



US 20150280929A1

(19) **United States**(12) **Patent Application Publication**  
**Balay et al.**(10) **Pub. No.: US 2015/0280929 A1**(43) **Pub. Date: Oct. 1, 2015**(54) **SCALABLE IP-SERVICES ENABLED  
MULTICAST FORWARDING WITH  
EFFICIENT RESOURCE UTILIZATION****Publication Classification**

(51) **Int. Cl.**  
*H04L 12/18* (2006.01)  
*H04L 12/773* (2006.01)  
*H04L 12/713* (2006.01)

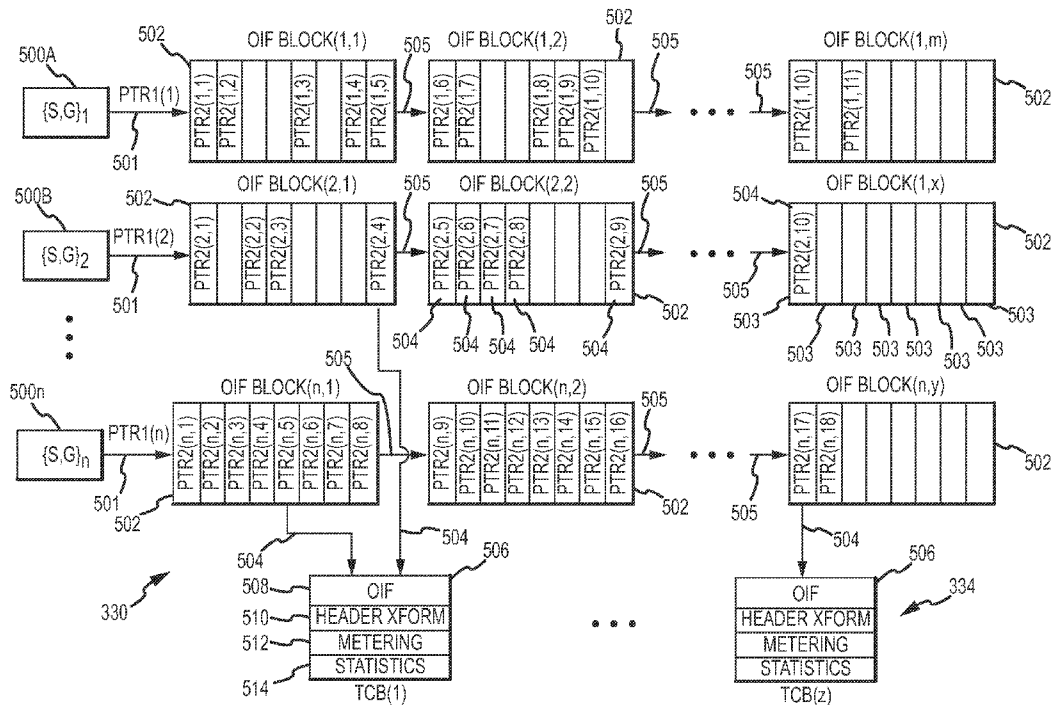
(52) **U.S. Cl.**  
CPC ..... *H04L 12/18* (2013.01); *H04L 45/586*  
(2013.01); *H04L 45/60* (2013.01)

(71) Applicant: **Fortinet, Inc.**, Sunnyvale, CA (US)(72) Inventors: **Rajesh I. Balay**, Cupertino, CA (US);  
**Girish Bhat**, San Diego, CA (US);  
**Gregory Lockwood**, Redwood City, CA  
(US); **Rama Krishnam Nagarajan**,  
Sunnyvale, CA (US)(73) Assignee: **Fortinet, Inc.**, Sunnyvale, CA (US)(21) Appl. No.: **14/714,270**(22) Filed: **May 16, 2015****Related U.S. Application Data**

(60) Continuation of application No. 14/616,521, filed on Feb. 6, 2015, which is a continuation of application No. 13/756,071, filed on Jan. 31, 2013, now Pat. No. 8,953,513, which is a continuation of application No. 13/015,880, filed on Jan. 28, 2011, now Pat. No. 8,369,258, which is a continuation of application No. 12/328,858, filed on Feb. 12, 2009, now Pat. No. 8,213,347, which is a division of application No. 10/949,943, filed on Sep. 24, 2004, now Pat. No. 7,499,419.

(57) **ABSTRACT**

Methods, apparatus and data structures are provided for managing multicast IP flows. According to one embodiment, a network switch module includes a memory and multiple processors partitioned among multiple virtual routers (VRs). Each VR maintains a data structure containing therein information regarding the multicast sessions, including a first value for each of the multicast sessions, at least one chain of one or more blocks of second values and one or more transmit control blocks (TCBs). Each first value is indicative of a chain of one or more blocks of second values. Each second value corresponds to an outbound interface (OIF) participating in the multicast session and identifies a number of times packets associated with the multicast session are to be replicated. The TCBs have stored therein control information to process or route packets. Each second value is indicative of a TCB that identifies an OIF of the network device through which packets are to be transmitted.



