r_y, T_y}} = \frac{T_x}{T_y}, \quad \text{for } r_x = r_y \quad (5)$$

Using this relationship, we can express (6) the burst period duration d_{r_i, T_i} for any session S_i as a function of its rate adjustment interval T_i and a reference burst period duration $d_{r_i, T}$.

$$d_{r_i, T_i} = d_{r_i, T} \left(\frac{T_i}{T} \right) \quad (6)$$

The values for $d_{r_i, T}$ can presumably be calculated based on information about the codec used for the specific media stream inside each sessions, together with assumption on per session server side capacity. Alternatively one could make measurements on a specific system and establish a $d_{r_i, T}$ matrix for all valid values of r_i and the reference T value.

However, if we assume that the server side capacity is not a limitation, and that it will always try to burst with a certain bitrate C_{burst} we can also express the burst period duration d_{r_i, T_i} as follows (7).

$$d_{r_i, T_i} = \left(r_i T / C_{burst} \right) \left(T_i / T \right) = \left(r_i T_i / C_{burst} \right) \quad (7)$$

The maximum value for C_{burst} is natural to think of as the access capacity for the user group / home network, as this is normally the end-to-end bottleneck. However, it is likely

that the actual C_{burst} is related to the maximum rate for the specific service.

B. Probability for Burst Period Overlap

For T_i values according to a uniform distribution, the probability $P_{i,r,t}$ for a session i at rate level r to be in its busy period at time t will be according to the following expression (8).

$$P_{i,r,t} = d_{r,T_i}/T_i = \frac{d_{r,T}(T_i/T)}{T_i} = d_{r,T}/T \quad (8)$$

From this, we see that all sessions at a specific rate level has the same probability of being in its busy period at time t . We can then express the probability that all n sessions are in their busy period at time t as follows (9).

$$P_{all\ busy,t} = \left(d_{r_1,T}/T\right)^{c_1} \left(d_{r_m,T}/T\right)^{c_m} \quad (9)$$

The parameter c_m represents the number of sessions at rate level r_m and the sum of all c_m values equals n . From this we see that the probability of any session to see all other sessions during its busy period depends on the session rate level mix, and this probability increases when more sessions are running at high rate levels.

Further on, we recognize that the probability for that a session i has an overlap with each of the other sessions sometimes during its busy period T_i is the integral of $P_{all\ busy,t}$ over the period $[0, T_i]$ which, is easily expressed as the constant $P_{all\ busy,t}$ multiplied by T_i .

We then let a specific session mix be described by the vector $R_{mix}=\{r_1, \dots, r_n\}$, whereas r_i represents the rate level for session i . Also, for a specific session i let A_i be the group of sessions which, has overlapping busy periods with session i at a specific time t_0 , and B_i be the group of sessions for which, it did not have an overlap. In the situation where all sessions have the same rate adjustment interval duration T_i , the probability of that session i has an overlapping busy period with any of the sessions in group B_i at time t_0+T_i is zero. This leads to that while R_{mix} remains unchanged, the view a specific session has of the total traffic will not change. The system state for session i in terms of busy period overlap with other sessions is independent of the state at t_0 and also t in general.

In the case where T_i is not equal for all sessions, but instead are chosen according to some stochastic distribution – the group of sessions which, overlap the busy period of session i at t_0+T_i is not independent of the state at t_0 . If we let C_i denote the sub-group of sessions from B_i which, has overlapping busy periods with session i at time t_0+T_i , it can be shown that there is a deterministic relationship between A_i , B_i and C_i .

If we then remember the assumed use of a smoothed average function we see the benefit of this potential additional burst period overlaps in subsequent periods.

C. Dynamics in Burst Period Overlap

When the starting times for each session and their respective rate adjustment intervals (T_i) are considered stochastic processes, the sessions will combine in time in different ways. In order to define the deterministic relationship between overlapping busy periods during subsequent intervals, we need to analyze scenarios where sessions with different rate levels and different rate adjustment interval are combined.

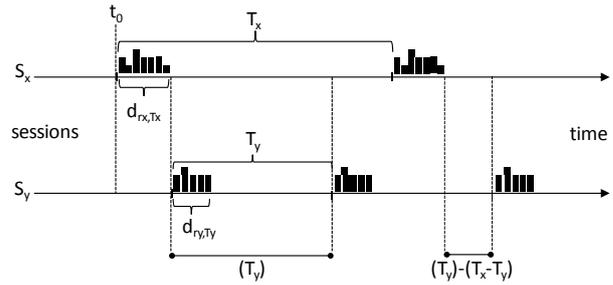


Figure 15. Session S_y starting after S_x ($T_y < T_x$)

The first scenario (a) to be studied is the one where two sessions (S_x, S_y) with different T_i values (T_x, T_y) are active at the same time. We assume $T_x > T_y$ and that S_y starts immediately after the busy period of S_x finishes as illustrated in Figure 15.

