



- (51) International Patent Classification: *G06F 15/16* (2006.01)
- (21) International Application Number: PCT/US2012/054689
- (22) International Filing Date: 11 September 2012 (11.09.2012)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data: 13/232,909 14 September 2011 (14.09.2011) US
- (71) Applicant (for all designated States except US): **MOB-ITV, INC.** [US/US]; 6425 Christie Avenue, 5th Floor, Emeryville, CA 94608 (US).
- (72) Inventors; and
- (75) Inventors/Applicants (for US only): **NOONEY, Charles** [US/US]; 6425 Christie Avenue, 5th Floor, Emeryville, CA 94608 (US). **KARLSSON, Kent** [US/US]; 6425 Christie Avenue, 5th Floor, Emeryville, CA 94608 (US).
- (74) Agent: **KWAN, Avery Audrey**; Kwan & Olynick LLP, 1300 Clay Street, Ste. 600, Oakland, CA 94612 (US).
- (81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM,

[Continued on next page]

(54) Title: INTELLIGENT DEVICE MEDIA STREAM CACHING

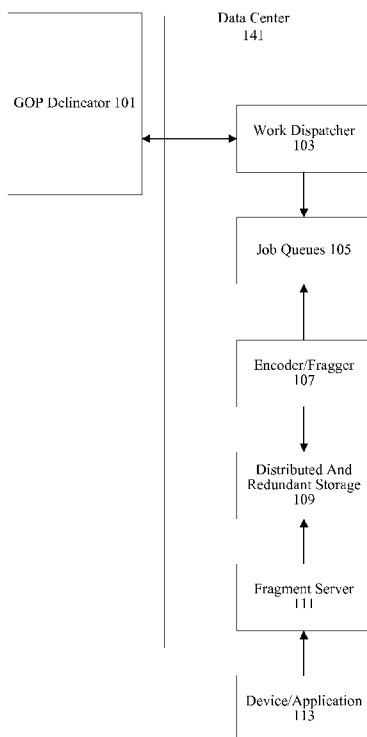


Figure 1

(57) Abstract: A media stream delivery system encodes numerous media streams into media stream fragments. The media stream delivery system may encode each media stream into a number of variants each having different resolutions, frame rates, audio quality levels, etc. Devices access the media stream fragments from a fragment server in order to reconstruct a particular media stream for playback. A device may perform caching of media stream fragments so that particular fragments need not be accessed from a fragment server. The device monitors and analyzes media streams and viewing characteristics to intelligently select fragments that will likely be needed again, such as fragments associated with repeated advertisements, introduction sequences, and end sequences.

WO 2013/039927 A1

TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG). **Published:**

— with international search report (Art. 21(3))

INTELLIGENT DEVICE MEDIA STREAM CACHING

RELATED APPLICATION DATA

- 5 [0001] This application claims priority to U.S. Patent Application No. 13/232,909, entitled INTELLIGENT DEVICE MEDIA STREAM CACHING, filed on September 14, 2011 (Attorney Docket No. MOBIP075US), the entire disclosure of which is incorporated herein by reference for all purposes.

TECHNICAL FIELD

- 10 [0002] The present disclosure relates to intelligent device media stream caching.

DESCRIPTION OF RELATED ART

- [0003] Media streams typically involve encoding or re-encoding prior to transmission to devices and users associated with the devices. In many instances, media streams are encoded into a format such as H.264 (MPEG-4 Part 15 10). H.264 is a block oriented motion compensation based codec that is widely used in Blu-ray Discs and streaming Internet sources. H.264 encoding can be resource intensive, and specialized hardware is often used to accelerate encoding particularly at high quality levels. In many implementations, live stream encoding servers are configured with application specific hardware to receive 20 one or more channels or media streams and encode the channels or media streams into particular formats. The encoding servers may have the capacity to perform real-time live encoding on up to half a dozen media streams simultaneously. Devices receive media streams and perform decoding for playback.
- 25 [0004] However, mechanisms for caching media streams for more efficient playback are limited. Consequently, the techniques and mechanisms of the present invention provide improved mechanisms for performing media stream caching.

BRIEF DESCRIPTION OF THE DRAWINGS

[0005] The disclosure may best be understood by reference to the following description taken in conjunction with the accompanying drawings, which illustrate particular embodiments.

5 [0006] Figure 1 illustrates one example of a distributed encoding system.

[0007] Figure 2 illustrates one example of a mechanism for implementing distributed encoding redundancy with live stream variant monitoring.

[0008] Figure 3 illustrates a technique for performing distributing encoding and fragmentation.

10 [0009] Figure 4 illustrates a technique for performing intelligent device media stream fragment caching.

[0010] Figure 5 illustrates one example of a fragment table.

[0011] Figure 6 illustrates a technique for performing fragment server directed device caching.

15 [0012] Figure 7 illustrates examples of files stored by a fragment writer.

[0013] Figure 8 illustrates one example of an exchange used with a fragmentation system.

[0014] Figure 9 illustrates one example of a system.

DESCRIPTION OF EXAMPLE EMBODIMENTS

[0015] Reference will now be made in detail to some specific examples of the invention including the best modes contemplated by the inventors for carrying out the invention. Examples of these specific embodiments are illustrated in the accompanying drawings. While the invention is described in conjunction with these specific embodiments, it will be understood that it is not intended to limit the invention to the described embodiments. On the contrary, it is intended to cover alternatives, modifications, and equivalents as may be included within the spirit and scope of the invention as defined by the appended claims.

[0016] For example, the techniques of the present invention will be described in the context of fragment servers. However, it should be noted that the techniques of the present invention may also apply to fragment server variations and media stream servers. In the following description, numerous specific details are set forth in order to provide a thorough understanding of the present invention. Particular example embodiments of the present invention may be implemented without some or all of these specific details. In other instances, well known process operations have not been described in detail in order not to unnecessarily obscure the present invention.

[0017] Various techniques and mechanisms of the present invention will sometimes be described in singular form for clarity. However, it should be noted that some embodiments include multiple iterations of a technique or multiple instantiations of a mechanism unless noted otherwise. For example, a system uses a processor in a variety of contexts. However, it will be appreciated that a system can use multiple processors while remaining within the scope of the present invention unless otherwise noted. Furthermore, the techniques and mechanisms of the present invention will sometimes describe a connection between two entities. It should be noted that a connection between two entities does not necessarily mean a direct, unimpeded connection, as a variety of other entities may reside between the two entities. For example, a processor may be connected to memory, but it will be appreciated that a variety of bridges and

controllers may reside between the processor and memory. Consequently, a connection does not necessarily mean a direct, unimpeded connection unless otherwise noted.

[0018] Overview

- 5 [0019] A media stream delivery system encodes numerous media streams into media stream fragments. The media stream delivery system may encode each media stream into a number of variants each having different resolutions, frame rates, audio quality levels, etc. Devices access the media stream fragments from a fragment server in order to reconstruct a particular media stream for playback.
- 10 A device may perform caching of media stream fragments so that particular fragments need not be accessed from a fragment server. The device monitors and analyzes media streams and viewing characteristics to intelligently select fragments that will likely be needed again, such as fragments associated with repeated advertisements, introduction sequences, and end sequences.

15 [0020] Example Embodiments

- [0021] A variety of mechanisms are used to deliver media streams to devices. Different devices and different networks may require different variants of a media stream. Some devices may request a higher bit rate or higher resolution stream while changes in network conditions may necessitate switching to a stream having a lower quality level. Some devices may be able to handle higher resolutions, while others may have limited processing resources or limited screen real estate. Consequently, many systems will encode numerous variants of each media stream. For example, a media provider covering 152 channels may encode 8 variants of each channel for a total of 1216 variants. In some instances,
- 20 a media provider may actually encode each channel into 8 variants of a particular
- 25 codec, such as MPEG-4 part 10 or H.264.