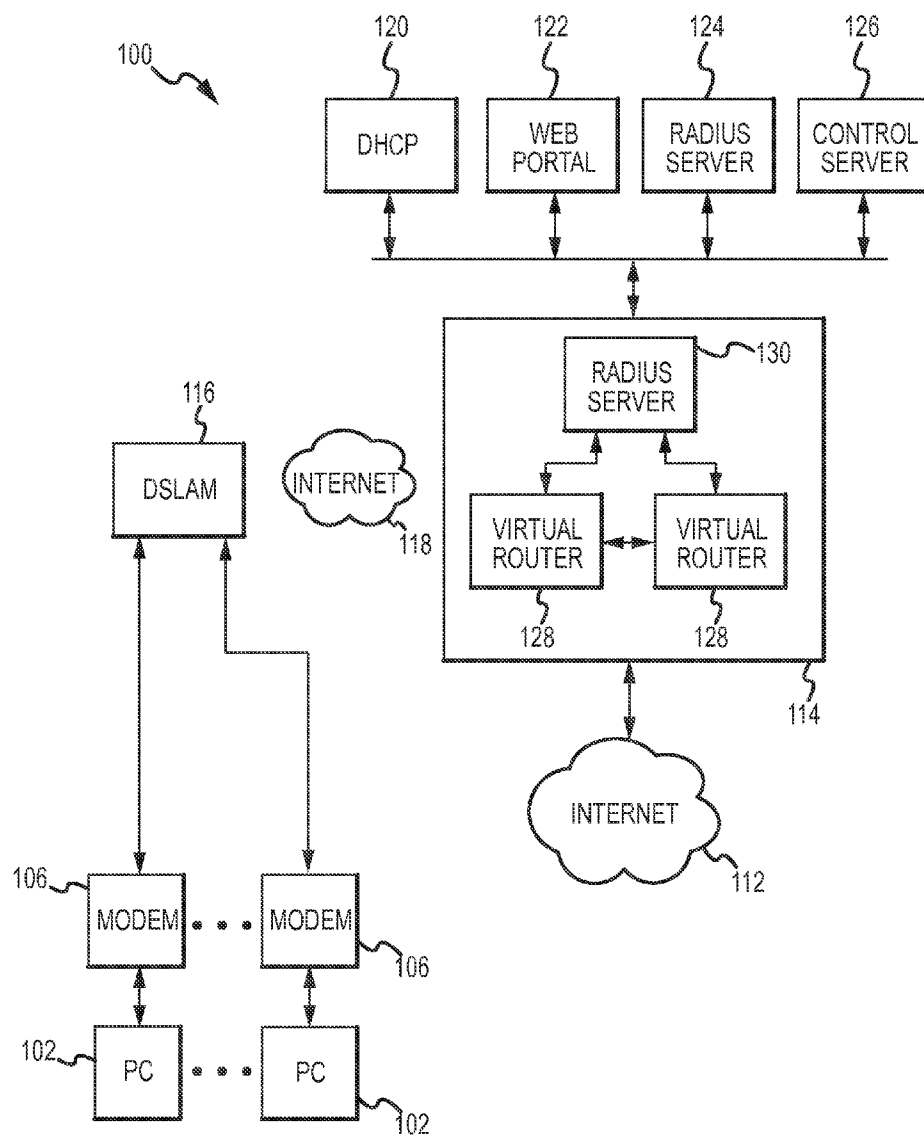


FIG.1

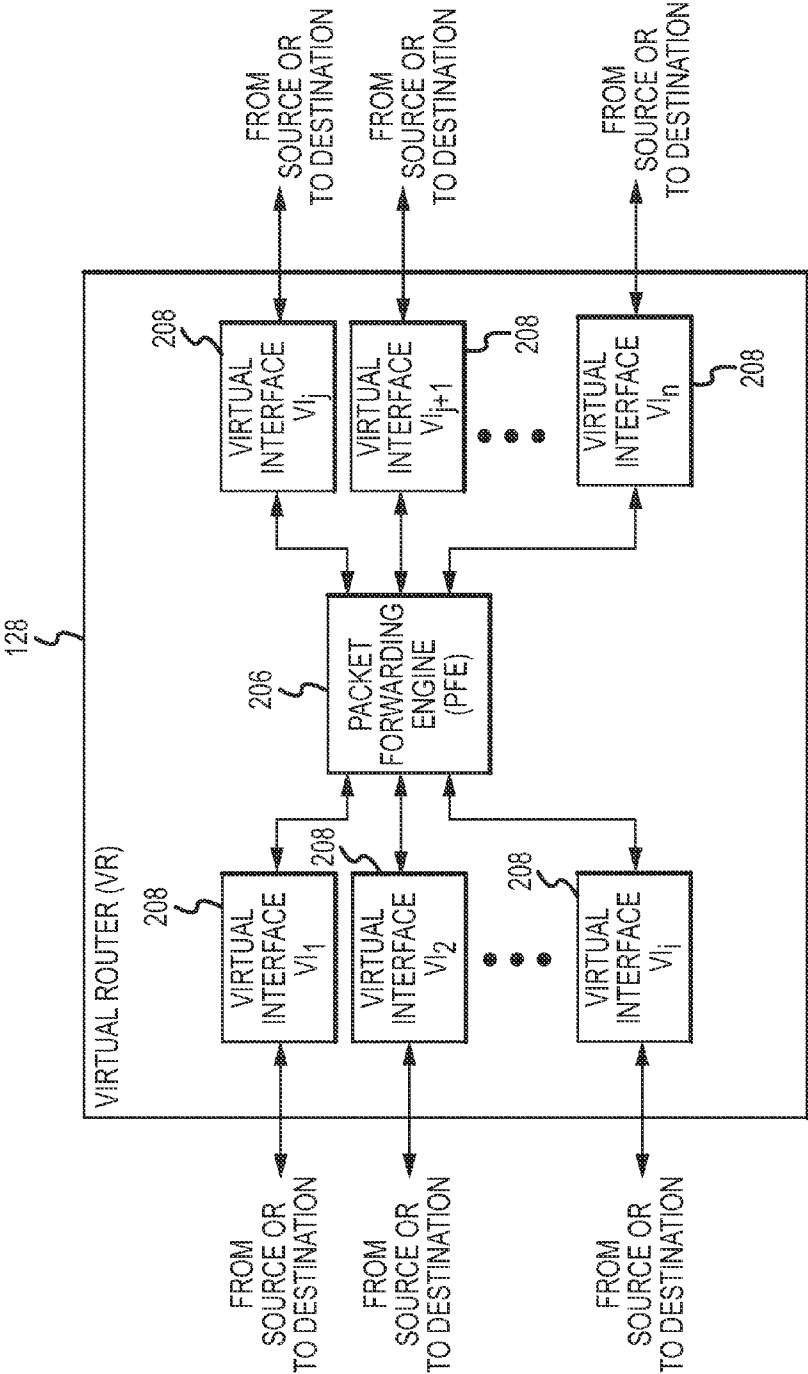


FIG.2

