

Iraj Sodagar
Microsoft
Corporation

The MPEG-DASH Standard for Multimedia Streaming Over the Internet

Watching the Olympics live over the Internet? Streaming last week's episode of your favorite TV show to your game console? Watching a 24-hour news TV channel on your mobile phone? These use cases might already seem possible as part of our daily lives. In fact, during the 2008 Olympics, NBC reported delivering 3.4 petabytes of video content over the Internet.¹ The truth is, however, that multimedia streaming over the Internet is still in its infancy compared to its potential market. One reason is that today every commercial platform is a closed system with its own manifest format, content formats, and streaming protocols. In other words, no interoperability exists between devices and servers of various vendors. A recent study indicated that in a few years video content could make up the vast majority of Internet traffic.² One of the main enablers of this would be an adopted standard that provides interoperability between various servers and devices. Achieving such interoperability will be instrumental for market growth, because a common ecosystem of content and services will be able to provision a broad range of devices, such as PCs, TVs, laptops, set-top boxes, game consoles, tablets, and mobiles phones. MPEG-Dynamic Adaptive Streaming (DASH) was developed to do just that.

HTTP streaming

Delivering video content over the Internet started in the 1990s with timely delivery and consumption of large amounts of data being the main challenge. The Internet Engineering Task Force's Real-Time Transport Protocol (RTP) was designed to define packet formats for audio and video content along with stream-session management, which allowed delivery of those packets with low overhead. RTP works well in managed IP networks. However, in today's Internet, managed networks have been replaced by content delivery networks (CDN), many of which don't support RTP streaming. In addition, RTP packets are often not allowed through firewalls. Finally, RTP streaming requires the server to manage a separate streaming session for each client, making large-scale deployments resource intensive.

With the increase of Internet bandwidth and the tremendous growth of the World Wide Web, the value of delivering audio or video data in small packets has diminished. Multimedia content can now be delivered efficiently in larger segments using HTTP. HTTP streaming has several benefits. First, the Internet infrastructure has evolved to efficiently support HTTP. For instance, CDNs provide localized edge caches, which reduce long-haul traffic. Also, HTTP is firewall friendly because almost all firewalls are configured to support its outgoing connections. HTTP server technology is a commodity and therefore supporting HTTP streaming for millions of users is cost effective. Second, with HTTP streaming the client manages the streaming without having to maintain a session state on the server. Therefore, provisioning a large number of streaming clients doesn't impose any additional cost on server resources beyond standard Web use of HTTP,

Editor's Note

MPEG has recently finalized a new standard to enable dynamic and adaptive streaming of media over HTTP. This standard aims to address the interoperability needs between devices and servers of various vendors. There is broad industry support for this new standard, which offers the promise of transforming the media-streaming landscape.