[0022] Conventional MPEG-4 files require that a player on a device parse the entire header before any of the data can be decoded. Parsing the entire header can take a notable amount of time, particularly on devices with limited network

and processing resources. Consequently, the techniques and mechanisms of the present invention provide a fragmented MPEG-4 framework that allows playback upon receiving a first MPEG-4 file fragment. A second MPEG-4 file fragment can be requested using information included in the first MPEG-4 file fragment. According to various embodiments, the second MPEG-4 file fragment requested may be a fragment corresponding to a higher or lower bit-rate stream than the stream associated with the first file fragment.

[0023] MPEG-4 is an extensible container format that does not have a fixed structure for describing media types. Instead, MPEG-4 has an object hierarchy that allows custom structures to be defined for each format. The format description is stored in the sample description ('stsd') box for each stream. The sample description box may include information that may not be known until all data has been encoded. For example, the sample description box may include an average bit rate that is not known prior to encoding.

[0024] According to various embodiments, MPEG-4 files are fragmented so that a live stream can be intelligently encoded in a distributed architecture on dynamically scalable hardware, recorded, and played back in a close to live manner. MPEG-4 files can be created without having to wait until all content is written to prepare the movie headers. To allow for MPEG-4 fragmentation without out of band signaling, a box structure is provided to include synchronization information, end of file information, and chapter information. According to various embodiments, synchronization information is used to synchronize audio and video when playback entails starting in the middle of a stream. End of file information signals when the current program or file is over. This may include information to continue streaming the next program or file. Chapter information may be used for video on demand content that is broken up into chapters, possibly separated by advertisement slots.

[0025] MPEG-4 fragments may be maintained on fragment servers. Devices request fragments from the fragment server in order to reconstruct particular media streams for playback. In some examples, fragments correspond to a

particular live stream variant that is reconstructed as additional fragments are made available.

[0026] According to various embodiments, devices maintain fragments in a caching layer and check the caching layer before accessing a fragment server to
5 obtain the fragment. Various network nodes between a fragment server and a device may also operate to cache particular fragments or even portions of fragments. Caching reduces latency required to obtain a particular fragment for playback and also reduces network resource consumption. Device caching schemes may include first in first out (FIFO) buffers that will maintain recent
10 fragments for a particular period of time. However, many conventional caching algorithms do not work particularly well for media stream fragments. Consequently, the techniques of the present invention provide mechanisms for devices to recognize characteristics of a media stream and viewing characteristics of a user in order to more intelligently select fragments or
15 portions of fragments for caching.

[0027] According to various embodiments, a device monitors and analyzes viewing characteristics of particular users as well as characteristics of media streams in order to identify fragments that may be requested more than once. The device may determine that particular fragments corresponding to
20 introduction and end sequences for particular programs are repeatedly accessed by a particular user on a device. Fragments corresponding to those repeated introduction and end sequences may be maintained for extended periods of time in a caching layer that can be efficiently accessed during playback. In other examples, repeated sequences such as repeated program and advertisement
25 sequences, etc., are automatically identified and maintained in a local caching layer.

[0028] According to various embodiments, a fragment server can also perform media stream and device characteristics monitoring and analysis. According to various embodiments, a fragment server has more information about what
30 fragments are more frequently accessed or what fragments will be more frequently accessed. In particular embodiments, a fragment server knows that a

particular advertisement sequence will be played a particular number of times over the next 24 hours and indicates to various devices that fragments corresponding to the advertisement sequence should be cached. In particular embodiments, flags and/or time period indicators may be included in individual
5 fragments to indicate how long fragments should be maintained at device caching layers.

[0029] Request for fragments are exposed as separate files to clients and files should play on players that handle fragmented MPEG-4. Live or near live, video on demand (VOD), and digital video record (DVR) content can all be encoded
10 on distributed and dynamically scalable encoding resources and fragmented for local caching and playback.

[0030] Figure 1 illustrates one example of a system for performing distributed encoding, fragmentation, and caching. According to various embodiments, a media stream is received from a content provider source such as a satellite. In particular embodiments, the media stream is provided in an MPEG-2 format.
15 The media stream is delineated into Groups of Pictures (GOPs) using a GOP delineator 101. The GOP is a group of pictures in coded media and typically includes key and predictive frames. A key frame may be an I-frame or intra-coded frame that represents a fixed image that is independent of other pictures.
20 According to various embodiments, each GOP begins with an I-frame. Predictive frames such as P-frames or predictive-coded frames and B-frames or bidirectionally predictive coded frames contain different information indicating distinctions from a reference frame such as a key frame or another predictive frame.

25 [0031] After the media stream is delineated into GOPs, a work dispatcher 103 is notified that a GOP is available. According to various embodiments, the work dispatcher 103 determines if it is the one assigned to work on it as well as what should be done with the GOP. According to various embodiments, the work dispatcher may determine that the GOP should be encoded into 8 different
30 variants. In particular embodiments, the work dispatcher 103 creates a description of what needs to be done, assigns a weight or priority level to the job,

and sends the job to job queues 105. According to various embodiments, job queues are first in first out (FIFO) queues that are empty most of the time. Encoders/fraggers 107 request jobs and obtain them from the job queues 105. According to various embodiments, jobs may be ordered in a job queue based on weight. In particular embodiments, encoders/fraggers 107 may select higher priority jobs first from the job queues.

[0032] In particular embodiments, different priority jobs are placed in different priority job queues. Multiple jobs may be taken from the higher priority job queues before a single job is taken from the lower priority job queues. According to various embodiments, highest priority jobs are processed before lower priority jobs. In particular embodiments, queues are assigned percentage of service values. A high priority queue may get serviced 40% of the time. A medium priority queue 30% of the time, and the remaining queues 20% and 10% of the time by the encoders/fraggers. According to various embodiments, hundreds or thousands of encoders/fraggers reside in a system. In particular embodiments, the same device performs both encoding and fragmentation, but it should be noted that separated devices can be used to perform these operations. According to various embodiments, additional encoder/fraggers can be dynamically brought online when resource usage reaches a particular threshold. Alternatively, encoder/fraggers can be taken offline when resources usage falls beneath a particular floor. According to various embodiments, encoder/fragger 107 is a virtual machine that may reside on one or more physical servers that may or may not have specialized encoding hardware. In particular embodiments, a cloud service determines how many of these virtual machines to use based on established thresholds.

[0033] According to various embodiments, a unique identifier is provided for each GOP and a log of each step is maintained. After the encoder/fragger 107 completes processing a job and outputs an encoded fragment, the encoded fragment is maintained in distributed and redundant storage 109. In one example, distributed and redundant storage 109 is a virtualized scale out network attached storage system. The distributed and redundant storage 109 allows a

system to maintain numerous fragments on any number of virtualized storage devices.

[0034] According to various embodiments, fragments on distributed and redundant storage 109 are accessible by fragment server 111. The fragment server 111 provides the caching layer with fragments for clients. The design philosophy behind the client/server API minimizes round trips and reduces complexity as much as possible when it comes to delivery of the media data to a client device. The fragment server 111 provides live streams and/or DVR configurations.

[0035] According to various embodiments, a client device uses a media component that requests fragmented MPEG-4 files, allows trick-play, and manages bandwidth adaptation. In particular embodiments, each client device receives a media stream that is behind a live stream by 12 seconds or more. There may also be server buffering. According to various embodiments, GOP delineation, encoding, fragmentation can occur within a server buffering timeframe. By having numerous encoder/fraggers, capacity can be increased or decreased by percentage points at any time. According to various embodiments, the encoding and fragmentation system at data center 141 can be accessed by devices during media stream playback. Device 113 may request fragments from fragment server 111. In particular embodiments, the device may perform local caching. It should be noted that a variety of nodes may reside between a device 113 and a fragment server 111.

[0036] Figure 2 illustrates one example of a distributed, scalable encoding system that provides for localized redundancy. According to various embodiments, a media stream is received from a content provider source such as a satellite. In particular embodiments, the media stream is provided in an MPEG-2 format. The media stream is delineated into Groups of Pictures (GOPs) using a GOP delineator 201. The GOP is a group of pictures in coded media and typically includes key and predictive frames. A key frame may be an I-frame or intra-coded frame that represents a fixed image that is independent of other pictures. According to various embodiments, each GOP begins with an I-

frame. Predictive frames such as P-frames or predictive-coded frames and B-frames or bidirectionally predictive coded frames contain different information indicating distinctions from a reference frame such as a key frame or another predictive frame. According to various embodiments, multiple GOP delineators
5 201 are active simultaneously. If a GOP delineator fails, other GOP delineators are available and all GOP delineators can send out notifications.