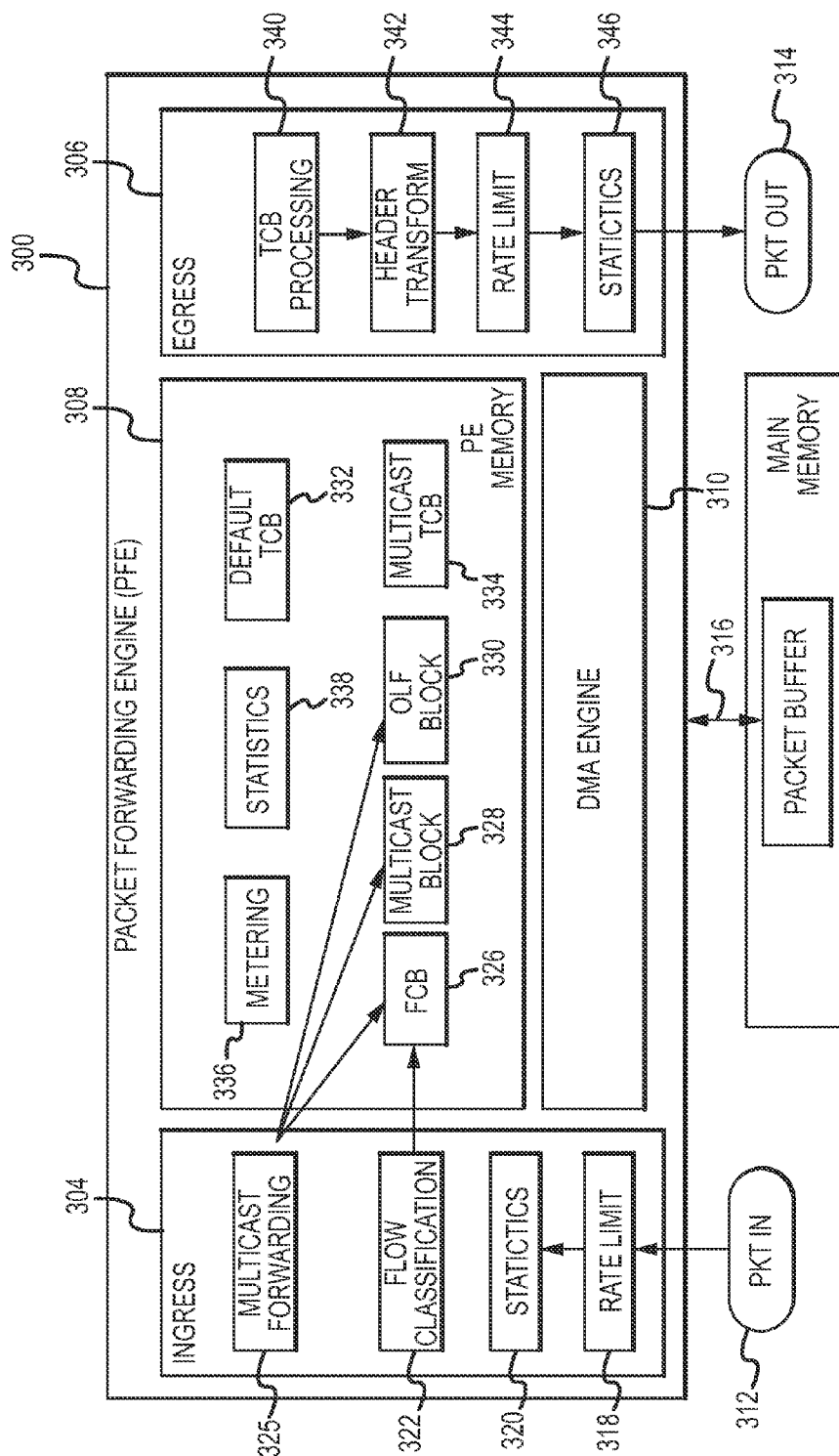


FIG.3

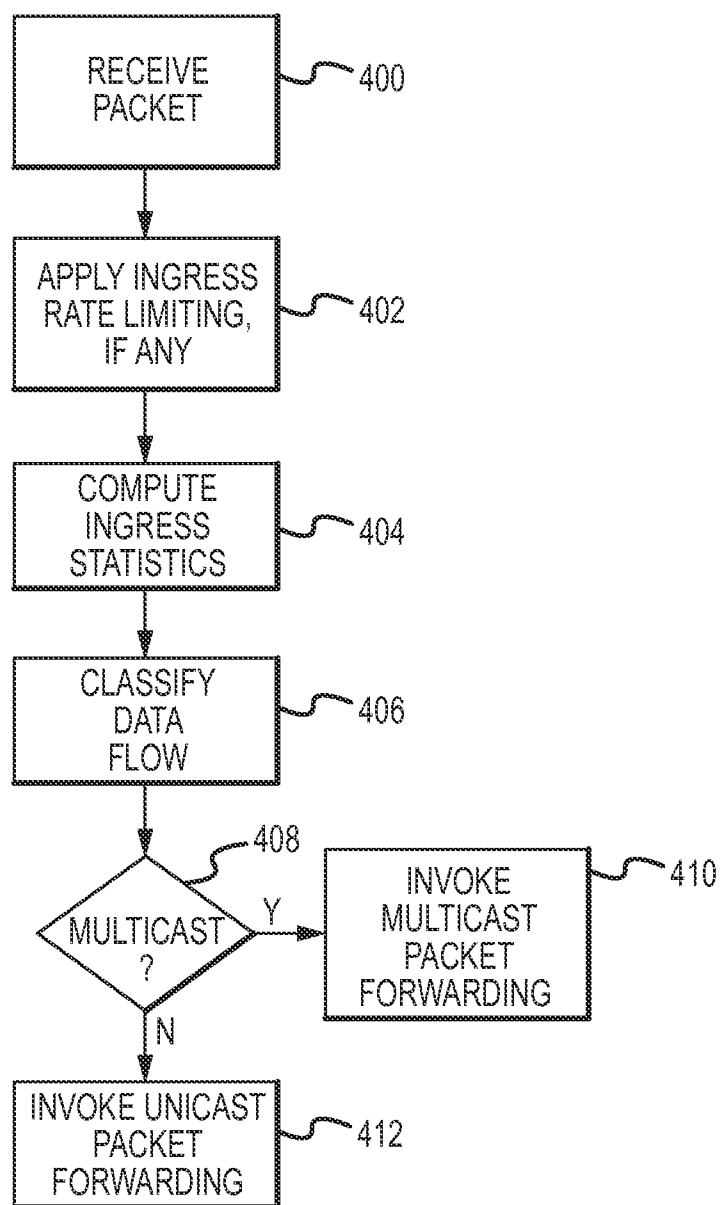
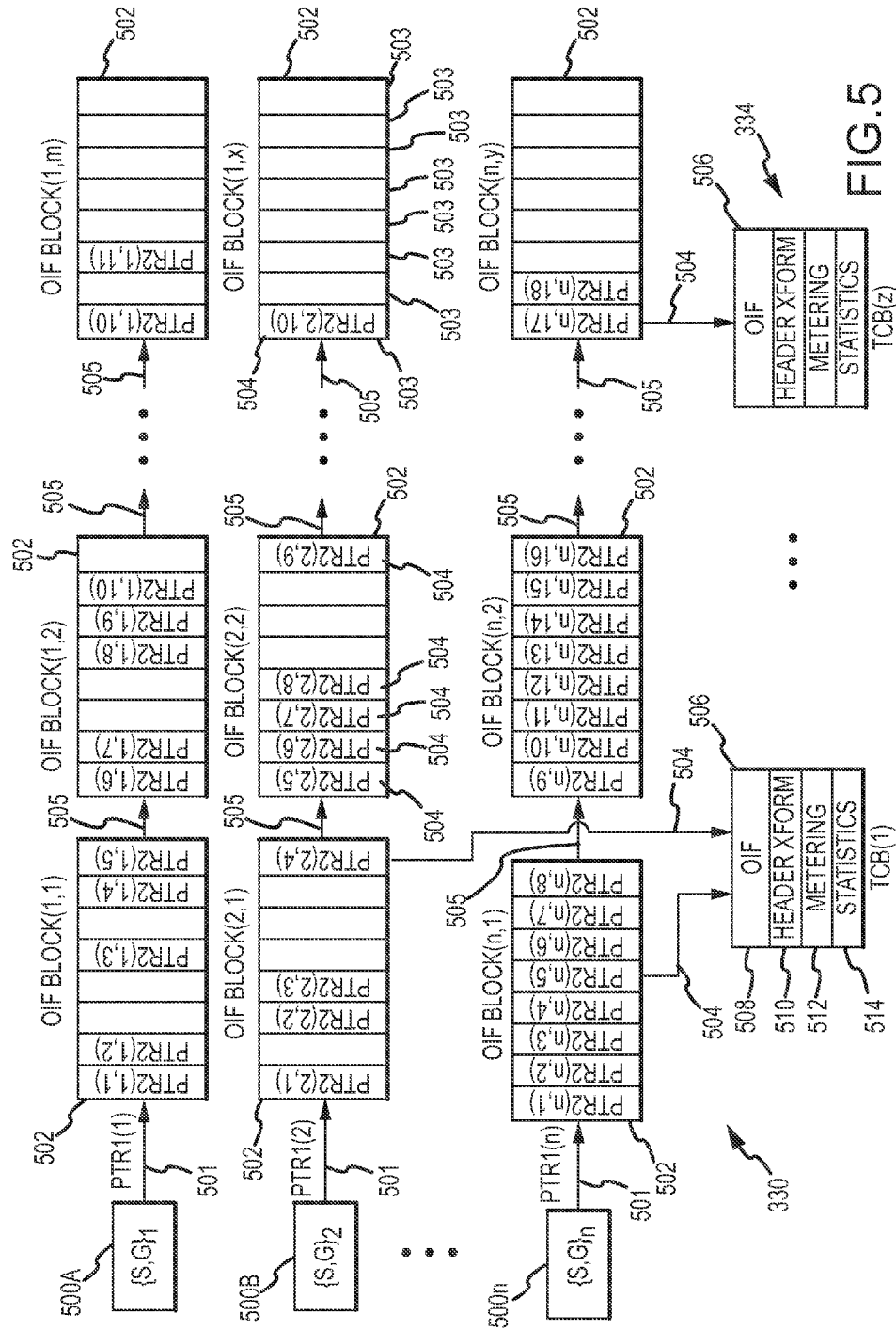


