

Buffer-based algorithm implementation to Rate Adaptation

Piergiorgio Ladisa
Nicola Sebastianelli

22/11/2018

Context

These days we are witnessing to the continuous growing demand of streamed content in network: anyone asks anywhere and in any-moment for contents to servers from their mobile devices. Hence, since the percentage of contents in the network is becoming everyday bigger, CDNs and ISPs are called to respond to this demand. The paradigm of the multimedia systems and for CDN is mainly the so-called Quality of Experience (QoE). This is very different from the paradigm on which Internet was designed, i.e. Quality of Service (QoS) and Best Effort, in fact the quality of experience is really subjective and very difficult to measure. Nevertheless, different protocols and strategies have been proposed in order to met the demand of users of streaming of contents. Firstly, different strategies of streaming exists:

- all-in-once, i.e. the content is entirely sent to the user. This approach is enforced by the server and the encoding rate is unintentional, since the client's buffer is filled up, in the same way in which the TCP buffer is filled up; since TCP performs the flow control, the client with this approach read at the encoding rate;
- throttling, i.e. the rate is adapted to the user with two possible approach:
 - on-off-S, in which the connection is maintained persistently;
 - on-off-M, in which the connection is not maintained persistently

These two approaches are caused by the client application, that periodically stops reading from the TCP socket. In a general case of pseudo-streaming, the video is buffered until the buffer is enough full and, once reached the steady state, the controller starts doing an on-off step, downloading the block size of the video.

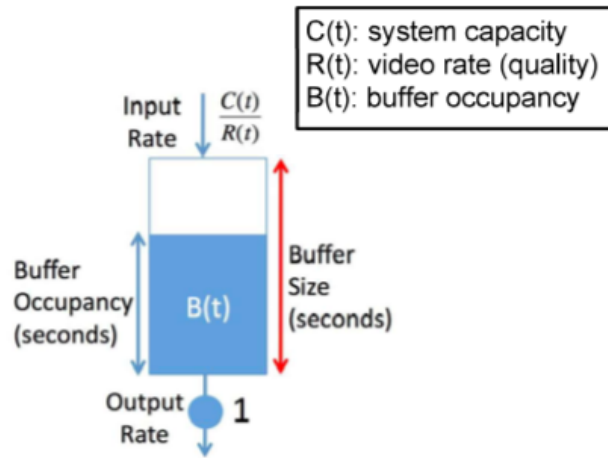
In this context, the mainly adopted protocol is HTTP Adaptive Streaming protocol (HAS). Here the idea is to make the requested bitrate of the video fit with the varying network resources, ensuring the best possible quality of experience for the user, based on the current situation of the network. This adaptation is usually done on client-side. Here the content is divided in chunks and available in different qualities. So HAS allows to encode each segment in different quality referring to the current bandwidth and so the stream is splitted into a sequence of segments, instead of downloading one entire large file (as in the all-in-once approach). The policies of HAS are:

- Rate-based, that select the highest possible video representation referring to the measured speed from the previous received chunk;
- Buffer-based, that uses different thresholds for the buffer in a way in which the more the buffer is filled, the more quality of the video is. Thus, with this approach, the quality can switch step-by-step at each new request of a video segment: it is impossible to jump from the lowest to the highest quality just in one step;
- Buffer and Rate based. This is an hybrid approach, i.e. is based on the rate-based approach basically, but the previous chunk download speed is

weighted using a factor that depends on the amount of the saturation of the buffer: if the buffer is depleting the previous chunk, the download speed will be considered to be less than the measured one and hence a lower quality is selected.

A buffer-based approach to rate adaptation

When the client want to access to a video-content for streaming, it chooses which the video rate to stream by monitoring network conditions and estimating the available network capacity. This process is referred to as *adaptive bit rate selection* or *ABR*. ABR algorithms try to balance two opposite goals. The first goal is to maximize the video quality choosing as video rate the highest supportable by the network. The second goal is to minimize re-buffering events, which cause the video to halt if the client's playback buffer goes empty. Data is requested in chunks and the buffer drains at 1 unit/second. So if the video-rate $R(t)$ is greater than the system capacity $C(t)$, new data is added at $C(t)/R(t) < 1$, depleting the buffer.