[0037] After the media stream is delineated into GOPs, an elastic load balancer 211 is used to distribute work to work dispatchers 221 and 225. According to various embodiments, a live stream variant encoding manager 213 monitors live
10 stream variant consumption. If particular variant are not being consumed, jobs for creating those variants are no longer performed. If particular not yet available variants are requested, then jobs creating those variants can be generated by the work dispatcher 225 at the request of the live stream variant encoding manager 213. If a work dispatcher fails right as it takes a notification,
15 another notification occurs to a different work dispatcher. Two notifications for the same GOP will end up on two different machines. At each work dispatcher 221 and 225, there may also be a proxy. According to various embodiments, the GOP delineator 201 resides on a different data center than the work dispatchers 221 and 225. Using proxies at work dispatchers 221 and 225 allows for a single
20 transfer of a media stream GOP between data centers.

[0038] According to various embodiments, the work dispatchers 221 and 225 determine characteristics of a particular job and what should be done with the GOP. According to various embodiments, the work dispatchers 221 and 225 may determine that the GOP should be encoded into 8 different variants. In
25 particular embodiments, the work dispatchers 221 and 225 create descriptions of what needs to be done and send jobs to job queues 223. According to various embodiments, job queues 223 include an active job queue and a standby job queue. According to various embodiments, job queues are first in first out (FIFO) queues that are empty most of the time. Timeouts may be associated
30 with each job in the queue. Encoders/fraggers 231, 233, and 235 request jobs and obtain them from the job queues 223. In particular embodiments, encoders/fraggers 231, 233, and 235 are identical and can be dynamically

activated or deactivated. According to various embodiments, hundreds or thousands of encoders/fraggers reside in a system.

[0039] In particular embodiments, the same device performs both encoding and fragmentation, but it should be noted that separated devices can be used to perform these operations. According to various embodiments, additional encoder/fraggers can be dynamically brought online when resource usage reaches a particular threshold. Alternatively, encoder/fraggers can be taken offline when resources usage falls beneath a particular floor. According to various embodiments, encoder/fragger 231, 233, and 235 is a virtual machine that may reside on one or more physical servers that may or may not have specialized encoding hardware. In particular embodiments, a cloud service determines how many of these virtual machines to use based on established thresholds.

[0040] According to various embodiments, encoders/fraggers 231, 233, and 235 are stateless. According to various embodiments, a unique identifier is provided for each GOP and a log of each step is maintained. If a particular encoder/fragger fails at any point in the process, another encoder/fragger can perform encoding and fragmentation operations. After the encoders/fraggers 231, 233, and 235 complete the jobs and generate encoded fragments, the encoded fragments are maintained in distributed and redundant storage 241. In one example, distributed and redundant storage 241 is a virtualized scale out network attached storage system. The distributed and redundant storage 241 includes nodes 243 and 245, allowing a system to maintain numerous fragments on any number of virtualized storage devices.

[0041] According to various embodiments, fragments on distributed and redundant storage 241 are accessible by fragment servers 251, 253, and 255. The fragment servers 251, 253, and 255 provide a caching layer with fragments for clients. The design philosophy behind the client/server API minimizes round trips and reduces complexity as much as possible when it comes to delivery of the media data to a client device. The fragment servers 251, 253, and 255 provide live streams and/or DVR configurations. According to various

embodiments, fragment servers also operate without state. In particular
embodiments, fragments servers operate using HTTP get requests. According to
various embodiments, each process allows a flow to continue without having a
centralized control point. An elastic load balancer 261 distributes fragment
5 requests from a cloud front 271 provided to devices 281, 283, and 285.
According to various embodiments, devices 281, 283, and 285 monitor and
analyze media streams to determine what fragments should be cached. In some
examples, devices 281, 283, and 285 cache any fragment that has been
determined to be redundant to any fragment previously requested. Fragments
10 can be compared using identifiers, hashes, etc.

[0042] It is determined that if the same fragment has been obtained at least
twice, it is likely that the fragment will need to be accessed again. In particular
embodiments, the device identifies sequences corresponding to advertisements,
program introductions, program endings, etc., and caches fragments associated
15 with those identified sequences. Particular advertisements, introductions, and
endings, etc., may be more likely repeated. If a viewer is identified as frequently
watching a particular program, fragments associated with introduction
sequences, ending sequences, and advertisements for that program may be
cached.

[0043] According to various embodiments, a client device uses a media
20 component that requests fragmented MPEG-4 files, allows trick-play, and
manages bandwidth adaptation. In particular embodiments, each client device
receives a media stream that is behind a live stream by 12 seconds or more.
Caching can reduce delay and/or network resource usage. According to various
25 embodiments, GOP delineation, encoding, fragmentation can occur within a
server buffering timeframe. By having numerous encoder/fraggers, capacity can
be increased or decreased by percentage points at any time. According to
various embodiments, a system provides not only localized redundancy but
geographic redundancy as well. A complete system including load balancers,
30 work dispatchers, encoders/fraggers, storage, fragment servers, etc., may be
replicated at a geographically separate data center.

[0044] Figure 3 illustrates one example of a technique for performing distributed encoding. At 301, a live media stream is received. According to various embodiments, the media stream is a live MPEG-2 media stream received from a satellite receiver. In particular embodiments, a live media stream refers to a media program received at a particular time that is designated for distribution at that particular time. For example, a program may be configured to run at 8pm PST, and the live media stream is received at the satellite receiver at 8pm PST. At 303, the media stream is delineated into GOPs. In particular embodiments, key frames are identified and used to begin groups of pictures. The GOPs may be optionally encrypted at this point before transmission at 305.

[0045] At 307, a work dispatcher determines that a GOP is available. According to various embodiments, the GOP delineator and the work dispatcher reside in different data centers and it is desirable to have a single GOP transfer instead of redundant GOP transfers. At 309, the work dispatcher creates descriptions of multiple jobs corresponding to different encoding quality levels and places the jobs in work queues. According to various embodiments, the jobs are assigned different weights and/or priority levels. An encoder/fragger at 311 pulls the job from the job queue and encodes the GOP into a format such as H.264. Encoded fragments are maintained on distributed storage at 313. A user can then continuously access the encoded fragments at a particular quality level through a fragment server at 315 that allows real-time or near real-time Internet Protocol (IP) distribution of the media stream fragments.

[0046] Figure 4 illustrates a particular example of a technique for performing device caching. At 401, a device monitors and analyzes user viewing characteristics. The device may determine at 403 that a user watches particular programs, particular channels, particular types of programs, or watches at particular times. Fragments that frequently occur during those particular programs, particular channels, particular types of programs, or particular time periods may be identified for caching at 405. At 407, characteristics of viewed programs are monitored and analyzed. Repeated sequences such as introductions, advertisements, endings, etc., are identified as candidates for local caching. Other sequences that are repeated may also be identified for caching at

409. In particular embodiments, fragments have associated identifiers, checksums, or hash codes that allow quick comparison to determine fragment identity. In one example, a fragment table is maintained showing fragment identifiers, access frequencies, and last accessed timestamps. When a fragment
5 has been accessed more than a particular number of times at 411, that fragment may be selected for caching at 413.

[0047] Figure 5 illustrates one example of a fragment table that may be used by a device or a caching layer between the device and a fragment server. According to various embodiments, a device may maintain identifiers associated with
10 particular media fragments to determine what fragments to cache. In some examples, a fragment table 501 includes media fragment identifiers 501, 511, and 521. In particular embodiments, media fragment identifiers may be hashes or checksums of media fragments. Media fragment identifiers may be compared with a newly received fragment to determine whether the device has previously
15 seen the newly received fragment. The fragment table 501 may also include media fragment frequency indicators 503, 513, and 523 which may be values indicating how often or how many times a particular fragment has been seen. In some examples, the media fragment frequency indicator may show that a particular fragment is required once every 24 hours or once every 7 days. In
20 other examples, the media fragment frequency indicator may merely indicate that a particular fragment has been seen once before.