FIG.4



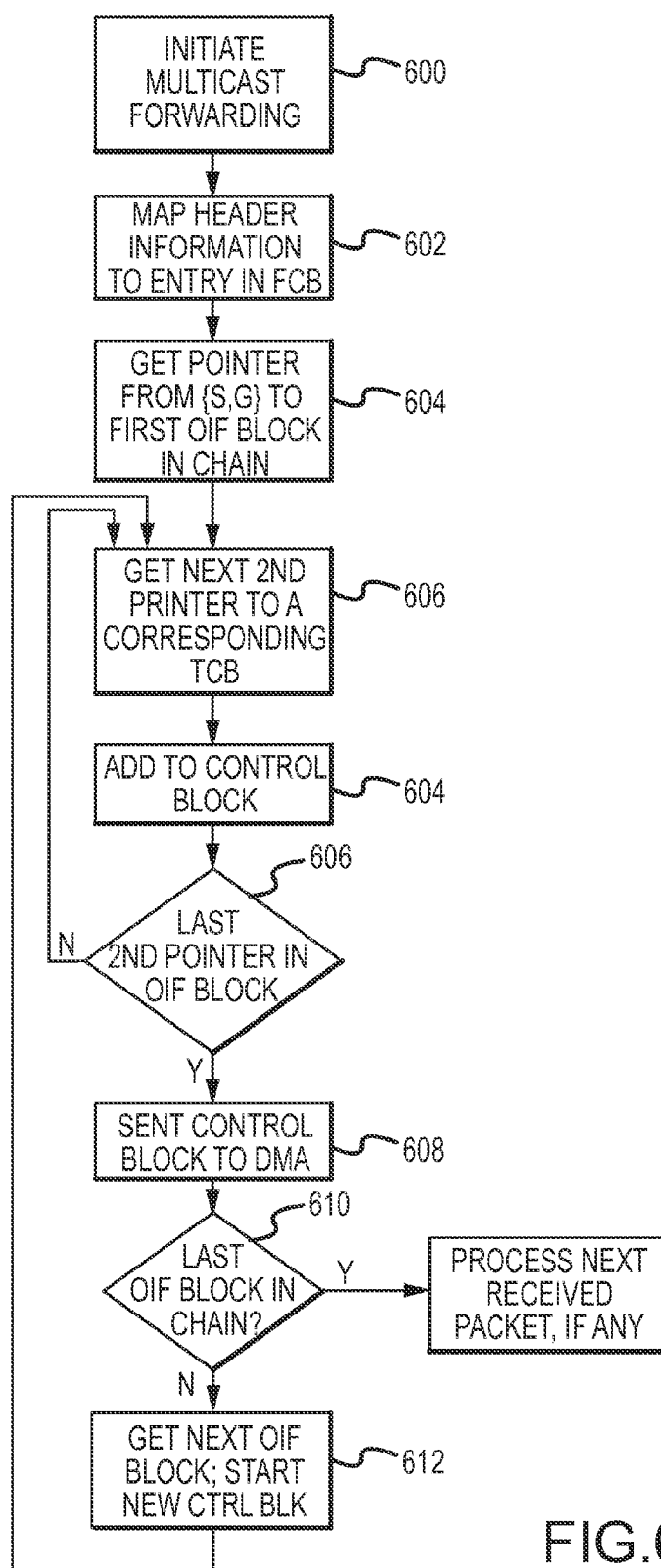


FIG. 6

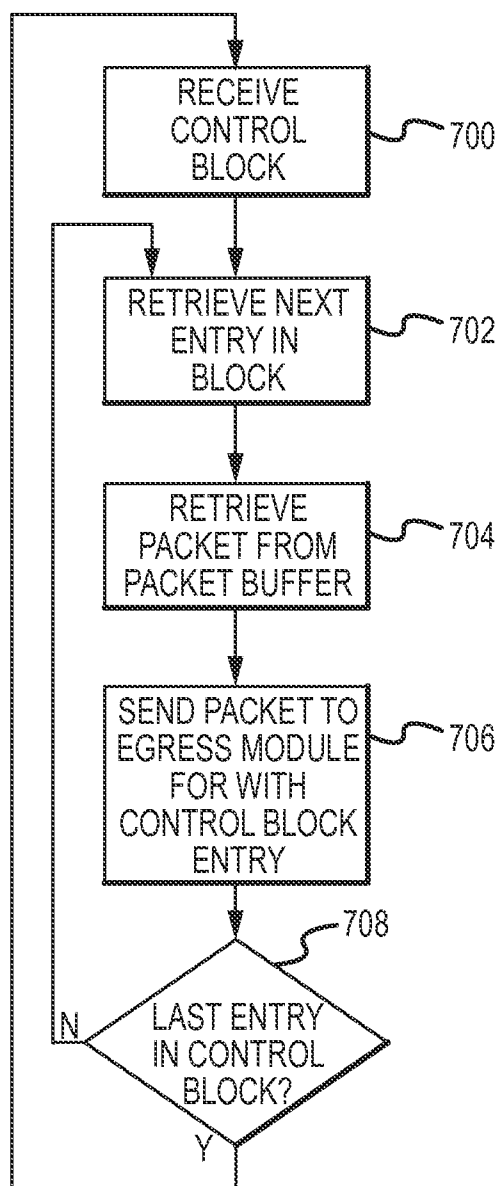


FIG.7



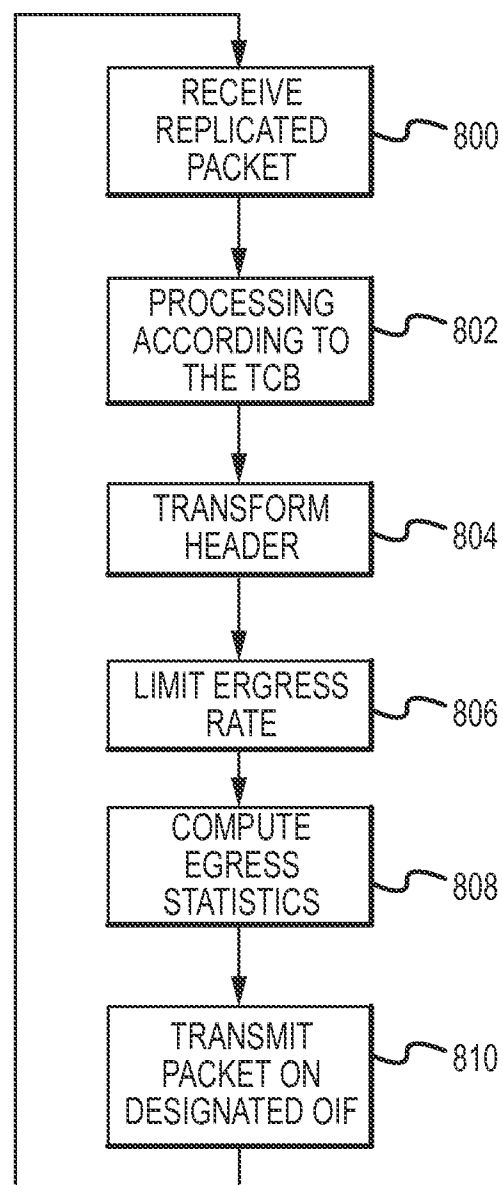
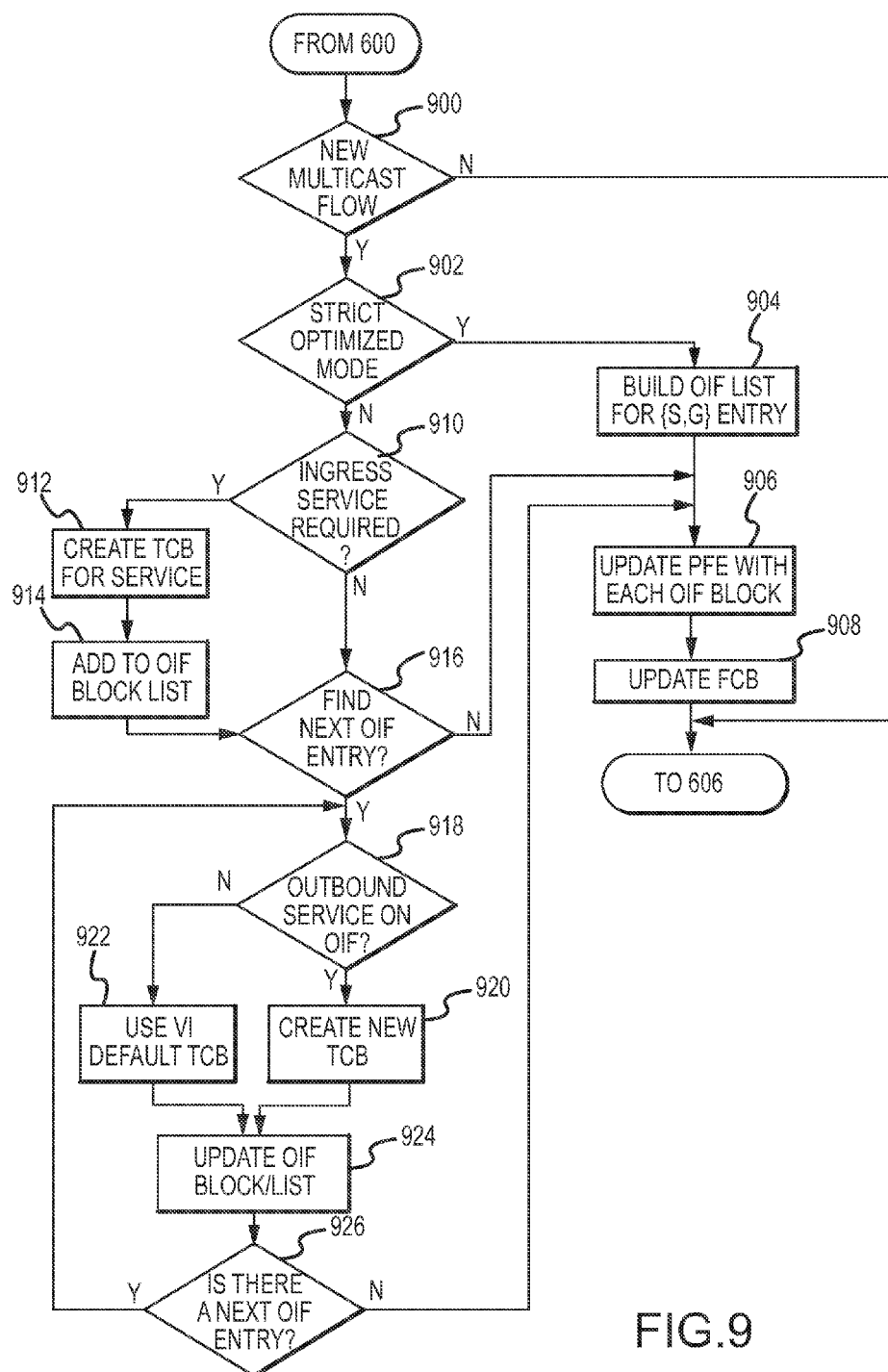


FIG.8



# SCALABLE IP-SERVICES ENABLED MULTICAST FORWARDING WITH EFFICIENT RESOURCE UTILIZATION

## CROSS-REFERENCE TO RELATED APPLICATIONS

**[0001]** This application is a continuation of U.S. patent application Ser. No. 14/616,521, filed Feb. 6, 2015, which is a continuation of U.S. patent application Ser. No. 13/756,071, filed Jan. 31, 2013, now U.S. Pat. No. 8,953,513, which is a continuation of U.S. patent application Ser. No. 13/015,880, filed Jan. 31, 2011, now U.S. Pat. No. 8,369,258, which is a continuation of U.S. patent application Ser. No. 12/328,858, filed Feb. 12, 2009, now U.S. Pat. No. 8,213,347, which is a divisional of U.S. patent application Ser. No. 10/949,943 filed Sep. 24, 2004, now U.S. Pat. No. 7,499,419, all of which are hereby incorporated by reference in their entirety for all purposes.

## COPYRIGHT NOTICE

**[0002]** Contained herein is material that is subject to copyright protection. The copyright owner has no objection to the facsimile reproduction of the patent disclosure by any person as it appears in the Patent and Trademark Office patent files or records, but otherwise reserves all rights to the copyright whatsoever. Copyright© 2004-2015, Fortinet, Inc.

## BACKGROUND

### **[0003]** 1. Field

**[0004]** Various embodiments of the present invention are generally related to the field of telecommunications and more particularly, but not by way of limitation, to network switches and systems and methods for multicast internet protocol (IP) forwarding.

### **[0005]** 2. Description of the Related Art

**[0006]** The use of computer or communications networks, including Local Area Networks (LANs), Wide-Area Networks (WANs), and the Internet continues to grow at ever increasing rates. Each day, more and more computer systems or communications devices are becoming interconnected in such wired or wireless networks, which typically communicate data in packets. This has created a need for high performance network switches, such as for use by network service providers. Many such switches comprise multiple modules, with many data flows between the modules themselves and between the interfaces to external networks. A data flow is sometimes called an “IP flow,” which refers to a stream of packets that enter and exit the same set of interfaces. The packets of a particular IP flow have the same values in the IP packet header for the following six attributes of the IP packet header: (1) Source IP Address, (2) Source L4 Port, (3) Type of Service (TOS), (4) Destination IP Address, (5) Destination L4 Port, and (6) Protocol.