For the two sessions (S_x, S_y) there will be shift in phase between them as a function of time which, makes them have a full or partial busy period overlap at some time. The question is then how many rounds it will take for S_x to see S_y and vice versa. It can be shown that we can express the number of rounds for S_x before it has an overlapping busy period with S_y as follows (10).

$$N_{a,x \rightarrow y} = 1 + \left\lceil \frac{T_y}{T_x - T_y} \right\rceil$$

when $\frac{T_x}{2} < T_y < (T_x - d_{r_x, T_x} - d_{r_y, T_y})$ (10)

$$N_{a,x \rightarrow y} = 2$$

when $(T_x - d_{r_x, T_x} - d_{r_y, T_y}) < T_y < T_x$

In the same way, we can express the number of rounds for S_y before the same overlap of busy period with S_x takes place (11).

$$N_{a,y \rightarrow x} = 1 + \left\lceil \frac{T_x}{T_x - T_y} \right\rceil$$

when $\frac{T_x}{2} < T_y < (T_x - d_{r_x, T_x} - d_{r_y, T_y})$ (11)

$$N_{a,y \rightarrow x} = 2$$

when $(T_x - d_{r_x, T_x} - d_{r_y, T_y}) < T_y < T_x$

The next scenario (b) to be studied is where the sessions (S_x, S_y) are running with different T_i values (T_x, T_y) but now S_y finishes its busy period before S_x (cf. Figure 16).

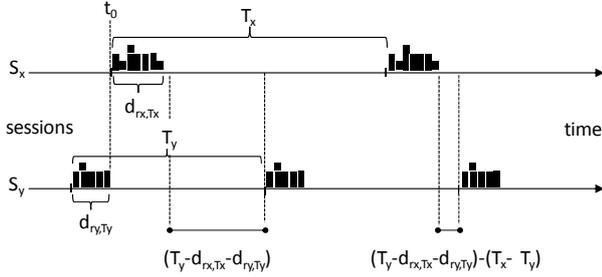


Figure 16. Session S_x starting after S_y ($T_y < T_x$)

The number of rounds it takes for S_x to see S_y is expressed as follows (12).

$$N_{b,x \rightarrow y} = 1 + \left\lceil \frac{T_y - d_x - d_y}{T_x - T_y} \right\rceil$$

when $\frac{T_x}{2} < T_y < (T_x - d_{r,Ty})$ (12)

$$N_{b,x \rightarrow y} = 2$$

when $(T_x - d_{r,Ty}) < T_y < T_x$

The number of rounds it takes for S_y to see S_x is expressed as follows (13).

$$N_{b,y \rightarrow x} = 1 + \left\lceil \frac{T_x - d_x - d_y}{T_x - T_y} \right\rceil$$

when $\frac{T_x}{2} < T_y < (T_x - d_{r,Ty})$ (13)

$$N_{b,y \rightarrow x} = 2$$

when $(T_x - d_{r,Ty}) < T_y < T_x$

It should be noted that for both scenarios there is a special case where $N_{a,y \rightarrow x}/N_{b,y \rightarrow x}$ and $N_{a,x \rightarrow y}/N_{b,x \rightarrow y}$ are always 2, i.e., two sessions which, did not have overlapping busy periods at t_0 is guaranteed to have overlapped during the next period for S_x and S_y . For a smoothed average function operating over two periods this is desirable, i.e., whatever it does not see in the first period it is guaranteed to see in the next.

D. Optimization Problem

The expressions for $N_{y \rightarrow x}$ and $N_{x \rightarrow y}$ contain many variables. These variables are the rate adjustment intervals T_i and the burst period durations d_{r_i, T_i} for all sessions. The latter are calculated based on the session rates r_x and r_y and C_{burst} as defined in Section V. These expressions can be used as input to a constrained optimization problem and analyzed as such in order to find maximum and minimum values.

As the starting point for this optimization problem we can focus on the worst case scenario, that would be the number of rounds for S_y before it has an overlap with S_x ($N_{a,y \rightarrow x}/N_{b,y \rightarrow x}$), which, will always be higher than the number of rounds for S_x before this has an overlap with S_y .