[0048] A media fragment last access time 505, 515, or 525 may indicate how recent activity for a particular fragment has been. In some examples, a fragment may have been regularly occurring at a distant point in the past, but may no
25 longer need caching as it has not been used for an extended period of time. In other examples, a fragment may be relatively infrequent but may be recently active, such as in the instance where the fragment occurs in a recent frequently repeated commercial. It should be noted that other field may also be included, such as fragment availability, fragment size, associated program, etc.

30 [0049] Figure 6 illustrates one example of a technique for fragment server directed device fragment caching. At 601, a fragment server monitors and

analyzes media streams to identify fragments appropriate for device caching. According to various embodiments, the fragment server knows that a particular sequence such as a promotion sequence will be played a particular number of times during the next 72 hours at 603 and identifies fragments associated with that sequence for caching at 605. In particular embodiments, the fragment server also identifies fragments most frequently repeatedly accessed by a particular device, a subgroup of devices, or all devices that can connect to the fragment server at 607. These frequently accessed fragments are selected for caching at 609. In some examples, content providers may also indicate to a fragment server what fragments to cache at 611.

[0050] At 613, the fragment server includes metadata such as a flag and/or a time period during which a particular fragment should be cached at a device. The metadata may be included in the fragment itself or may be included out of band.

[0051] Figure 7 illustrates examples of files stored by the fragment writer associated with a fragment server. According to various embodiments, the fragment writer is a component in the overall fragmenter. It is a binary that uses command line arguments to record a particular program based on either NTP time from the encoded stream or wallclock time. In particular embodiments, this is configurable as part of the arguments and depends on the input stream. When the fragment writer completes recording a program it exits. For live streams, programs are artificially created to be short time intervals e.g. 5-15 minutes in length.

[0052] According to various embodiments, the fragment writer command line arguments are the SDP file of the channel to record, the start time, end time, name of the current and next output files. The fragment writer listens to RTP traffic from the live video encoders and rewrites the media data to disk as fragmented MPEG-4. According to various embodiments, media data is written as fragmented MPEG-4 as defined in MPEG-4 part 12 (ISO/IEC 14496-12). Each broadcast show is written to disk as a separate file indicated by the show ID (derived from EPG). Clients include the show ID as part of the channel name

when requesting to view a prerecorded show. The fragment writer consumes each of the different encodings and stores them as a different MPEG-4 fragment.

[0053] In particular embodiments, the fragment writer writes the RTP data for a particular encoding and the show ID field to a single file. Inside that file, there is metadata information that describes the entire file (MOOV blocks). Atoms are stored as groups of MOOF/MDAT pairs to allow a show to be saved as a single file. At the end of the file there is random access information that can be used to enable a client to perform bandwidth adaptation and trick play functionality.

[0054] According to various embodiments, the fragment writer includes an option which encrypts fragments to ensure stream security during the recording process. The fragment writer will request an encoding key from the license manager. The keys used are similar to that done for DRM. The encoding format is slightly different where MOOF is encoded. The encryption occurs once so that it does not create prohibitive costs during delivery to clients.

[0055] The fragment server responds to HTTP requests for content. According to various embodiments, it provides APIs that can be used by clients to get necessary headers required to decode the video, seek to any desired time frame within the fragment and APIs to watch channels live. Effectively, live channels are served from the most recently written fragments for the show on that channel. The fragment server returns the media header (necessary for initializing decoders), particular fragments, and the random access block to clients. According to various embodiments, the APIs supported allow for optimization where the metadata header information is returned to the client along with the first fragment. The fragment writer creates a series of fragments within the file. When a client requests a stream, it makes requests for each of these fragments and the fragment server reads the portion of the file pertaining to that fragment and returns it to the client.

[0056] According to various embodiments, the fragment server uses a REST API that is cache friendly so that most requests made to the fragment server can be cached. The fragment server uses cache control headers and ETag headers to

provide the proper hints to caches. This API also provides the ability to understand where a particular user stopped playing and to start play from that point (providing the capability for pause on one device and resume on another).

[0057] In particular embodiments, client requests for fragments follow the following format:

5 http://{HOSTNAME}/frag/{CHANNEL}/{BITRATE}/{ID}/{COMMAND}/{ARG} e.g.

http://frag.hosttv.com/frag/1/H8QVGAH264/1270059632.mp4/fragment/42.

10 According to various embodiments, the channel name will be the same as the backend-channel name that is used as the channel portion of the SDP file. VoD uses a channel name of "vod". The BITRATE should follow the BITRATE/RESOLUTION identifier scheme used for RTP streams. The ID is dynamically assigned. For live streams, this may be the UNIX timestamp; for DVR this will be a unique ID for the show; for VoD this will be the asset ID.

15 The ID is optional and not included in LIVE command requests. The command and argument are used to indicate the exact command desired and any arguments. For example, to request chunk 42 this portion would be "fragment/42".

[0058] The URL format makes the requests content delivery network (CDN) friendly because the fragments will never change after this point so two separate clients watching the same stream can be serviced using a cache. In particular, the headend architecture leverages this to avoid too many dynamic requests arriving at the Fragment Server by using an HTTP proxy at the head end to cache requests.

25 [0059] According to various embodiments, the fragment controller is a daemon that runs on the fragmenter and manages the fragment writer processes. We propose that it uses a configured filter that is executed by the Fragment Controller to generate the list of broadcasts to be recorded. This filter integrates with external components such as a guide server to determine which shows to record and the broadcast ID to use.

30

[0060] According to various embodiments, the client includes an application logic component and a media rendering component. The application logic component presents the UI for the user and also communicates to the front-end server to get shows that are available for the user and to authenticate. As part of this process, the server returns URLs to media assets that are passed to the media rendering component.

[0061] In particular embodiments, the client relies on the fact that each fragment in a fragmented MPEG-4 file has a sequence number. Using this knowledge and a well defined URL structure for communicating with the server, the client requests fragments individually as if it was reading separate files from the server simply by requesting URLs for files associated with increasing sequence numbers. In some embodiments, the client can request files corresponding to higher or lower bit rate streams depending on device and network resources.

[0062] Since each file contains the information needed to create the URL for the next file, no special playlist files are needed, and all actions (startup, channel change, seeking) can be performed with a single HTTP request. After each fragment is downloaded the client assesses among other things the size of the fragment and the time needed to download it in order to determine if downshifting is needed, or if there is enough bandwidth available to request a higher bitrate.

[0063] Because each request to the server looks like a request to a separate file, the response to requests can be cached in any HTTP Proxy, or be distributed over any HTTP based CDN.

[0064] Figure 8 illustrates an interaction for a client receiving a live stream. The client starts playback when fragment plays out from the server. The client uses the fragment number so that it can request the appropriate subsequence file fragment. An application such as a player application 807 sends a request to mediakit 805. The request may include a base address and bit rate. The mediakit 805 sends an HTTP get request to caching layer 803. According to various embodiments, the live response is not in cache, and the caching layer

803 forward the HTTP get request to a fragment server 801. The fragment server 801 performs processing and sends the appropriate fragment to the caching layer 803 which forwards to the data to mediakit 805.

[0065] The fragment may be cached for a short period of time at caching layer
5 803. The mediakit 805 identifies the fragment number and determines whether resources are sufficient to play the fragment. In some examples, resources such as processing or bandwidth resources are insufficient. The fragment may not have been received quickly enough, or the device may be having trouble decoding the fragment with sufficient speed. Consequently, the mediakit 805
10 may request a next fragment having a different data rate. In some instances, the mediakit 805 may request a next fragment having a higher data rate. According to various embodiments, the fragment server 801 maintains fragments for different quality of service streams with timing synchronization information to allow for timing accurate playback.

[0066] The mediakit 805 requests a next fragment using information from the
15 received fragment. According to various embodiments, the next fragment for the media stream may be maintained on a different server, may have a different bit rate, or may require different authorization. Caching layer 803 determines that the next fragment is not in cache and forwards the request to fragment server
20 801. The fragment server 801 sends the fragment to caching layer 803 and the fragment is cached for a short period of time. The fragment is then sent to mediakit 805.