**[0007]** In some cases, the network switch modules, including the processors residing in the modules, can be partitioned into virtual routers (VRs), that is, software running on the processors that emulates the functioning of an individual physical hardware router. As a result of the combination of hundreds of thousands of data flows for the virtual routers in these network switches, there is a need for efficiently processing packet data flows, and for controlling the resources consumed within the network switch.

**[0008]** As broadband network access becomes more available, individual subscribers of network service providers have more available options for different services and service levels. Even the same subscriber may have different service needs at different times. As an illustrative example, a first subscriber may desire high definition television (HDTV) service over a network. A second subscriber may desire mobile telephone service over the network. The first subscriber may occasionally desire video-on-demand (VOD). The second subscriber may need to switch between voice communication and high-speed digital data communication.

**[0009]** A “unicast” communication typically refers to a communication from a single source device to a single destination device over a network. By contrast, a “multicast” communication typically refers to a communication to a group of destination devices from one or more source devices. Multicast packet forwarding raises additional complexity because of the many destination devices. Many existing router devices will be unable to provide the desired scalability to accommodate such additional destination devices. This is particularly true when each individual data flow may require “per-flow” services for the multicast traffic. Allocating resources efficiently for a large number of multicast data flows is a challenging problem. Moreover, multicast broadcasting of content presents additional complexity because individual users may join or leave a particular multicast group at will and often. Such “channel surfing” creates an additional burden for keeping track of the participants of a multicast group so that the content can be routed appropriately.

## SUMMARY

**[0010]** Methods, apparatus and data structures for managing multicast Internet Protocol (IP) flows are described. According to one embodiment, a network switch module includes a memory and multiple processors partitioned among multiple virtual routers (VRs). Each VR maintains a data structure containing therein information regarding the active multicast IP sessions, including a first value (e.g., a pointer) for each of the active multicast IP sessions, at least one chain of one or more blocks of second values (e.g., pointers) and one or more transmit control blocks (TCBs). Each first value is indicative of a chain of one or more blocks of second values. Each second value is indicative of an out-bound interface (OIF) of the network device participating in the active multicast IP session defined by the first value and identifies a number of times packets associated with the active multicast IP session are to be replicated. The TCBs have stored therein control information to process or route packets. Each second value is indicative of a TCB. Each TCB identifies an OIF of the network device through which packets are to be transmitted.

**[0011]** Other features of embodiments of the present invention will be apparent from the accompanying drawings and from the detailed description that follows.

## BRIEF DESCRIPTION OF THE DRAWINGS

**[0012]** Embodiments of the present invention are illustrated by way of example, and not by way of limitation, in the figures of the accompanying drawings and in which like reference numerals refer to similar elements and in which:

**[0013]** FIG. 1 is a block diagram of one example of an operating environment for the present system and methods.

**[0014]** FIG. 2 is a block diagram of one example of a Virtual Router (VR) in accordance with an embodiment of the present invention.

**[0015]** FIG. 3 is a block diagram of one example of a Packet Forwarding Engine (PFE) and a main memory of a VR in accordance with embodiments of the present invention.

**[0016]** FIG. 4 is a flow chart illustrating a multicast internet protocol (IP) packet forwarding method in accordance with an embodiment of the present invention.

**[0017]** FIG. 5 is a block diagram of an OIF module and multicast TCB module for a set of multicast sessions in accordance with an embodiment of the present invention.

**[0018]** FIG. 6 is a flow chart illustrating a process that is invoked in accordance with an embodiment of the present invention if multicast packet forwarding is invoked in FIG. 4.

**[0019]** FIG. 7 is a flow chart illustrating packet retrieval and replication in accordance with an embodiment of the present invention.

**[0020]** FIG. 8 is a flow chart illustrating processing by an egress module in accordance with an embodiment of the present invention.

**[0021]** FIG. 9 is a flow chart illustrating more detail of acts included in the multicast forwarding in accordance with an embodiment of the present invention.

#### DETAILED DESCRIPTION

**[0022]** Methods, apparatus and data structures for multicast internet protocol (IP) forwarding are described herein. In the following description, numerous specific details are set forth. However, it is understood that embodiments of the invention may be practiced without these specific details. In other instances, well-known circuits, structures and techniques have not been shown in detail in order not to obscure the understanding of this description. Note that in this description, references to “one embodiment” or “an embodiment” mean that the feature being referred to is included in at least one embodiment of the invention. Further, separate references to “one embodiment” in this description do not necessarily refer to the same embodiment; however, neither are such embodiments mutually exclusive, unless so stated and except as will be readily apparent to those of ordinary skill in the art. Thus, the present invention can include any variety of combinations and/or integrations of the embodiments described herein. Moreover, in this description, the phrase “exemplary embodiment” means that the embodiment being referred to serves as an example or illustration.

**[0023]** Herein, block diagrams illustrate exemplary embodiments of the invention. Also herein, flow diagrams illustrate operations of the exemplary embodiments of the invention. The operations of the flow diagrams will be described with reference to the exemplary embodiments shown in the block diagrams. However, it should be understood that the operations of the flow diagrams could be performed by embodiments of the invention other than those discussed with reference to the block diagrams, and embodiments discussed with references to the block diagrams could perform operations different than those discussed with reference to the flow diagrams. Moreover, it should be understood that although the flow diagrams may depict serial operations, certain embodiments could perform certain of those operations in parallel.

**[0024]** The following detailed description includes references to the accompanying drawings, which form a part of the detailed description. The drawings show, by way of illustrat-

tion, specific embodiments in which the invention may be practiced. These embodiments, which are also referred to herein as “examples,” are described in enough detail to enable those skilled in the art to practice the invention. The embodiments may be combined, other embodiments may be utilized, or structural, logical and electrical changes may be made without departing from the scope of the present invention. The following detailed description is, therefore, not to be taken in a limiting sense, and the scope of the present invention is defined by the appended claims and their equivalents.

**[0025]** In this document, the terms “a” or “an” are used, as is common in patent documents, to include one or more than one. In this document, the term “or” is used to refer to a nonexclusive or, unless otherwise indicated. Furthermore, all publications, patents, and patent documents referred to in this document are incorporated by reference herein in their entirety, as though individually incorporated by reference. In the event of inconsistent usages between this document and those documents so incorporated by reference, the usage in the incorporated reference(s) should be considered supplementary to that of this document; for irreconcilable inconsistencies, the usage in this document controls.

**[0026]** Some portions of the following detailed description are presented in terms of algorithms and symbolic representations of operations on data bits within a computer memory. These algorithmic descriptions and representations are the ways used by those skilled in the data processing arts to most effectively convey the substance of their work to others skilled in the art. An algorithm includes a self-consistent sequence of steps leading to a desired result. The steps are those requiring physical manipulations of physical quantities. Usually, though not necessarily, these quantities take the form of electrical or magnetic signals capable of being stored, transferred, combined, compared, and otherwise manipulated. It has proven convenient at times, principally for reasons of common usage, to refer to these signals as bits, values, elements, symbols, characters, terms, numbers, or the like. It should be borne in mind, however, that all of these and similar terms are to be associated with the appropriate physical quantities and are merely convenient labels applied to these quantities. Unless specifically stated otherwise as apparent from the following discussions, terms such as “processing” or “computing” or “calculating” or “determining” or “displaying” or the like, refer to the action and processes of a computer system, or similar computing device, that manipulates and transforms data represented as physical (e.g., electronic) quantities within the computer system’s registers and memories into other data similarly represented as physical quantities within the computer system memories or registers or other such information storage, transmission or display devices.

**[0027]** FIG. 1 is a block diagram of one example of an operating environment for the present system and methods. In the example of FIG. 1, a system 100 typically includes personal computers (PCs) 102 that are respectively connected to modems 106. The modems 106 are typically respectively connected to a digital subscriber line access module (DSLAM) 116. The DSLAM 116 multiplexes signals from the modems 106 onto the Internet Protocol (IP) network 118. The IP network 118 is typically connected to a router box 114 that includes virtual routers (VRs) 128. The router box 114 is typically connected to the Internet 112. The router box 114 is

also typically connected to a dynamic host configuration protocol (DHCP) server **120**, a web portal **122**, a RADIUS server **124**, and a control server **126**.

[0028] Although, in this example, the router **114** includes three VRs **128**, other examples call for any number of VRs **128**. In one example, one or more of the VRs **128** can establish subscriber connections, such as to users of the PCs **102**. When establishing such connections, the VRs **128** can use the DHCP server **120** for assigning IP network addresses to the PCs **102**. The VRs **128** can use the RADIUS server **124** to authenticate subscribers. After authenticating subscribers, the VRs **128** can configure subscriber connections according to service profiles, which refer to subscriber-specific services that individual subscribers receive during connections. In one example, the VRs **128** can receive service profiles information from the control server **126** or the RADIUS server **224**.