We also see that $N_{a,y \rightarrow x}$ will always be greater than $N_{b,y \rightarrow x}$ since $T_x > T_y$. This gives us only one expression to analyze for the worst case scenario as follows (14).

Maximize: $N_{a,y \rightarrow x}$

where

$$N_{a,y \rightarrow x} = \begin{cases} 1 + \left\lceil \frac{T_x}{T_x - T_y} \right\rceil, & \text{if } T_y < (T_x - d_{r_x, T_x} - d_{r_y, T_y}) \\ 2, & \text{if } (T_x - d_{r_x, T_x} - d_{r_y, T_y}) < T_y < T_x \end{cases}$$

subject to: (14)

$$1.6 < T_y, T_x < 2.4 \text{ and } T_x/2 < T_y$$

$$R_x, R_y \in \{250, 750, 1500, 2500, 3500, 5000\}$$

$$d_{r_x, T_x} = r_x T_x / C_{burst}$$

$$d_{r_y, T_y} = r_y T_y / C_{burst}$$

The above maximization can then be done for different values of C_{burst} . In the simulations the access speeds used were between 10 and 40Mbps and the maximum session rate was 5Mbps. Based on measurements of real traffic we can see that the C_{burst} is lower than the actual access speed and therefore values of respectively 5Mbps, 7.5Mbps and 10Mbps were used for C_{burst} .

For the two different alternatives of choosing values for T_i used in the simulations, the uniform approach is easiest to work with in the optimization context since it gives a min and max value for T_i . For the negexp alternative the corresponding range would be $[0, \infty]$ and for this scenario the optimization problem does not have a useful solution.

The result from solving the optimization problem is shown in Figure 17. The three different burst bitrates (C_{burst}) give surfaces which, are plotted, whereas the highest capacity gives the highest values for $N_{a,y \rightarrow x}$.

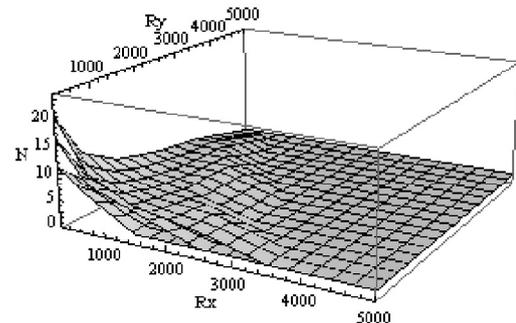


Figure 17. Maximum $N_{a,y \rightarrow x}$ for different burst bitrates

We see that in many cases we get an overlap already in the second round, and thereby we improve the basis for the

available bandwidth estimation algorithm. This analysis then strengthens the findings of the simulations.

The more likely explanation to the negexp behavior in terms of good perceived fairness is the somewhat extreme proportional unfairness. By allowing some sessions to be very greedy, one prevents others from increasing at all. This is a stable but proportionally very unfair situation.

In order to improve the available bandwidth estimations further one may consider the well known PASTA principle [22] from queuing theory which, states that a Poisson based Arrival process See Time Averages. This implies that the bandwidth probing should take place not only during the burst periods, but as a process taking samples throughout the whole rate adjustment period.

VII. CONCLUSIONS AND FUTURE WORK

The results show that there is a significant potential of improving perceived fairness as defined and associated QoE for adaptive streams of the category studied. The positive effect of the suggested enhancement to the rate adaption scheme, i.e., using a stochastically determined duration of rate adjustment intervals rather than fixed values is supported by the simulation results and theoretical analysis.

The results also illustrate that when studying the performance of adaptive streaming solutions, it is not enough to only focus on the network centric QoS domain. A change in this domain does not necessary lead to a corresponding change in the QoE domain, and vice versa. The significant improvement in Perceived Fairness, while proportional fairness remained the same for the uniform T case supports this statement.

As future work in this field it is planned to further study objective and no-reference based QoE metrics such as Perceived Fairness which, is possible to correlate over to the QoS and network domain. It is also planned to verify the simulation and analytical results by means of measurements.

VIII. ACKNOWLEDGEMENTS

The reported work is done as part of the Road to media-aware user-Dependant self-adaptive NETWORKS - R2D2 project. This project is funded by The Research Council of Norway. The work has also been actively supported by TV2, the leading commercial TV broadcaster in Norway. TV2 is among the pioneers in providing a full commercial TV offering over the Internet based on ABR technology.