[0067] Figure 9 illustrates one example of a computer system. According to
25 particular embodiments, a system 900 suitable for implementing particular embodiments of the present invention includes a processor 901, a memory 903, an interface 911, and a bus 915 (e.g., a PCI bus or other interconnection fabric) and operates as a streaming server. When acting under the control of appropriate software or firmware, the processor 901 is responsible for modifying and transmitting live media data to a client. Various specially configured devices can
30 also be used in place of a processor 901 or in addition to processor 901. The

interface 911 is typically configured to send and receive data packets or data segments over a network.

[0068] Particular examples of interfaces supports include Ethernet interfaces, frame relay interfaces, cable interfaces, DSL interfaces, token ring interfaces, and the like. In addition, various very high-speed interfaces may be provided such as fast Ethernet interfaces, Gigabit Ethernet interfaces, ATM interfaces, HSSI interfaces, POS interfaces, FDDI interfaces and the like. Generally, these interfaces may include ports appropriate for communication with the appropriate media. In some cases, they may also include an independent processor and, in some instances, volatile RAM. The independent processors may control such communications intensive tasks as packet switching, media control and management.

[0069] According to various embodiments, the system 900 is a fragment server that also includes a transceiver, streaming buffers, and a program guide database. The fragment server may also be associated with subscription management, logging and report generation, and monitoring capabilities. In particular embodiments, functionality for allowing operation with mobile devices such as cellular phones operating in a particular cellular network and providing subscription management. According to various embodiments, an authentication module verifies the identity of devices including mobile devices. A logging and report generation module tracks mobile device requests and associated responses. A monitor system allows an administrator to view usage patterns and system availability. According to various embodiments, the fragment server handles requests and responses for media content related transactions while a separate streaming server provides the actual media streams.

[0070] Although a particular fragment server is described, it should be recognized that a variety of alternative configurations are possible. For example, some modules such as a report and logging module and a monitor may not be needed on every server. Alternatively, the modules may be implemented on another device connected to the server. In another example, the server may not

include an interface to an abstract buy engine and may in fact include the abstract buy engine itself. A variety of configurations are possible.

[0071] In the foregoing specification, the invention has been described with reference to specific embodiments. However, one of ordinary skill in the art
5 appreciates that various modifications and changes can be made without departing from the scope of the invention as set forth in the claims below. Accordingly, the specification and figures are to be regarded in an illustrative rather than a restrictive sense, and all such modifications are intended to be included within the scope of invention.

10

CLAIMS

1. A method, comprising:
 - requesting a first media stream from a fragment server, the first media stream including a plurality of encoded fragments;
 - 5 receiving a first encoded fragment and generating an identifier corresponding to the first encoded fragment to access a fragment table;
 - determining whether to cache the first encoded fragment using the fragment table at a device;
 - requesting a second encoded fragment using metadata from the first
 - 10 encoded fragment;
 - constructing the first media stream using the first encoded fragment and the second encoded fragment.
2. The method of claim 1, wherein the device monitors and analyzes a plurality of media streams to identify sequences of media likely to be repeated.
- 15 3. The method of claim 1, wherein frames associated with the sequences of media likely to be repeated are identified for caching at the device.
4. The method of claim 3, wherein the sequences of media include advertisement sequences.
5. The method of claim 3, wherein the sequences of media include
- 20 program introduction sequences.
6. The method of claim 3, wherein the sequences of media include program end sequences.
7. The method of claim 1, wherein the fragment table includes a plurality of fragment identifiers associated with the number of times each fragment has
- 25 been requested.
8. The method of claim 1, wherein the first encoded fragment is an H.264 encoded fragment processed using a plurality of distributed and dynamically scalable encoder resources.
9. The method of claim 8, wherein the plurality of distributed and
- 30 dynamically scalable encoder resources are a plurality of virtual machines configured to perform encoding and fragmentation.

10. The method of claim 8, wherein the H.264 encoded fragment is encoded from an MPEG-2 group of pictures (GOP).

11. A device, comprising:

an interface operable to request a first media stream from a fragment server, the first media stream including a plurality of encoded fragment, wherein the interface is operable to receive a first encoded fragment;

a processor operable to generate an identifier corresponding to the first encoded fragment to access a fragment table and determine whether to cache the first encoded fragment using the fragment table at a device;

10 wherein the device requests a second encoded fragment using metadata from the first encoded fragment and constructs the first media stream using the first encoded fragment and the second encoded fragment.

12. The device of claim 11, wherein the device monitors and analyzes a plurality of media streams to identify sequences of media likely to be repeated.

15 13. The device of claim 11, wherein frames associated with the sequences of media likely to be repeated are identified for caching at the device.

14. The device of claim 13, wherein the sequences of media include advertisement sequences.

20 15. The device of claim 13, wherein the sequences of media include program introduction sequences.

16. The device of claim 13, wherein the sequences of media include program end sequences.

25 17. The device of claim 11, wherein the fragment table includes a plurality of fragment identifiers associated with the number of times each fragment has been requested.

18. The device of claim 11, wherein the first encoded fragment is an H.264 encoded fragment processed using a plurality of distributed and dynamically scalable encoder resources.

30 19. The device of claim 18, wherein the plurality of distributed and dynamically scalable encoder resources are a plurality of virtual machines configured to perform encoding and fragmentation.

20. An apparatus, comprising:

means for requesting a first media stream from a fragment server, the first media stream including a plurality of encoded fragments;

5 means for receiving a first encoded fragment and generating an identifier corresponding to the first encoded fragment to access a fragment table;

means for determining whether to cache the first encoded fragment using the fragment table at a device;

means for requesting a second encoded fragment using metadata from the first encoded fragment;

10 means for constructing the first media stream using the first encoded fragment and the second encoded fragment.