—Anthony Vetro

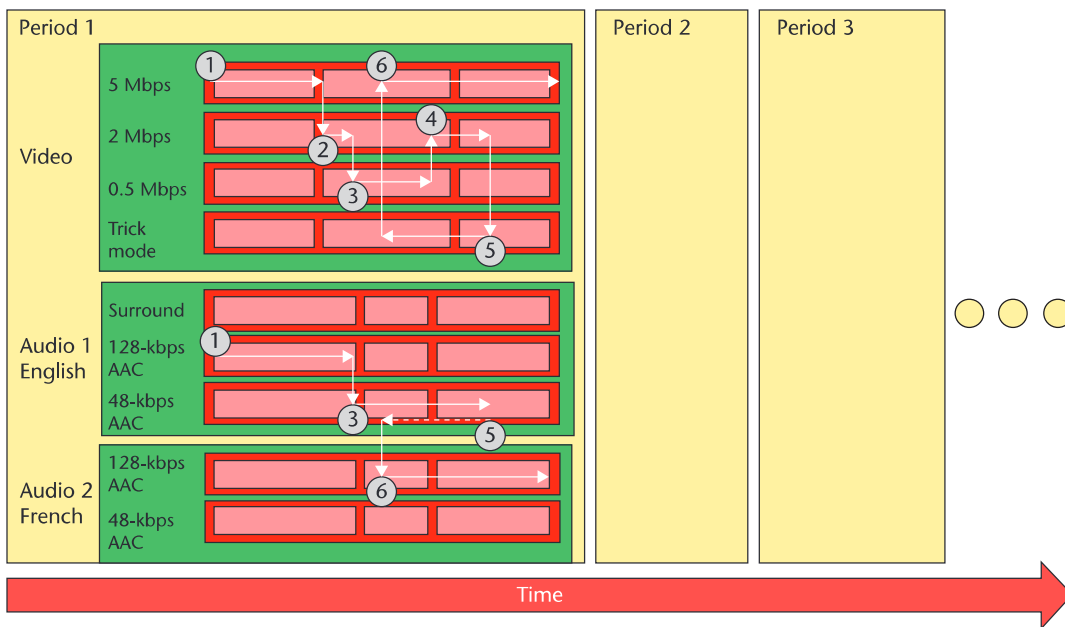


Figure 1. Simple example of dynamic adaptive streaming. Numbered circles demonstrate the action points taken by the device.

and can be managed by a CDN using standard HTTP optimization techniques.

For all of these reasons, HTTP streaming has become a popular approach in commercial deployments. For instance, streaming platforms such as Apple's HTTP Live Streaming,³ Microsoft's Smooth Streaming,⁴ and Adobe's HTTP Dynamic Streaming (see <http://www.adobe.com/products/httpdynamicstreaming>) all use HTTP streaming as their underlying delivery method. However, each implementation uses different manifest and segment formats and therefore, to receive the content from each server, a device must support its corresponding proprietary client protocol. A standard for HTTP streaming of multimedia content would allow a standard-based client to stream content from any standard-based server, thereby enabling interoperability between servers and clients of different vendors.

Observing the market prospects and requests from the industry, MPEG issued a Call for Proposal for an HTTP streaming standard in April 2009. Fifteen full proposals were received by July 2009, when MPEG started the evaluation of the submitted technologies. In the two years that followed, MPEG developed the specification with participation from many experts and with collaboration from other standard groups, such as the Third Generation Partnership Project (3GPP).⁵ The resulting standard, known as MPEG-DASH over HTTP, is currently at the Draft International Standard stage.⁶ Note that,

at the time of publishing this article, only the referenced draft is publically available. The specification was further revised in August 2011 and is expected to be published as ISO/IEC 23009-1.

A simple case of adaptive streaming

Figure 1 illustrates a simple example of on-demand, dynamic, adaptive streaming. In this figure, the multimedia content consists of video and audio components. The video source is encoded at three different alternative bitrates: 5 Mbytes, 2 Mbytes, and 500 kilobits per second (Kbps). Additionally, an I-frame-only bitstream with a low frame rate is provided for streaming during the trick-mode play. The accompanying audio content is available in two languages: audio 1 is a dubbed English version of the audio track and is encoded in surround sound, Advanced Audio Coding (AAC) with 128-Kbyte and 48-Kbps alternatives; while audio 2 is the original French version, encoded in AAC 128-Kbyte and 48-Kbps alternatives only.

Assume that a device starts streaming the content by requesting segments of the video bitstream at the highest available quality (5 Mbytes) and the English audio at 128 Kbytes AAC because, for instance, the device doesn't support surround audio (label 1 in Figure 1). After streaming the first segments of video and audio, and monitoring the effective network bandwidth, the device realizes that the actual available bandwidth is lower than

