

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
12 August 2010 (12.08.2010)

(10) International Publication Number
WO 2010/090650 A2

- (51) International Patent Classification: Not classified
- (21) International Application Number: PCT/US2009/038904
- (22) International Filing Date: 31 March 2009 (31.03.2009)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data: 12/060,173 31 March 2008 (31.03.2008) US
- (71) Applicant (for all designated States except US): **AMAZON TECHNOLOGIES, INC.** [US/US]; PO Box 8102, Reno, NV 89507 (US).
- (72) Inventors; and
- (75) Inventors/Applicants (for US only): **RICHARDSON, David, R.** [US/US]; 1200 12th Avenue South, Suite 1200, Seattle, WA 98144 (US). **SIVASUBRAMANIAN, Swaminathan** [IN/US]; 1200 12th Avenue South, Suite 1200, Seattle, WA 98144 (US). **MARSHALL, Bradley, E.** [US/US]; 1200 12th Avenue South, Suite 1200, Seattle, WA 98144 (US). **SCOFIELD, Christopher, L.** [US/US]; 1200 12th Avenue South, Suite 1200, Seattle, WA 98144 (US). **POPE, Elmore Eugene** [US/US]; 1200 12th Avenue South, Suite 1200, Seattle, WA 98144 (US).

- (74) Common Representative: **AMAZON TECHNOLOGIES, INC.**; 8329 West Sunset Road, Suite 220, Las Vegas, NV 89113 (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, BR, BW, BY, BZ, CA, CH, 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, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, 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, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (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, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:
— without international search report and to be republished upon receipt of that report (Rule 48.2(g))

(54) Title: REQUEST ROUTING BASED ON CLASS

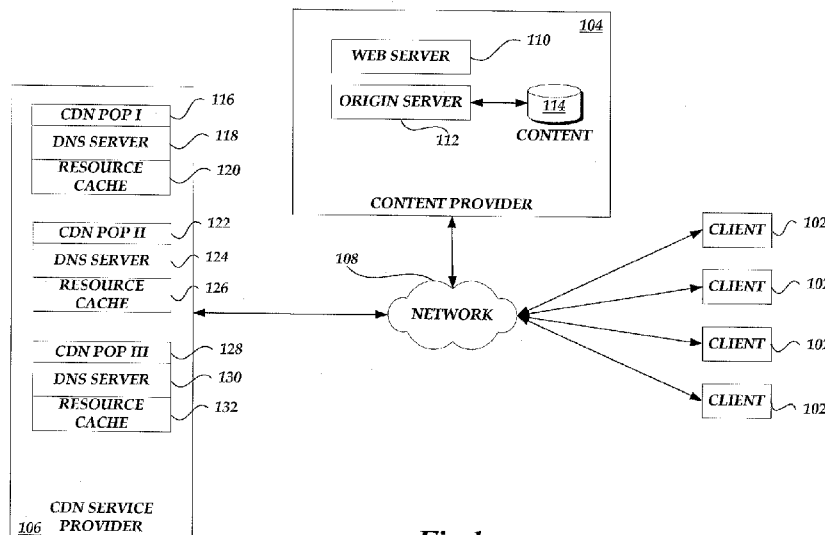


Fig. 1.

(57) **Abstract:** A system and method for management and processing of resource requests is provided. A content delivery network service provider receives a DNS query from a client computing device. The DNS query corresponds to a requested resource from the client computing device. The content delivery network service provider associates the client computing device with a cluster of other client computing devices. Based on routing information for the cluster, the content delivery network service provider routes the DNS query. The process can further include monitoring performance data associated with the delivery of the requested resource and updating the routing information for the cluster based on the performance data for use in processing subsequent requests from client computing devices in the cluster.

WO 2010/090650 A2

REQUEST ROUTING BASED ON CLASS

BACKGROUND

[0001] Generally described, computing devices and communication networks can be utilized to exchange information. In a common application, a computing device can request content from another computing device via the communication network. For example, a user at a personal computing device can utilize a software browser application to request a Web page from a server computing device via the Internet. In such embodiments, the user computing device can be referred to as a client computing device and the server computing device can be referred to as a content provider.