15

20

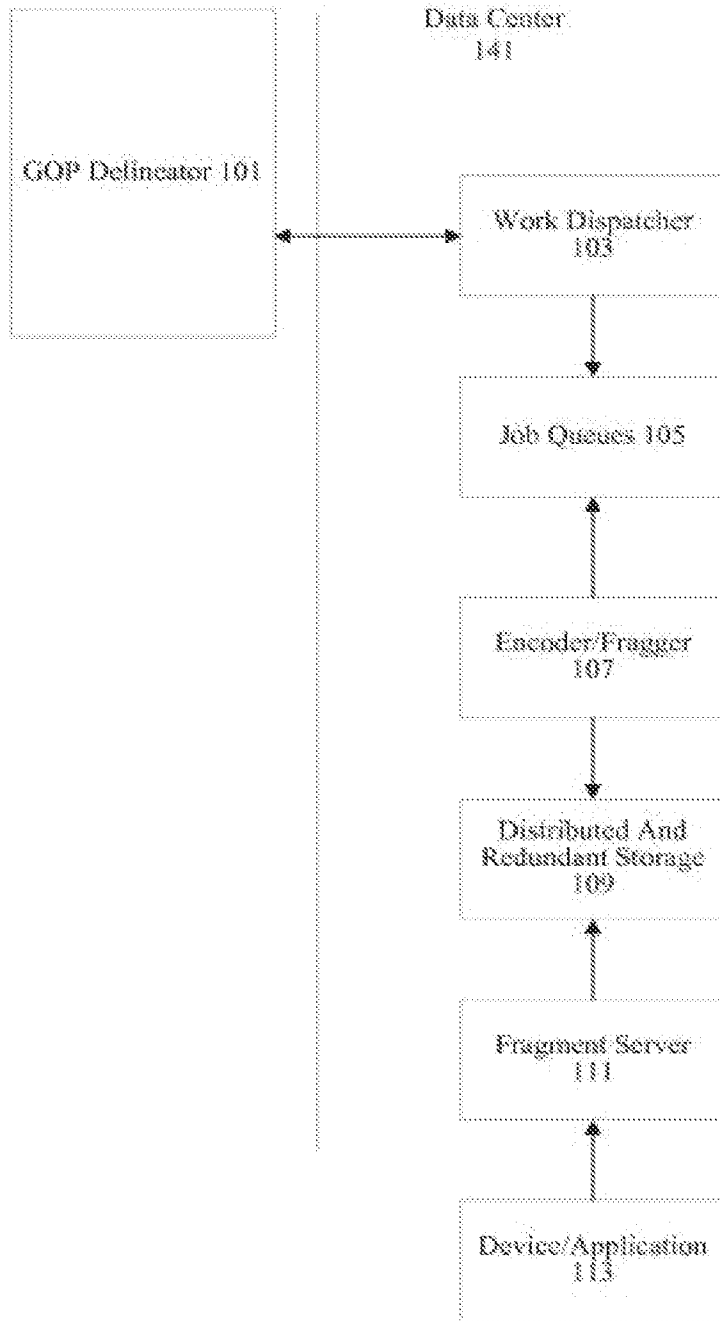


Figure 1

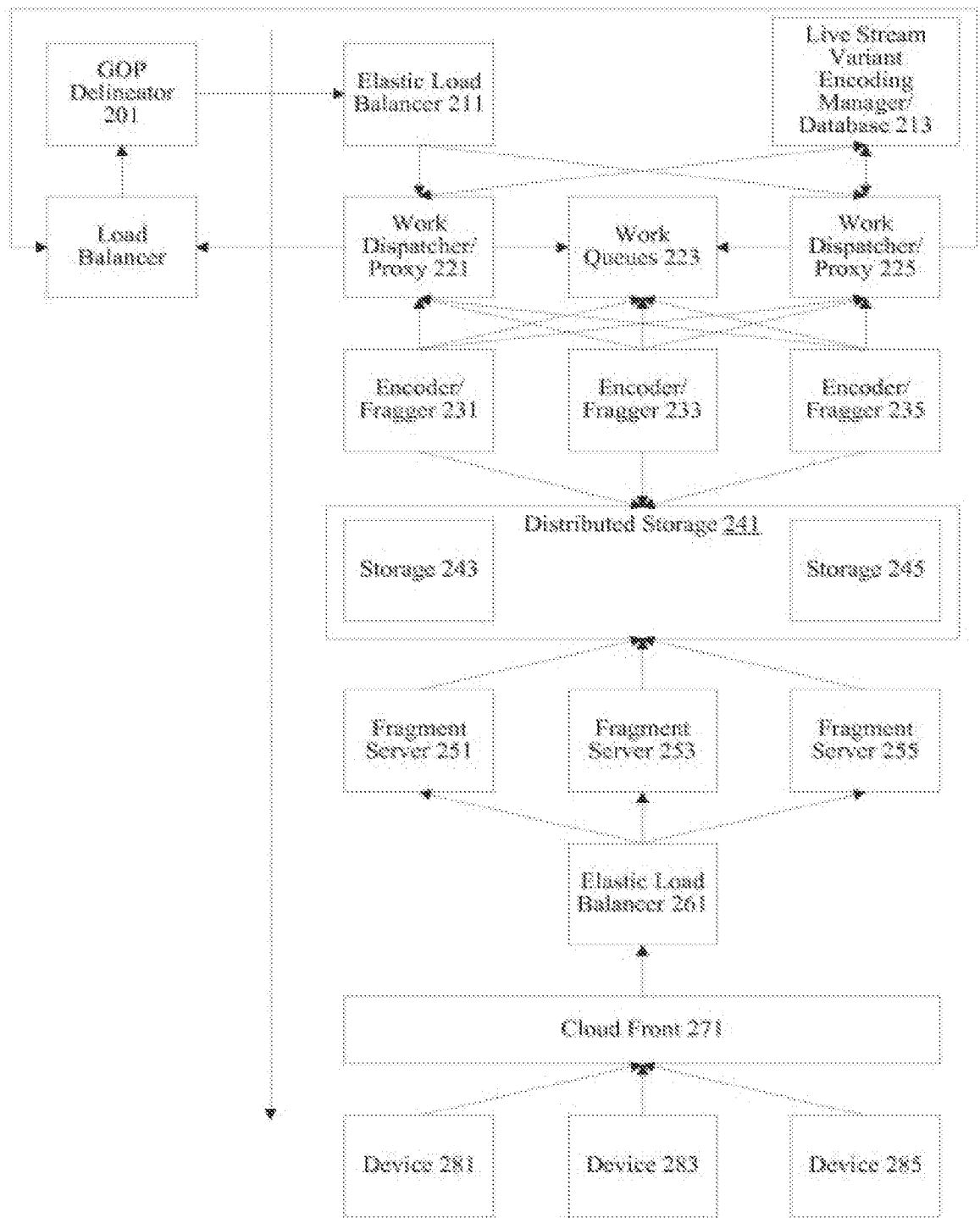


Figure 2

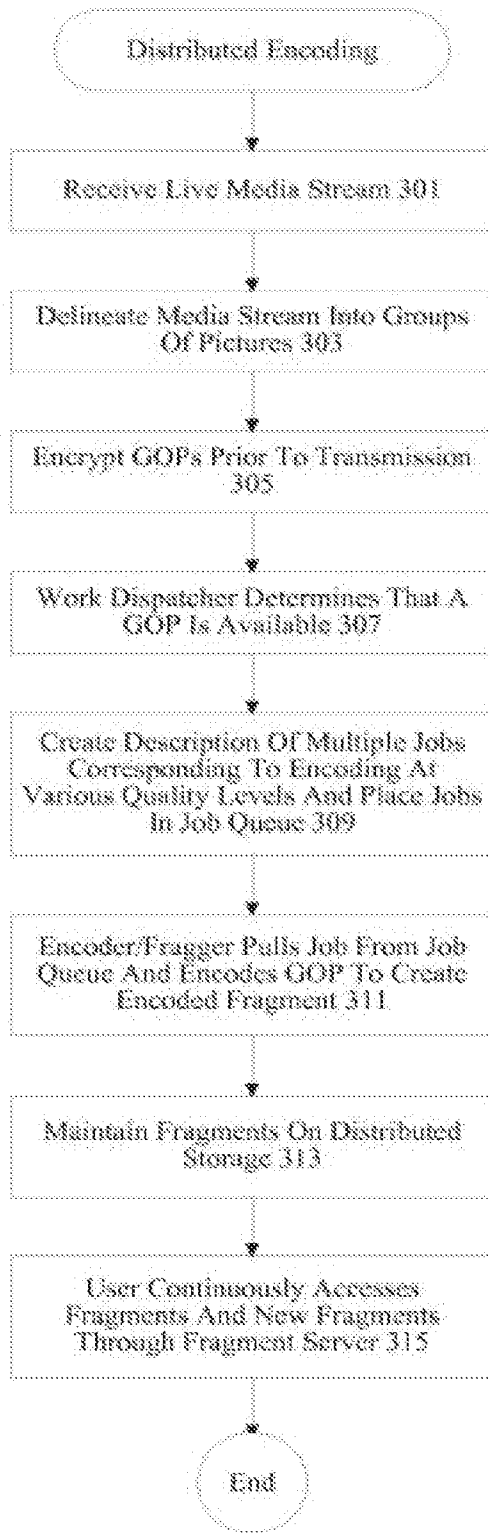


Figure 3

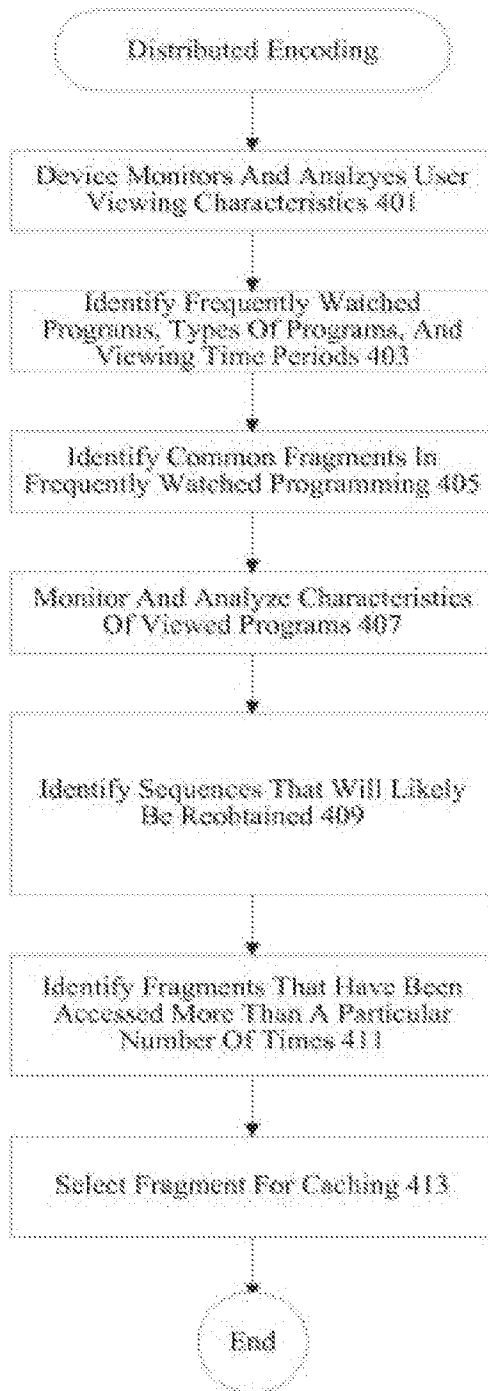


Figure 4

Fragment Table 501		
Media Fragment Identifier 501	Media Fragment Frequency Indicator 503	Media Fragment Last Access Time 505
Media Fragment Identifier 511	Media Fragment Frequency Indicator 513	Media Fragment Last Access Time 515
Media Fragment Identifier 521	Media Fragment Frequency Indicator 523	Media Fragment Last Access Time 525