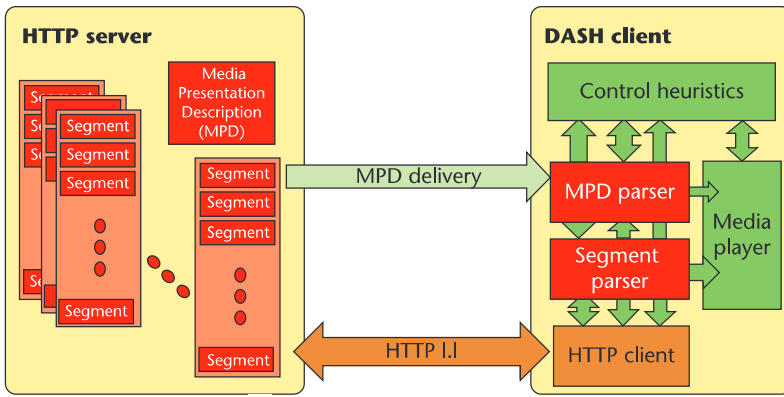


Figure 2. Scope of the MPEG-DASH standard. The formats and functionalities of the red blocks are defined by the specification. The clients control heuristics and media players, which aren't within the standard's scope.

5 Mbps. So, at the next available switching point, it switches the video down to 2 Mbps by streaming the next segments from the mid-quality track while continuing streaming of the 128-Kbyte AAC English audio (label 2 in Figure 1). The device continues to monitor the actual network bandwidth and realizes that the network bandwidth has further decreased to a value lower than 2 Mbps. Therefore, to maintain continuous playback, the device further switches the streams down to 500-Kbps video and 48-Kbps audio (label 3 in Figure 1). It continues playing the content at these rates until the network bandwidth increases and then it switches the video up to 2 Mbytes (label 4 in Figure 1). After a while, the user decides to pause and rewind. At this point, the device starts streaming the video from the trick-mode track to play the video in reverse order, while audio is muted (label 5 in Figure 1). At the desired point, the user clicks to play the content with the original French audio. At this point, the device resumes streaming the video from the highest quality (5 Mbytes) and audio from 128-Kbyte French audio (label 6 in Figure 1).

This example perhaps is one of the most simple use cases of dynamic streaming of multimedia content. More advanced use cases might include switching between multiple camera views, 3D multimedia content streaming, video streams with subtitles and captions, dynamic ad insertion, low-latency live streaming, mixed-streaming and prestored content playback, and others.

Scope of MPEG-DASH

Figure 2 illustrates a simple streaming scenario between an HTTP server and a DASH client. In this figure, the multimedia content is captured and stored on an HTTP server and is

delivered using HTTP. The content exists on the server in two parts: Media Presentation Description (MPD), which describes a manifest of the available content, its various alternatives, their URL addresses, and other characteristics; and segments, which contain the actual multimedia bitstreams in the form of chunks, in single or multiple files.

To play the content, the DASH client first obtains the MPD. The MPD can be delivered using HTTP, email, thumb drive, broadcast, or other transports. By parsing the MPD, the DASH client learns about the program timing, media-content availability, media types, resolutions, minimum and maximum bandwidths, and the existence of various encoded alternatives of multimedia components, accessibility features and required digital rights management (DRM), media-component locations on the network, and other content characteristics. Using this information, the DASH client selects the appropriate encoded alternative and starts streaming the content by fetching the segments using HTTP GET requests.

After appropriate buffering to allow for network throughput variations, the client continues fetching the subsequent segments and also monitors the network bandwidth fluctuations. Depending on its measurements, the client decides how to adapt to the available bandwidth by fetching segments of different alternatives (with lower or higher bitrates) to maintain an adequate buffer.

The MPEG-DASH specification only defines the MPD and the segment formats. The delivery of the MPD and the media-encoding formats containing the segments, as well as the client behavior for fetching, adaptation heuristics, and playing content, are outside of MPEG-DASH's scope.

Multimedia Presentation Description

Dynamic HTTP streaming requires various bitrate alternatives of the multimedia content to be available at the server. In addition, the multimedia content might consist of several media components (for example, audio, video, and text), each of which might have different characteristics. In MPEG-DASH, these characteristics are described by MPD, which is an XML document.