[0029] After the VRs **128** establish subscriber connections, they typically provide access to the web portal **122**, where users can select new services. Additionally, after establishing subscriber connections, the VRs **128** typically process and forward packets over the IP network **118** and the Internet **112**. Although FIG. 1 illustrates an example in which the users accessing the Internet **112** via PCs, this is merely an illustrative example. In other examples, the individual users may access the Internet **112** or other computer or communications network wirelessly, such as by using a 3rd Generation (3G) or other mobile phone or other handheld or portable device, or by using a laptop or other portable computing device with Wireless Fidelity (WiFi) (e.g., using IEEE 802.11b wireless networking) capability or the like. In still other examples, the individual users may access the Internet **112** or other communications or computer network using an Ethernet connection, a Very High bit rate Digital Subscriber Line (VHDSL) connection, a Fiber To The Premises (FTTP) connection, or a cable TV line or like connection, or the like. Thus, the present systems and methods are not limited to any particular devices or techniques for accessing the Internet **112**, such as through the router box **114**. The exemplary system **100** typically provides network services to thousands of subscribers. Each subscriber can receive a particular set of services upon establishing a connection with the system **100**.

[0030] FIG. 2 is a block diagram of one example of a VR **128**. In this example, the VR **128** includes a packet forwarding engine (PFE) **206**, and one or more virtual interfaces (VIs) **208** from a source or to a destination. A VI over which multicast packets are forwarded is sometimes referred to as an Outbound Interface (OIF). Different services can be applied to multicast packet forwarding traffic as well as unicast packet forwarding traffic. Which services are applied to a particular packet are determined, in one example, by an inbound policy or an outbound policy associated with a particular VI **208**. In one example, the packet header (e.g., one or more of the above-described six packet header attributes defining an IP flow) is examined (e.g., such as by comparing such attribute (s) to one or more matching criteria in an access control list (ACL)) to determine whether any services should be applied to the packet and, if so, which services should be applied.

[0031] FIG. 3 is a block diagram of one example of a Packet Forwarding Engine (PFE) **300** and a main memory **302** of a VR **128**. In this example, PFE **300** includes an ingress module **304**, an egress module **306**, a PFE memory **308**, a Direct Memory Access (DMA) engine **310**, a packet input interface **312**, a packet output interface **314**, and a main memory interface **316**. In this example, the ingress module **304** includes an

ingress rate limit module **318**, an ingress statistics module **320**, a flow classification module **322**, and a multicast forwarding module **324**. In this example, the PFE memory **308** includes a Flow Control Block (FCB) **326**, a multicast block **328**, an Outbound InterFace (OIF) module **330**, a default Transmit Control Block (TCB) **332**, a multicast TCB module **334**, a metering block **336**, and a statistics block **338**. In this example, the egress module **306** includes a TCB processing module **340**, a header transform module **342**, an egress rate limit module **344**, and an egress statistics module **346**.

[0032] FIG. 4 is a flow chart of one example of a multicast internet protocol (IP) packet forwarding method such as can be performed, for example, by using the PFE **300** and main memory **302** of FIG. 3. At **400**, a packet is received at the packet input interface **312** on a particular virtual interface (VI) of a particular VR in the router box **114**. When a packet is received at **400**, it is not known whether the packet is part of a multicast data flow or a unicast data flow. At **402**, in one example, the ingress rate limit module **318** performs a rate limiting function to control a packet ingress data rate through the ingress module **304**. At **404**, in one example, the ingress statistics module **320** computes packet ingress statistics, such as packet count, byte count, etc. Such ingress statistics may be important for managing subscriber service levels, among other things.

[0033] At **406**, in one example, the flow classification module **322** is used to classify the data flow, for example, as a unicast flow or a multicast flow. The flow classification module **322** typically uses a predefined portion of the packet header to classify the data flow, and to identify the particular FCB associated with the flow. For example, the “destination address” portion of the packet header is used to identify the packet as a multicast packet. In one example, in a first mode (sometimes referred to as a “strict-optimized mode”), the data flow classification uses the source IP address and the destination IP address portions of the packet header to classify the data flow. In a second mode (sometimes referred to as an “adaptive-optimized mode”), in which subscriber-specific services are needed, additional portions of the packet header are used to further classify the data flow in accordance with the appropriate services.

[0034] In one example, the flow classification at **406** uses the information extracted from the packet header to look up a corresponding FCB entry in FCB **326**. If the data flow is a multicast data flow then, in one example, the corresponding FCB entry will have a “multicast” flag set, and a “forwarding action” field of the FCB entry will indicate that hardware forwarding of packets is to be used for the multicast data flow. At **408**, if the classification indicates a multicast data flow, then, at **410**, multicast packet forwarding is invoked. Otherwise, at **412**, unicast packet forwarding is invoked.

[0035] Each FCB entry in FCB **326** includes information identifying a particular multicast session. Each multicast session is defined by a {Source, Group} pair, which is sometimes referred to as an {S, G} pair. The Source field of the {S, G} pair defines the source of the multicast transmission. In one example, this is a single multicast transmission source. In another example, there are multiple (e.g., redundant) transmission sources for the same multicast transmission. The Group field of the {S, G} pair defines a group corresponding to the multicast session. In one example, the group can be conceptualized as a “channel” of content. There may be one recipient or a very large number of recipients of the content. Such recipients of the multicast content can join or leave the

Group at will, such as by issuing the appropriate Internet Group Management Protocol (IGMP) request or using one or more other protocols. Thus, scalability and the ability to easily update the Group are desirable qualities of the present multicast forwarding systems and methods.

**[0036]** Since each multicast session can have multiple IP flows associated with that particular multicast session, there can be multiple FCBs associated with the same  $\{S, G\}$ , where each FCB corresponds to one of these IP flows, and the  $\{S, G\}$  defines the particular multicast session. This may be true, for example, in the adaptive-optimized mode case, where because of the different services levels needed, there are different IP flows associated with the same multicast session.

**[0037]** FIG. 5 is a block diagram of one example of an OIF module 330 and multicast TCB module 334 for a set of multicast sessions defined by respective  $\{S, G\}$  pairs  $\{S, G\}_1$  through  $\{S, G\}_n$ . The  $\{S, G\}$  pair 500 of a particular multicast session includes a first pointer 501 that points to a dynamically allocated set of OIF Blocks 502. The particular number of OIF Blocks 502 depends on how many OIFs are then participating in that multicast session. For a particular multicast session, each OIF block 502 points to a subsequent OIF block 502 (with the exception of the last OIF block 502 in this conceptual “chain” of OIF blocks).

**[0038]** Each OIF block includes a reasonably small number of slots 503 for storing corresponding second pointers 504 to a TCB 506 for a particular OIF. The example of FIG. 5 illustrates eight slots 503 per OIF block 502, each slot for storing a corresponding second pointer 504 to a TCB 506. Another example includes six second pointer slots 503 per OIF block 502. Each second pointer 504 points to a particular TCB 506 for a particular OIF, which may service one or more users participating in the corresponding multicast session. Each OIF block 502 typically has the same number of second pointer slots 503 as every other OIF block 502, however, the number of OIF blocks 502 can vary between different  $\{S, G\}$  pairs, or even for the same  $\{S, G\}$  pair, such as at different points in time when different numbers of OIFs are part of that particular multicast session. More particularly, as OIFs are added or removed from a multicast session (such as may happen when users join or leave the multicast session) corresponding second pointers 504 are added or removed, respectively. If needed, additional OIF blocks 502 are added or removed, such as to accommodate the addition or removal of the second pointers 504. Using the present systems and methods, dynamically adding or removing such OIF blocks 502 as needed is easy because, among other things, each multicast session includes OIF blocks 502 that are chained together by third pointers 505 from one OIF block 502 to another OIF block 502 (except for the last OIF block 502 in the chain). When a user joins or leaves a multicast session under circumstances that require adding or removing an OIF to that multicast session, the OIF list can be updated by simply updating a single OIF block 502, during which time the other OIF blocks 502 in that chain are still available and usable for performing multicast forwarding. Although FIG. 5 illustrates a typical example in which each multicast session (defined by a particular  $\{S, G\}$  pair) points to its own chain of OIF blocks 502, it is possible that, in one example implementation, different multicast sessions point to the same (shared) chain of OIF blocks 502. This will likely be the less typical case, for example, in which these two or more different multicast sessions each have the same OIFs participating in that multicast session. This can be conceptualized as two or more different

channels that are being “watched” by the same OIFs. When this occurs, the pointers from each such multicast session can point to the same (shared) chain of OIF blocks 502, if desired. Alternatively, separate chains of OIF blocks 502 can be maintained for each multicast session, for example, if such a simplified implementation is desired.

**[0039]** Each second pointer 504 points to a particular TCB 506, which typically includes information relevant to processing or routing packets to the particular OIF that is associated with that second pointer 504, or to services associated with the particular OIF that is associated with that second pointer 504. For example, if the packet header matches particular services in the ACL, attributes in the TCB are adjusted accordingly to obtain such services. Each second pointer 504 corresponds to a particular outbound interface (OIF) through which multicast packets are being forwarded, such as from the packet output interface 314 of the VR out over the network.