REFERENCES

- [1] A. Zambelli, "IIS smooth streaming technical overview," <http://www.microsoft.com/silverlight/whitepapers/>, Tech. Rep., March 2009 (last accessed 10.01.2012).
- [2] E. Areizaga, L. Perez, C. Verikoukis, N. Zorba, E. Jacob, and P. Odling, "A road to media-aware user-dependent self-adaptive networks," in *Proc. IEEE Int. Symp. BMSB '09*, 2009, pp. 1–6.
- [3] S. Akhshabi, A. C. Begen, and C. Dovrolis, "An experimental evaluation of rate-adaptation algorithms in adaptive streaming over http," in *ACM Multimedia Systems (MMSys)*, 2011.
- [4] L. De Cicco and S. Mascolo, "An experimental investigation of the akamai adaptive video streaming," in *Proceedings of the 6th international conference on HCI in work and learning, life and leisure: workgroup human-computer interaction and usability engineering*, ser. USAB'10. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 447–464.
- [5] R. Kuschnig, I. Kofler, and H. Hellwagner, "An evaluation of tcp-based rate-control algorithms for adaptive internet streaming of h.264/svc," in *MMSys '10*. New York, NY, USA: ACM, 2010, pp. 157–168.
- [6] B. J. Villa and P. E. Heegaard, "Monitoring and control of QoE in media streams using the click software router," in *NIK2010, Norway*. ISBN 978-82-519-2702-4., vol. 1, November 2010, pp. 24–33.
- [7] B. J. Villa and P. E. Heegaard, "Towards knowledge-driven QoE optimization in home gateways," *ICNS 2011*, May 2011.
- [8] S. Bhatti, M. Bateman, and D. Miras, "Revisiting inter-flow fairness," in *Broadband Communications, Networks and Systems, 2008. BROADNETS 2008. 5th International Conference on*, sept. 2008, pp. 585–592.
- [9] S. Floyd, "Metrics for the Evaluation of Congestion Control Mechanisms," RFC 5166 (Informational), Internet Engineering Task Force, Mar. 2008.
- [10] R. Jain, D. Chiu, and W. Hawe, "A quantitative measure of fairness and discrimination for resource allocation in shared systems," DEC, Tech. Rep., 1984.
- [11] D. Mitra and J. B. Seery, "Dynamic adaptive windows for high speed data networks with multiple paths and propagation delays," *Computer Networks and ISDN Systems*, vol. 25, no. 6, pp. 663–679, 1993, high Speed Networks.
- [12] Y. Zhang, S.-R. Kang, and D. Loguinov, "Delayed stability and performance of distributed congestion control," in *Proceedings of the 2004 conference on Applications, technologies, architectures, and protocols for computer communications*, ser. SIGCOMM '04. New York, NY, USA: ACM, 2004, pp. 307–318.
- [13] Z. Cao and E. Zegura, "Utility max-min: an application-oriented bandwidth allocation scheme," in *INFOCOM '99*, vol. 2, mar 1999, pp. 793–801 vol.2.
- [14] F. Kelly, "Charging and rate control for elastic traffic," *European Transactions on Telecommunications*, 1997.
- [15] S. Kunniyur and R. Srikant, "End-to-end congestion control schemes: utility functions, random losses and ecn marks," *IEEE/ACM Trans. Netw.*, vol. 11, pp. 689–702, October 2003.
- [16] A. Jdidi and T. Chahed, "Flow-level performance of proportional fairness with hierarchical modulation in ofdma-based networks," *Comput. Netw.*, vol. 55, pp. 1784–1793, June 2011.
- [17] H. Levy, B. Avi-Itzhak, and D. Raz, "Network performance engineering," D. D. Kouvatso, Ed. Berlin, Heidelberg: Springer-Verlag, 2011, ch. Principles of fairness quantification in queueing systems, pp. 284–300.
- [18] W. Sandmann, "Quantitative fairness for assessing perceived service quality in queues," *Journal Operational Research*, pp. 1–34, April 2011.
- [19] R. Prasad, C. Dovrolis, M. Murray, and K. Claffy, "Bandwidth estimation: metrics, measurement techniques, and tools," *Network, IEEE*, vol. 17, no. 6, pp. 27–35, nov.-dec. 2003.
- [20] R. J. Pooley, *An Introduction to Programming in Simula*. Blackwell Scientific Publications. ISBN: 0632014229, 1987.
- [21] G. Birtwistle, *Demos - A system for Discrete Event Modelling on Simula*, G. Birtwistle, Ed. School of Computer Science, University of Sheffield, July 1997.
- [22] R. W. Wolff, "Poisson Arrivals See Time Averages," *Operations Research*, vol. 30, no. 2, pp. 223–231, 1982.