In order to maximize video quality, a service could just stream at the maximum video rate R_{max} all the time, but this would risk extensive re-buffering. On the other hand, to minimize rebuffering, the service could just stream at the minimum video rate R_{min} all the time but this extreme would lead to low video quality. Hence, the design goal of an ABR algorithm is to simultaneously obtain high performance on both metrics in order to give users a good quality of experience. Firstly, the client measures how fast chunks arrive to estimate capacity, let us say $C(t)$. The estimate enriched with knowledge of

the buffer occupancy, which is represented with an adjustment factor $F(B(t))$, i.e. a function of the playback buffer occupancy. So, the selected video rate is $R(t) = F(B(t))C(t)$; different designs use different adjustment functions $F(\cdot)$. When the buffer contains many chunks, $R(t)$ can safely deviate from $C(t)$ without triggering a rebuffer. The client can aggressively try to maximize the video quality by picking $R(t) = C(t)$. But when the buffer is low, the client should be more conservative.

We say that an ABR algorithm is *buffer-based* if it picks the video rate as a function of the current buffer occupancy $B(t)$. The region between:

- $[0, B_{max}]$ on the buffer-axis;
- $[R_{min}, R_{max}]$ on the rate-axis;

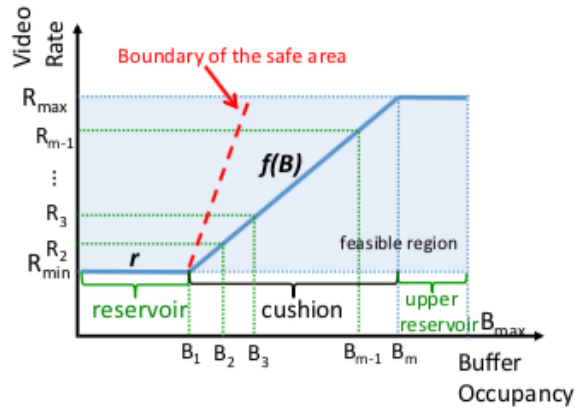
defines the feasible region. So, any curve $f(B)$ on the plane within the feasible region defines a *rate map*, i.e. a function that produces a video rate between R_{min} and R_{max} given the current buffer occupancy. Let us assume that:

- the chunk-size is infinitesimal;
- video rates within $[R_{min}, R_{max}]$ are continuous;
- videos are encoded at a constant bitrate;
- videos are infinitely long;

Thus, as long as $C(t) \geq R_{min}$ and $f(B)$ tends to R_{min} when B tends to 0, there will not be unnecessary rebuffering events. So, as long as $f(B)$ increases to R_{max} , the average video rate will match the average capacity when it is true that $R_{min} < C(t) < R_{max}$.

BBA0-Algorithm

The rate map implemented by the BBA0 algorithm is the following:



This algorithm add reservoir r to pad for finite chunk sizes and some $C(t)$ variation. So, while filling reservoir, it request R_{min} and must have r greater

or equal to the minimum chunk size. In addition, an upper reservoir is added to reach R_{max} before the buffer B is full and above the "safe area" is where the buffer will not deplete into r if $C(t)$ suddenly drops, but does not reach values lower than R_{min} .

BBA0 Algorithm

Input:

- $Rate_{prev}$: previously used video;
- Buf_{now} : current buffer occupancy;
- r : size of reservoir;
- cu : size of cushion.

Output: $Rate_{next}$: next video rate.

```

if Rate_prev = R_max then
    Rate+ = R_max
else
    Rate+ = min{Ri : Ri > Rate_prev}

if Rate_prev = R_min then
    Rate- = R_min
else
    Rate- = max{Ri : Ri < Rate_prev}

if Buf_now <= r then
    Rate_next = R_min

else if Buf_now >= (r + cu) then
    Rate_next = R_max

else if f(Buf_now) >= Rate+ then
    Rate_next = max{Ri : Ri < f(Buf_now)};

else if f(Buf_now) <= Rate- then
    Rate_next = min{Ri : Ri > f(Buf_now)};

else
    Rate_next = Rate_prev;

return Rate_next;

```