Figure 3 demonstrates the MPD hierarchical data model. The MPD consists of one or multiple periods, where a period is a program

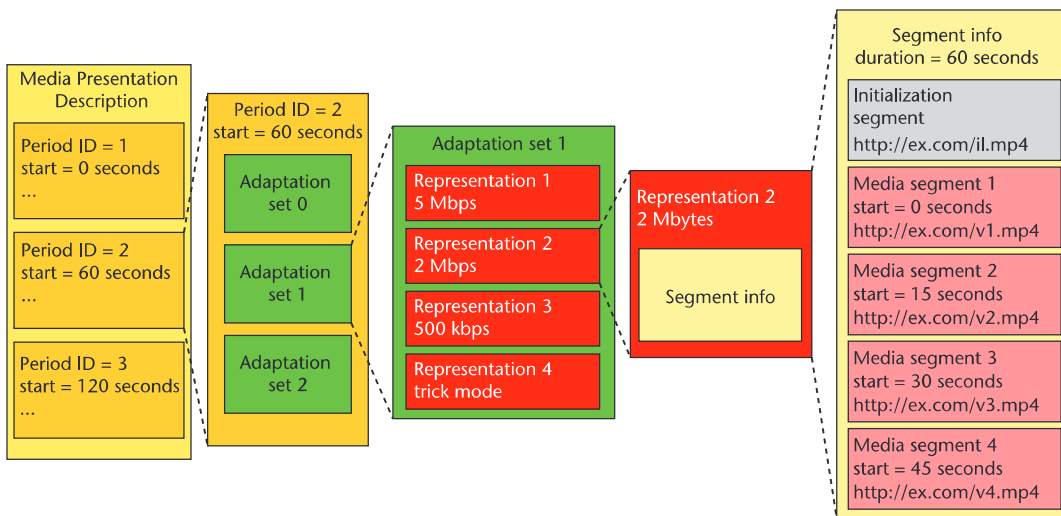


Figure 3. The Multimedia Presentation Description hierarchical data model. In this example, MPD contains three periods, period 2 contains three adaptation sets and the adaptation set 1 contains four representations, including three representations with various bit rates and one representation for trick mode. Finally, representation 2 consists of its segment info, which consequently includes its initialization segment and four media segments' information.

interval along the temporal axis. Each period has a starting time and duration and consists of one or multiple adaptation sets. An adaptation set provides the information about one or multiple media components and its various encoded alternatives. For instance, an adaptation set might contain the different bitrates of the video component of the same multimedia content. Another adaptation set might contain the different bitrates of the audio component (for example, lower-quality stereo and higher-quality surround sound) of the same multimedia content. Each adaptation set usually includes multiple representations.

A representation is an encoded alternative of the same media component, varying from other representations by bitrate, resolution, number of channels, or other characteristics. Each representation consists of one or multiple segments. Segments are the media stream chunks in temporal sequence. Each segment has a URI—that is, an addressable location on a server that can be downloaded using HTTP GET or HTTP GET with byte ranges.

To use this data model, the DASH client first parses the MPD XML document. The client then selects the set of representations it will use based on descriptive elements in the MPD, the client's capabilities, and user's choices. The client then builds a timeline and starts playing the multimedia content by requesting appropriate media segments. Each representation's description includes information about its segments, which enables requests for each segment to be formulated in terms of the HTTP URL and byte range. For live presentations, the MPD also provides segment availability

start time and end time, approximate media start time, and the fixed or variable duration of segments.

Segment format

The multimedia content can be accessed as a collection of segments. A segment is defined as the entity body of the response to the DASH client's HTTP GET or a partial HTTP GET. A media component is encoded and divided in multiple segments. The first segment might be an initialization segment containing the required information for initialization of the DASH client's media decoder. It doesn't include any actual media data.

The media stream then is divided to one or multiple consecutive media segments. Each media segment is assigned a unique URL (possibly with byte range), an index, and explicit or implicit start time and duration. Each media segment contains at least one stream access point, which is a random access or switch-to point in the media stream where decoding can start using only data from that point forward.