Figure 5

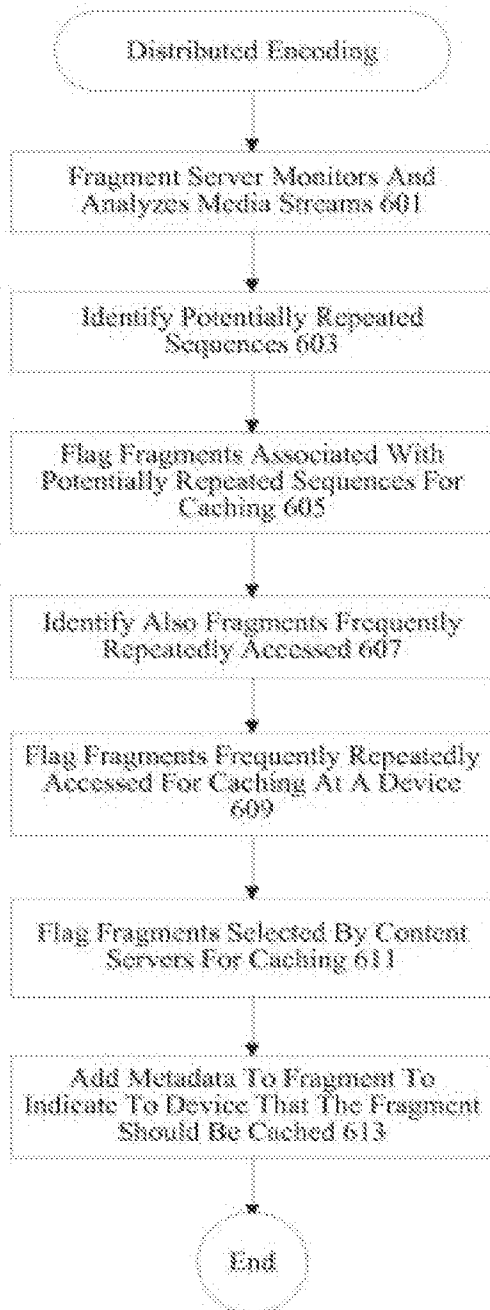


Figure 6

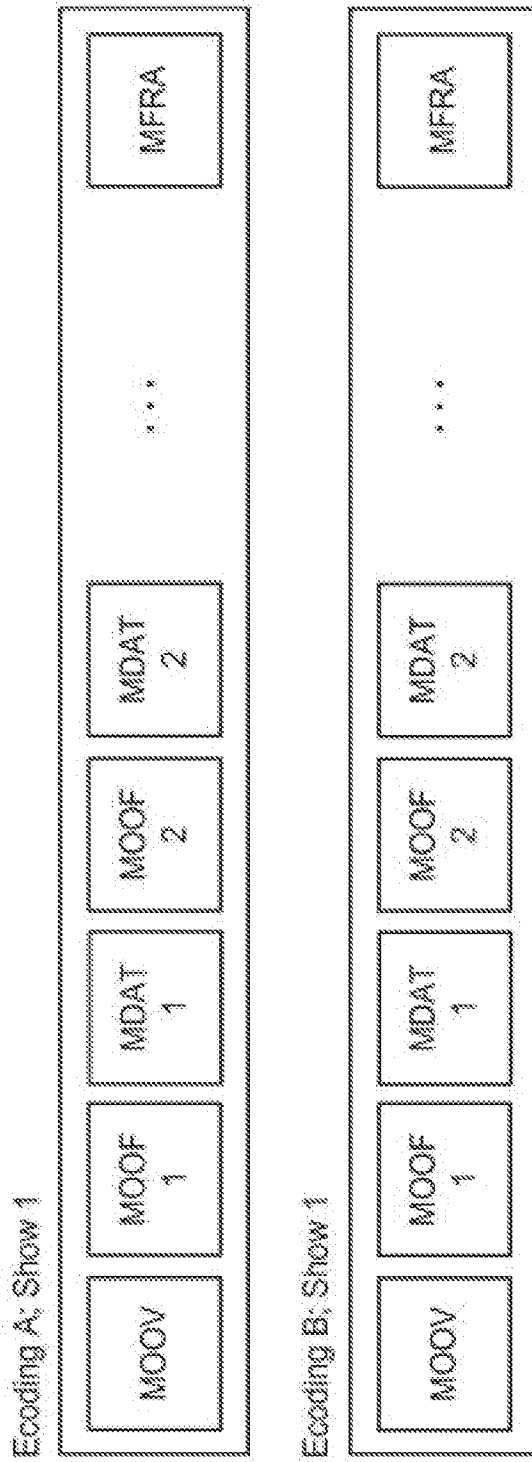


Figure 7

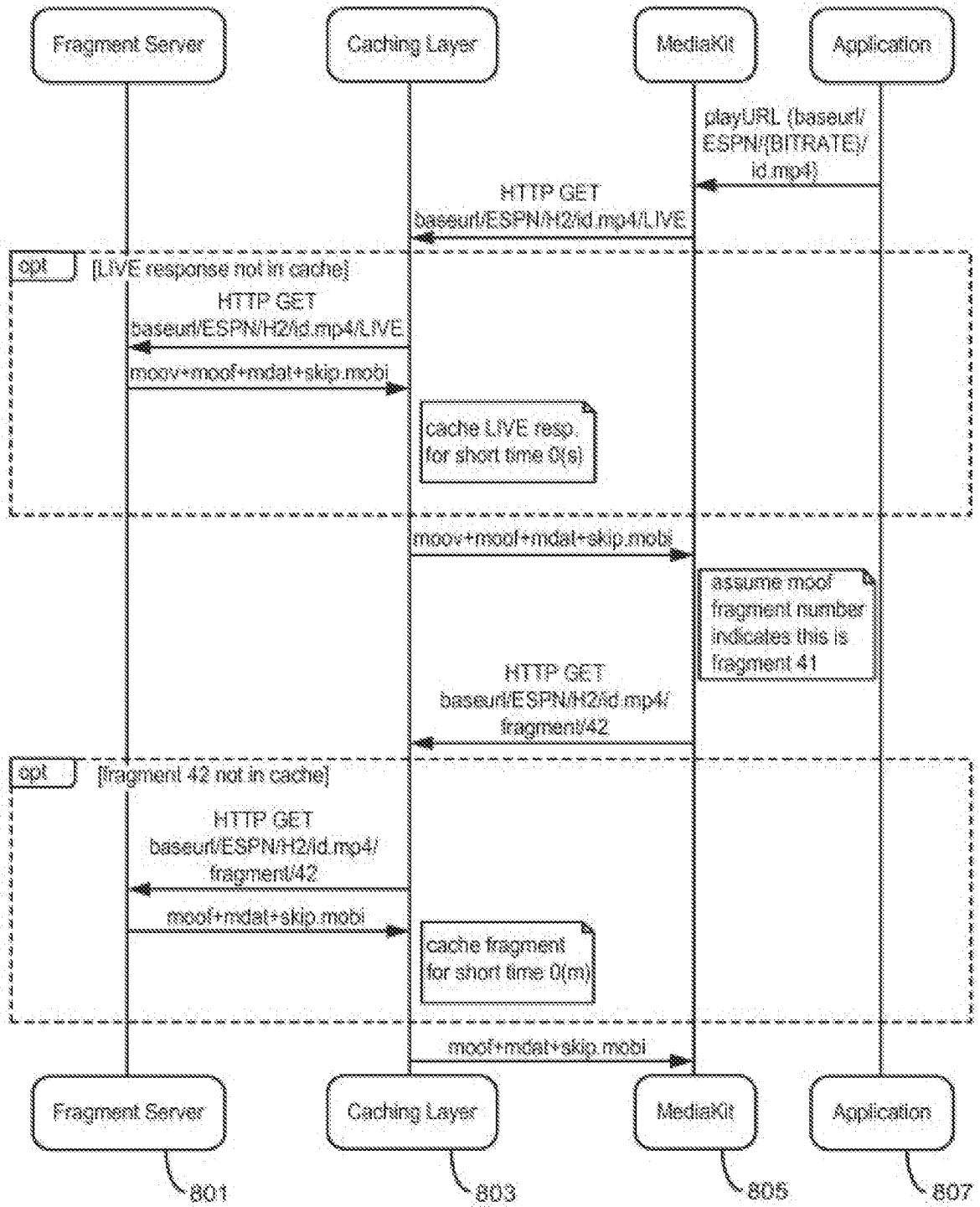


Figure 8

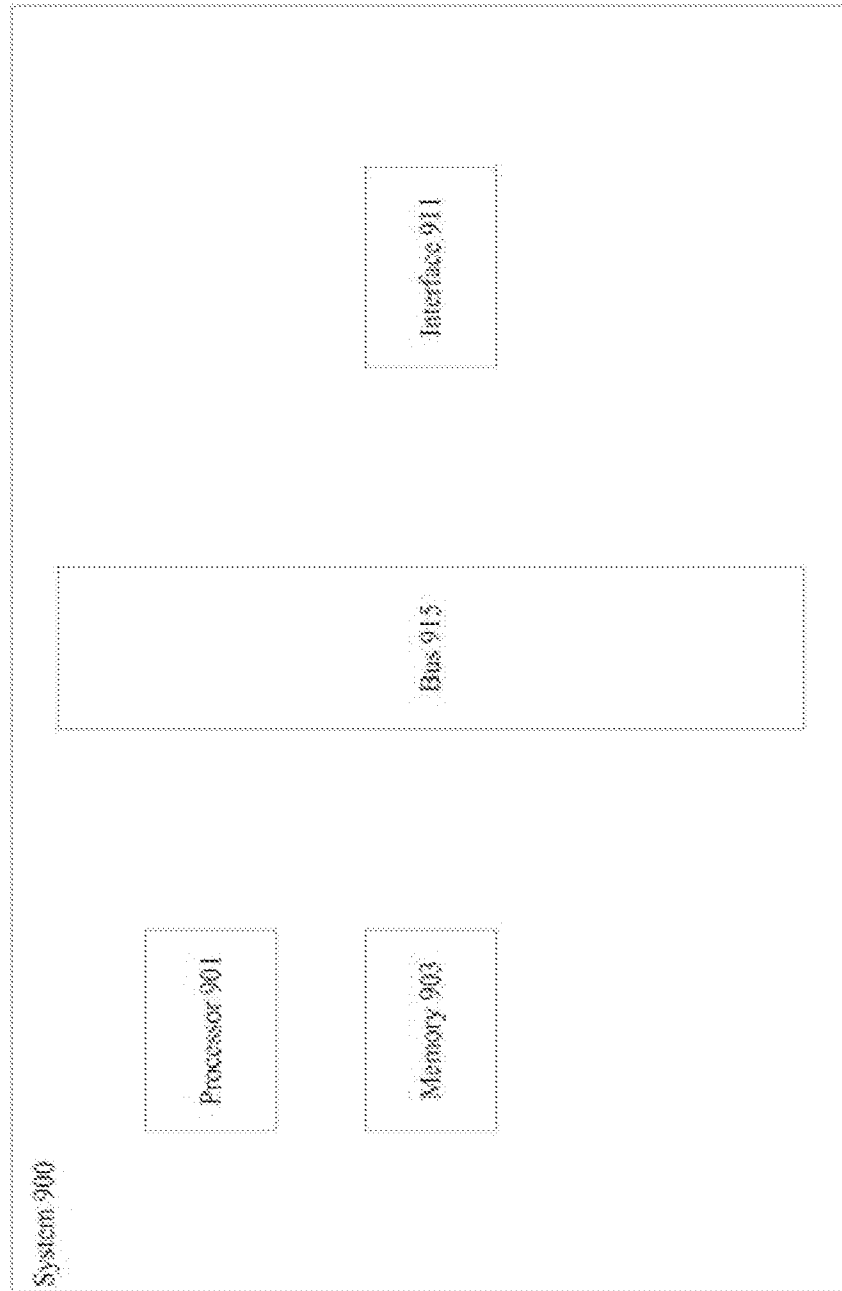


Figure 9

INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 12/54689

A. CLASSIFICATION OF SUBJECT MATTER
 IPC(8) - G06F 15/16 (2012.01)
 USPC - 709/231
 According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED
 Minimum documentation searched (classification system followed by classification symbols)
 USPC: 709/231; IPC(8): G06F 15/16 (2012.01)

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
 USPC: 709/219; 709/201; IPC(8): G06F 15/16 (2012.01)

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
 DialogWeb (347,348,349,351,371,652,654); Google Scholar; Google Web; Google Patents; PatBase.
 Search Terms: STREAM, MEDIA, FRAGMENT, SERVER, ENCODE, ENCRYPT, IDENTIFY, TABLE, SEQUENCE, CACHE, ADVERTISE, COMMERCIAL, PROGRAM, SHOW, INTRO, BEGIN, START, END, FINISH, H 264, SCALABLE

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X -- Y	US 2011/0080940 A1 (Bocharov et al.) 07 April 2011 (07.04.2011), entire document, especially para [0003], [0015], [0022], [0025], [0031], [0061], [0069], [0083]	1-4, 7, 8, 10-14, 17, 18, 20 ----- 5, 6, 9, 15, 16, 19
Y	US 2010/0122030 A1 (Peters et al.) 13 May 2010 (13.05.2010), entire document, especially para [0007], [0083]-[0084]	5, 6, 15, 16
Y	US 6,704,925 B1 (Bugnion) 09 March 2004 (09.03.2004), entire document, especially col. 4, ln 13-18; col. 6, ln 11-25	9, 19
A	US 7,925,774 B2 (Zhang et al.) 12 April 2011 (12.04.2011), entire document	1-20
A	US 5,148,209 A (Subbarao) 15 September 1992 (15.09.1992), entire document	1-20

Further documents are listed in the continuation of Box C.

* Special categories of cited documents:	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"A" document defining the general state of the art which is not considered to be of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"E" earlier application or patent but published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&" document member of the same patent family
"O" document referring to an oral disclosure, use, exhibition or other means	
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 19 October 2012 (19.10.2012)	Date of mailing of the international search report 16 NOV 2012
---	--

Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US, Commissioner for Patents P.O. Box 1450, Alexandria, Virginia 22313-1450 Facsimile No. 571-273-3201	Authorized officer: Lee W. Young PCT Helpdesk: 571-272-4300 PCT OSP: 571-272-7774
---	--