**[0040]** Because more than one multicast session can use the same OIF of the VR, second pointers 504 from different multicast sessions can point to the same (shared) TCB 506 for that OIF. In the illustrative example of FIG. 5, the second pointer PTR2(2,4) from the second multicast session points to the shared TCB(1) as the second pointer PTR2(n,5) from the nth multicast session. Thus, second pointers 504 from different multicast sessions may share the same TCB 506.

**[0041]** Similarly, because multiple IP flows can use the same OIF, there can be multiple TCBS 506 for the same OIF, such as for multiple IP flows on the same OIF, where such multiple flows use different services and, therefore, have different corresponding TCBS 506.

**[0042]** In FIG. 5, for example, a particular TCB 506 typically includes, among other things, OIF information 508, header transformation information 510, metering information 512, and statistics information 514. The OIF information 508 includes, for example, information identifying which OIF will be used by the packet output interface 314 to output the packets from the VR. The header transformation information 510 includes, for example, Media Access Control (MAC) address generation information and protocol independent multicast (PIM) encapsulation information for that particular OIF. The metering information 512 includes, for example, egress rate limiting or other like information for that particular OIF. The statistics information 514 includes egress statistics collection information for that particular OIF.

**[0043]** The schema depicted in FIG. 5 provides numerous advantages. As discussed above, scalability from one to very many users is a desirable property. The ability to update the multicast forwarding schema as many users join or leave different multicast sessions (which sometimes results in adding or removing OIFs) is another desirable property. For example, when a user joins or leaves a multicast session under circumstances that require adding or removing an OIF to that multicast session, the OIF list can be updated by simply updating a single OIF block 502, during which time the other OIF blocks 502 in that chain are still available and usable for performing multicast forwarding.

**[0044]** The schema depicted in FIG. 5 allows many users to be managed very efficiently because of, among other things, its use of first pointers 501 from  $\{S, G\}$  pairs to shared or independent chains of OIF blocks 502, and per-OIF second pointers 504 to shared or independent TCBS 506. Moreover, each OIF block 502 is typically apportioned into a small number of second pointer slots 503. Each OIF block 502 is

typically independently addressable and updatable when updating that OIF block to add or remove a particular OIF's second pointer 504. As an illustrative example, if the OIF corresponding to second pointer PTR2(1,3) in OIF Block (1,1) was removed from the multicast session of {S, G}<sub>1</sub> (for example, because all of the one or more users of that OIF left that multicast session), then the second pointer PTR2(1, 3) in that OIF Block (1, 1) is removed, opening one second pointer slot 503 in OIF Block (1, 1) that could later be filled by another second pointer for another OIF being added (e.g., to service one or more users joining that multicast session).

[0045] While such updating of a particular OIF block 502 is occurring, other OIF blocks 502 in the same or a different chain of OIF blocks 502 are still usable to carry out multicast forwarding to the users represented by the second pointers 504 in those other OIF blocks 502. This improves the ability to multicast content, without interruption, to a large number of recipient users on different OIFs of a particular multicast session, even as other second pointers 504 are added or removed, such as to accommodate other recipient users of that multicast session that are joining or leaving that multicast session. In one example, both OIF blocks 502 and TCBs 506 are capable of being dynamically allocated as needed. Together with the sharing of TCBs 506 or even of OIF chains, as discussed above, the schema illustrated in FIG. 5 typically offers one or more of the advantages of scalability, updatability, efficiency in memory usage, and high throughput performance with reduced interruptions.

[0046] FIG. 6 is a flow chart of one example of a process that is invoked if multicast packet forwarding is invoked at 410 of FIG. 4. At 600, a multicast forwarding operation is initiated, such as by calling executable or interpretable instructions of the multicast forwarding module 324. At 602, information in the packet header is mapped to an FCB entry in FCB module 326. Each FCB entry in the FCB module 326 identifies an IP flow or a group of IP flows for an {S, G} pair 500 corresponding to a particular multicast session. At 604, from that FCB entry, a first pointer 501 is extracted to the first OIF block 502 in the chain of one or more OIF blocks 502 corresponding to that multicast session.

[0047] At 606 the next second pointer 504 in the current OIF block 502 is retrieved. At 606, the retrieved second pointer 504 to a TCB 506 is used to build a portion of a control block that will be sent to the DMA engine 310. At 606, if other second pointers 504 exist in the current OIF block 502, then process flow returns to 606. Otherwise, process flow proceeds to 606 and the control block that was constructed for the completed OIF block 502 is sent to the DMA engine 310. In this manner, one control block corresponding to each OIF block 502 is sent to the DMA engine 310 after that control block is constructed from the corresponding OIF block 502. At 610, if other OIF blocks 502 exist in that chain, then the next OIF block 502 is retrieved and made the current OIF block, a new control block is initiated, and process flow returns to 606. Otherwise, at 610, if no other OIF blocks 501 exist in the chain, then process flow proceeds to 614 to process (or wait for) the next received packet (e.g., at 400 of FIG. 4).

[0048] FIG. 7 is a flow chart of one example of packet retrieval and replication. At 700, the DMA engine 310 receives a control block (such as from 608 in FIG. 6). At 702, the next entry in the received control block is retrieved. At 704, the stored packet is retrieved from a packet buffer in the main memory 302 by DMA engine 310. At 706, the retrieved

packet is sent to the egress module 306 for egress transmission, along with the corresponding control block entry, which provides information to the egress module 306 about how that particular packet is to be processed for the particular recipient user corresponding to the control block entry, which, in turn, corresponded to a particular second pointer 504, as discussed above. At 708, if there are more entries in the control block, then process flow returns to 702 to retrieve the next entry in the control block. Otherwise, process flow returns to 700 to receive (or wait for) another control block. In the manner illustrated in FIG. 7, a packet is held in the packet buffer in the main memory 302 so that it can be replicated. The replicated packets are sent to the egress module 306 for further processing (particular to the user that is to receive that replicated packet) and transmission that OIF.

[0049] FIG. 8 is a flow chart of one example of processing by the egress module 306. At 800, a replicated packet is received from the DMA engine 310. At 802, the replicated packet is processed according to the TCB 506 corresponding to the particular OIF's second pointer 504. As discussed above, such information is encapsulated into the control block that was submitted to the DMA engine 310, and communicated to the egress module 306 along with the replicated packet. At 804, header transformation occurs. In one example, this includes MAC address generation or encapsulation appropriate for the designated OIF over which the replicated packet will be transmitted. At 806, egress rate limiting, if any, is applied. At 808, egress statistics, if any, are computed. At 810, the replicated packet is transmitted out over the computer network on the designated OIF.

[0050] FIG. 9 is a flow chart of one example of more detail of acts included in the multicast forwarding, such as at 602 of FIG. 6 or elsewhere, as appropriate. Among other things, the flow chart of FIG. 9 illustrates one example of how TCBs 506 are created and, where possible, shared.

[0051] At 900, the system determines whether a received packet represents a new IP flow. This can be determined by looking at the above-described attributes in the packet header that identify a particular IP flow. If the packet corresponds to a previously identified multicast IP flow, then process flow proceeds to 606, and a previously defined FCB entry and a previously defined TCB 506 are used for further multicast forwarding processing. If a new flow is detected at 900, there will be no matching FCB entry in FCB 326. Therefore, for a new flow detected at 900, a new FCB entry will be created in FCB 326, as discussed below.

[0052] If a new flow is detected at 900, then, at 902, is it determined whether the new flow is a strict optimized mode or, instead, is in an adaptive optimized mode that provides one or more services for that particular flow. This determination is typically made using a configurable attribute.

[0053] At 902, if in the strict optimized mode, then, at 904, an OIF list (e.g., a chain of OIF blocks, as illustrated in FIG. 5) is built for the {S, G} entry corresponding to the newly identified multicast flow. Because no flow-specific services are required, this OIF list includes second pointers 504 to a default TCB corresponding to each OIF in that multicast session. This default TCB does not include any per-flow of ACL-based service-specific attributes. Instead this default TCB typically depends only on attributes and services that are applicable to all flows associated with the particular VI serving as the OIF. Each OIF participating in the multicast session will have a corresponding second pointer 504 to its corresponding default TCB. Then, at 906, the OIF module 330 of

the PFE 300 is updated with each OIF block 502 in the chain of OIF blocks that make up the OIF list of that particular multicast flow. The OIF module 330 of the PFE 300 is typically updated with such OIF blocks 502 on a block-by-block basis. Then, at 908, the FCB 326 of the PFE 300 is updated to include, for example, a pointer to the OIF list for the newly identified multicast flow. Then, process flow proceeds to 606 of FIG. 6 (or elsewhere, if appropriate).

[0054] At 902, if in the adaptive optimized mode instead of the strict optimized mode, then, at 910 it is determined whether any ingress services are needed. In one example, this includes checking for such ingress services on the VI 208 at which the packet is received. At 910, if one or more such ingress services are needed, then, at 912, a TCB 506 is created to control the providing of any such ingress services. (otherwise process flow proceeds to 916). Then, at 914, a second pointer 504 is created to point to this newly created TCB 506. This newly created TCB 506 for the ingress services includes an OIF field 508 that specifies a null OIF (the PFE 300 does not actually forward any packets out any such null OIF).