To enable downloading segments in multiple parts, the specification defines a method of signaling subsegments using a segment index box.⁷ This box describes subsegments and stream access points in the segment by signaling their durations and byte offsets. The DASH client can use the indexing information to request subsegments using partial HTTP GETs. The indexing information of a segment can be put in the single box at the beginning of that segment, or spread among many indexing boxes in the segment. Different methods of spreading are possible, such as hierarchical,

daisy chain, and hybrid. This technique avoids adding a large box at the beginning of the segment and therefore prevents a possible initial download delay.

MPEG-DASH defines segment-container formats for both ISO Base Media File Format⁸ and MPEG-2 Transport Streams.⁹ MPEG-DASH is media codec agnostic and supports both multiplexed and unmultiplexed encoded content.

Multiple DRM and common encryption

In MPEG-DASH, each adaptive set can use one content-protection descriptor to describe the supported DRM scheme. An adaptive set can also use multiple content-protection schemes and as long as the client recognizes at least one, it can stream and decode the content.

In conjunction with the MPEG-DASH standardization, MPEG is also developing a common encryption standard, ISO/IEC 23001-7, which defines signaling of a common encryption scheme of media content. Using this standard, the content can be encrypted once and streamed to clients, which support different DRM license systems. Each client gets the decryption keys and other required information using its particular supported DRM system, which is signaled in the MPD, and then streams the commonly encrypted content from the same server.

Additional features

The MPEG-DASH specification is a feature-rich standard. Some of the additional features include:

- *Switching and selectable streams.* The MPD provides adequate information to the client for selecting and switching between streams, for example, selecting one audio stream from different languages, selecting video between different camera angles, selecting the subtitles from provided languages, and dynamically switching between different bitrates of the same video camera.
- *Ad insertion.* Advertisements can be inserted as a period between periods or segment between segments in both on-demand and live cases.
- *Compact manifest.* The segments' address URLs can be signaled using a template scheme resulting in a compact MPD.

- *Fragmented manifest.* The MPD can be divided into multiple parts or some of its elements can be externally referenced, enabling downloading MPD in multiple steps.

- *Segments with variable durations.* The duration of segments can be varied. With live streaming, the duration of the next segment can also be signaled with the delivery of the current segment.

- *Multiple base URLs.* The same content can be available at multiple URLs—that is, at different servers or CDNs—and the client can stream from any of them to maximize the available network bandwidth.

- *Clock-drift control for live sessions.* The UTC time can be included with each segment to enable the client to control its clock drift.

- *Scalable Video Coding (SVC) and Multiview Video Coding (MVC) support.* The MPD provides adequate information regarding the decoding dependencies between representations, which can be used for streaming any multilayer coded streams such as SVC and MVC.

- *A flexible set of descriptors.* These describe content rating, components' roles, accessibility features, camera views, frame packing, and audio channels' configuration.

- *Subsetting adaptation sets into groups.* Grouping occurs according to the content author's guidance.

- *Quality metrics for reporting the session experience.* The standard has a set of well-defined quality metrics for the client to measure and report back to a reporting server.

Most of these features are provided in flexible and extensible ways enabling the possibility of deploying MPEG-DASH for unforeseeable use cases in the future.

What's next?

The specification defines five specific profiles, each addressing a different class of applications. Each profile defines a set of constraints, limiting the MPD and segment formats to a subset

of the entire specification. Therefore, a DASH client conforming to a specific profile is only required to support those required features and not the entire specification. Some profiles are specifically designed to use the legacy content and therefore provide a migration path for the existing nonstandard solutions to a standard one.

Several other standard organizations and consortia are collaborating with MPEG to reference MPEG-DASH in their own specifications. At the same time, it seems that industry is moving quickly to provide solutions based on MPEG-DASH. Some open source implementations are also on the way. It's believed that the next two years will be a crucial time for the industry—including content and service providers, platform providers, software vendors, CDN providers, and device manufacturers—to adopt this standard and build an interoperable ecosystem for multimedia streaming over the Internet.