[0055] At 916, it is determined whether there is a next OIF entry (that is, a second pointer 504) in the OIF list for the new multicast flow. If there is no such next OIF entry (e.g., upon specification of an invalid {S, G} entry or a null OIF), then process flow proceeds to 906. Otherwise, at 918, it is determined whether any outbound services are needed on the next OIF entry in the OIF module 330. If so, then, at 920, a new TCB 506 is created for that OIF entry to control the providing of any such outbound services, otherwise, at 922, the VI default TCB 332 is used for that OIF entry. Then, at 924, a second pointer 504 is created to point to the new TCB 506 or the default TCB 332, as appropriate, and the OIF list for that multicast session is updated accordingly. Then, at 926, it is determined if there is a next OIF entry in the OIF list for the multicast session. If so, process flow returns to 918, otherwise process flow proceeds to 906.

[0056] Using the above process described with respect to FIG. 9, some TCBs 506 may be shared by multiple second pointers 504. For example, the default TCB 332 is likely shared by multiple second pointers 504. Other TCBs 506 may correspond to individual second pointers 504.

[0057] Although the above examples have been discussed with respect to a router box providing virtual routers (e.g., VRs 128), the present systems and methods are not so limited. For example, certain aspects of the present systems and methods are also applicable to alternative systems using hardware routers instead of the virtual routers.

[0058] It is to be understood that the above description is intended to be illustrative, and not restrictive. For example, the above-described embodiments (and/or aspects thereof) may be used in combination with each other. Many other embodiments will be apparent to those of skill in the art upon reviewing the above description. The scope of the invention should, therefore, be determined with reference to the appended claims, along with the full scope of equivalents to which such claims are entitled. In the appended claims, the terms “including” and “in which” are used as the plain-English equivalents of the respective terms “comprising” and “wherein.” Also, in the following claims, the terms “including” and “comprising” are open-ended, that is, a system, device, article, or process that includes elements in addition to those listed after such a term in a claim are still deemed to fall within the scope of that claim. Moreover, in the following

claims, the terms “first,” “second,” and “third,” etc. are used merely as labels, and are not intended to impose numerical requirements on their objects.

What is claimed is:

1. A network switch module comprising:

a memory partitioned among a plurality of virtual routers (VRs);

a plurality of processors partitioned among the plurality of VRs; and

wherein each VR of the plurality of VRs maintains a data structure in the memory, the data structure including information relating to a set of multicast sessions being handled by the VR and including:

a plurality of pairs of a source field and a group field ({S, G} pairs) stored in the memory, in which each pair of the plurality of {S, G} pairs defines a multicast session of the set of multicast sessions and wherein the source field defines a source of a multicast transmission and the group field defines a group corresponding to the multicast session;

a first pointer associated with each of the plurality of {S, G} pairs that points to a dynamically allocated set of OIF blocks, wherein a number of outbound interface (OIF) blocks in the set of OIF blocks is dependent upon how many of a plurality of OIFs of the VR are currently participating in the multicast session and the number of OIF blocks in the set of OIF blocks defines how many times packets associated with the multicast session are to be replicated;

a set of slots associated with each OIF block of the set of OIF blocks, each slot of the set of slots configured to store a second pointer to a transmit control block (TCB) which services one or more users participating in the multicast session and which represents a data structure configured to store control information relevant to processing or routing packets, including information regarding an OIF of the plurality of OIFs through which the packets are to be transmitted;

a third pointer associated with each OIF block of the set of OIF blocks to chain together the set of OIF blocks, wherein only one OIF block of the set of OIF blocks is updated responsive to users joining or leaving the multicast session.

2. The network switch module of claim 1, wherein each OIF block of the set of OIF blocks is independently accessible without affecting processing of other OIF blocks of the set of OIF blocks.

3. The network switch module of claim 1, wherein the TCB includes control information for processing replicated packets.

4. The network switch module of claim 3, wherein the control information includes one or more of header transformation control information, metering control information and statistics control information.

5. The network switch module of claim 3, wherein the second pointer associated with a first OIF block of the set of OIF blocks of a first {S, G} pair of the plurality of {S, G} pairs is permitted to point to and thus share a first TCB with another OIF block of the set of OIF blocks of the first {S, G} pair or with a second OIF block of the set of OIF blocks of a second {S, G} pair of the plurality of {S, G} pairs.



6. A method of managing multicast Internet Protocol (IP) sessions, the method comprising:

identifying, by a router, active multicast IP sessions; and maintaining, by the router, a data structure within a memory of the router containing therein information regarding the active multicast IP sessions;

wherein the data structure includes:

a plurality of pairs of a source field and a group field ( $\{S, G\}$  pairs), in which each pair of the plurality of  $\{S, G\}$  pairs defines a multicast IP session of the active multicast IP sessions, wherein the source field defines a source of a multicast transmission of the multicast IP session and the group field defines a group corresponding to the multicast IP session;

a first pointer associated with each of the plurality of  $\{S, G\}$  pairs that points to a dynamically allocated set of outbound interface (OIF) blocks, wherein a number OIF blocks in the dynamically allocated set of OIF blocks is dependent upon a number of OIFs of the router that are participating in the IP multicast session and the number of OIF blocks in the dynamically allocated set of OIF blocks defines a number of times packets of the IP multicast session are to be replicated;

a set of slots for each OIF block of the set of dynamically allocated OIF blocks, each slot of the set of slots having stored therein a second pointer to a transmit control block (TCB) data structure which services one or more users participating in the IP multicast session and which has stored therein control information to process or route packets of the IP multicast session, including information regarding an OIF of the router through which the packets are to be transmitted;

a third pointer associated with each OIF block of the set of dynamically allocated OIF blocks that links together the set of dynamically allocated OIF blocks and updates only one of OIF block of the OIF blocks responsive to users joining or leaving the IP multicast session.

7. The method of claim 6, wherein each OIF block of the set of OIF blocks is independently accessible without affecting processing of other OIF blocks of the set of OIF blocks.

8. The method of claim 6, wherein the TCB includes control information for processing replicated packets.

9. The method of claim 8, wherein the control information includes one or more of header transformation control information, metering control information and statistics control information.

10. The method of claim 6, wherein the second pointer associated with a first OIF block of the set of OIF blocks of a first  $\{S, G\}$  pair of the plurality of  $\{S, G\}$  pairs is permitted to point to and thus share a first TCB with another OIF block of the set of OIF blocks of the first  $\{S, G\}$  pair or with a second OIF block of the set of OIF blocks of a second  $\{S, G\}$  pair of the plurality of  $\{S, G\}$  pairs.

11. A method of managing multicast Internet Protocol (IP) sessions, the method comprising:

identifying, by a network device, active multicast IP sessions; and

maintaining, by the network device, a data structure within a memory of the network device containing therein information regarding the active multicast IP sessions; wherein the data structure includes:

a plurality of pairs of a source field and a group field ( $\{S, G\}$  pairs), in which each pair of the plurality of  $\{S, G\}$  pairs defines a multicast IP session of the active multicast IP sessions, wherein the source field defines a source of a multicast transmission of the multicast IP session and the group field defines a group corresponding to the multicast IP session;

a first value associated with each of the plurality of  $\{S, G\}$  pairs that is indicative of a dynamically allocated set of outbound interface (OIF) blocks, wherein a number OIF blocks in the dynamically allocated set of OIF blocks is dependent upon a number of OIFs of the network device that are participating in the IP multicast session and the number of OIF blocks in the dynamically allocated set of OIF blocks defines a number of times packets of the IP multicast session are to be replicated;

a set of slots for each OIF block of the set of dynamically allocated OIF blocks, each slot of the set of slots having stored therein a second value indicative of a transmit control block (TCB) data structure which services one or more users participating in the IP multicast session and which has stored therein control information to process or route packets of the IP multicast session, including information regarding an OIF of the network device through which the packets are to be transmitted;

a third value associated with each OIF block of the set of dynamically allocated OIF blocks that links together the set of dynamically allocated OIF blocks and update only one OIF block of the OIF blocks responsive to users joining or leaving the IP multicast session.

12. The method of claim 11, wherein each OIF block of the set of OIF blocks is independently accessible without affecting processing of other OIF blocks of the set of OIF blocks.

13. The method of claim 11, wherein the TCB includes control information for processing replicated packets.

14. The method of claim 13, wherein the control information includes one or more of header transformation control information, metering control information and statistics control information.

15. The method of claim 11, wherein the second value associated with a first OIF block of the set of OIF blocks of a first  $\{S, G\}$  pair of the plurality of  $\{S, G\}$  pairs comprises a pointer and is permitted to point to and thus share a first TCB with another OIF block of the set of OIF blocks of the first  $\{S, G\}$  pair or with a second OIF block of the set of OIF blocks of a second  $\{S, G\}$  pair of the plurality of  $\{S, G\}$  pairs.

\* \* \* \* \*