MM

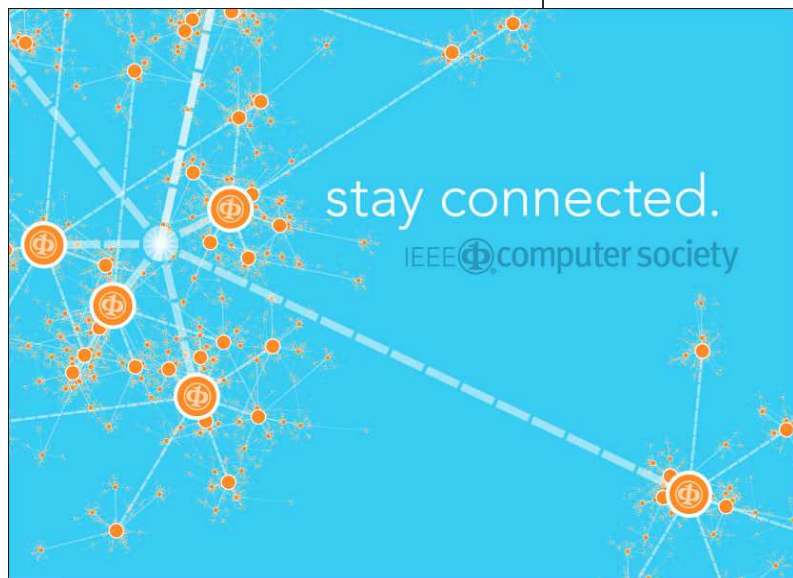
References

1. T. Sneath, *MIX09 Day 1 Keynote Pt 2: Scott Guthrie on Advancing User Experiences*, blog, 18 Mar. 2009; <http://blogs.msdn.com/b/tims/archive/2009/03/18/mix09-day-1-keynote-pt-2-scott-guthrie-on-advancing-user-experiences.aspx>.
2. Cisco Networks, *Cisco's Visual Networking Index Global IP Traffic Forecast 2010-2015*, tech. report, June 2011; http://www.cisco.com/en/US/netsol/ns827/networking_solutions_sub_solution.html#~forecast.
3. R. Pantos and E.W. May, "HTTP Live Streaming," IETF Internet draft, work in progress, Mar. 2011.
4. Microsoft, *IIS Smooth Streaming Transport Protocol*, Sept. 2009; [http://www.iis.net/community/files/media/smoothspecs/\[MS-SMTH\].pdf](http://www.iis.net/community/files/media/smoothspecs/[MS-SMTH].pdf).
5. T. Stockhammer, TS 26.247 Transparent End-to-End Packet-Switched Streaming Service (PSS); Progressive Download and Dynamic Adaptive Streaming over HTTP, 3GPP, June 2011; <http://www.3gpp.org/ftp/Specs/html-info/26247.htm>.
6. ISO/IEC FCD 23001-6, Part 6: *Dynamic Adaptive Streaming Over HTTP (DASH)*, MPEG Requirements Group, Jan. 2011; http://mpeg.chiariglione.org/working_documents/mpeg-b/dash/dash-dis.zip.
7. ISO/IEC 14496-12:2008/DAM 3, *Information Technology—Coding Of Audio-Visual Objects—Part 12: ISO Base Media File Format—Amendment 3: DASH Support and RTP Reception Hint Track Processing*, Jan. 2011.
8. *Information Technology—Coding of Audio-Visual Objects—Part 12: ISO Base Media File Format*, ISO/IEC 14496-12, 2008.
9. ITU-T Rec. H.222.0|ISO/IEC 13818-1, *Information Technology—Generic Coding of Moving Pictures and Associated Audio Information: Systems*, ITU-T/ISO/IEC, 2007; http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=44169.

Contact author Iraj Sodagar at irajs@microsoft.com.

Contact editor Anthony Vetro at avetro@merl.com.

cn Selected CS articles and columns are also available for free at <http://ComputingNow.computer.org>.



Keep up with the latest IEEE Computer Society publications and activities wherever you are.

twitter | @ComputerSociety
 | @ComputingNow
facebook | facebook.com/IEEEComputerSociety
 | facebook.com/ComputingNow
LinkedIn | IEEE Computer Society
 | Computing Now
You Tube | youtube.com/ieeecompulersociety