[0002] Content providers are generally motivated to provide requested content to client computing devices often with consideration of efficient transmission of the requested content to the client computing device and/or consideration of a cost associated with the transmission of the content. For larger scale implementations, a content provider may receive content requests from a high volume of client computing devices which can place a strain on the content provider's computing resources. Additionally, the content requested by the client computing devices may have a number of components, which can further place additional strain on the content provider's computing resources.

[0003] With reference to an illustrative example, a requested Web page, or original content, may be associated with a number of additional resources, such as images or videos, which are to be displayed with the Web page. In one specific embodiment, the additional resources of the Web page are identified by a number of embedded resource identifiers, such as uniform resource locators ("URLs"). In turn, software on the client computing devices typically processes embedded resource identifiers to generate requests for the content. Often, the resource identifiers associated with the embedded resources reference a computing device associated with the content provider such that the client computing device would transmit the request for the additional resources to the referenced content provider computing device. Accordingly, in order to satisfy a content request, the content provider would provide client computing devices data associated with the Web page as well as the data associated with the embedded resources.

[0004] Some content providers attempt to facilitate the delivery of requested content, such as Web pages and/or resources identified in Web pages, through the utilization of a content delivery network ("CDN") service provider. A CDN server provider typically maintains a number of computing devices in a communication network

that can maintain content from various content providers. In turn, content providers can instruct, or otherwise suggest to, client computing devices to request some, or all, of the content provider's content from the CDN service provider's computing devices.

[0005] As with content providers, CDN service providers are also generally motivated to provide requested content to client computing devices often with consideration of efficient transmission of the requested content to the client computing device and/or consideration of a cost associated with the transmission of the content. Accordingly, CDN service providers often consider factors such as latency of delivery of requested content in order to meet service level agreements or to generally improve the quality of delivery service.

DESCRIPTION OF THE DRAWINGS

[0006] The foregoing aspects and many of the attendant advantages of this invention will become more readily appreciated as the same become better understood by reference to the following detailed description, when taken in conjunction with the accompanying drawings, wherein:

[0007] FIGURE 1 is a block diagram illustrative of content delivery environment including a number of client computing devices, content provider, and a content delivery network service provider;

[0008] FIGURE 2 is a block diagram of the content delivery environment of FIGURE 1 illustrating the registration of a content provider with a content delivery service provider;

[0009] FIGURE 3 is a block diagram of the content delivery environment of FIGURE 1 illustrating the generation and processing of a content request from a client computing device to a content provider;

[0010] FIGURE 4 is a block diagram of the content delivery environment of FIGURE 1 illustrating one embodiment of the generation and processing of a DNS query corresponding to an embedded resource from a client computing device to a content delivery network service provider;

[0011] FIGURES 5A-5C are block diagrams of the content delivery environment of FIGURE 1 illustrating another embodiment of the generation and processing of a DNS query corresponding to an embedded resource from a client computing device to a content delivery network service provider and the subsequent

generation and processing of DNS queries corresponding to a first and a second alternative resource identifier from a client computing device to a content delivery network;

[0012] FIGURE 6 is a block diagram of the content delivery environment of FIGURE 1 illustrating the generation and processing of embedded resource requests from a client computing device to a content delivery network service provider;

[0013] FIGURE 7 is a flow diagram illustrative of a request routing routine implemented by a content delivery network service provider for selecting a cache server component; and

[0014] FIGURE 8 is a flow diagram illustrative a request routing routine implemented by a content delivery network service provider for updating routing information.

DETAILED DESCRIPTION

[0015] Generally described, the present disclosure is directed to the management and processing of resource requests made to a content delivery network (“CDN”) service provider from client computing devices. Specifically, aspects of the disclosure will be described with regard to routing information associated with a resource request based on routing information for a cluster of client computing devices. Although various aspects of the disclosure will be described with regard to illustrative examples and embodiments, one skilled in the art will appreciate that the disclosed embodiments and examples should not be construed as limiting.

[0016] FIGURE 1 is a block diagram illustrative of content delivery environment 100 for the management and processing of content requests. As illustrated in FIGURE 1, the content delivery environment 100 includes a number of client computing devices 102 (generally referred to as clients) for requesting content from a content provider and/or a CDN service provider. In an illustrative embodiment, the client computing devices 102 can correspond to a wide variety of computing devices including personal computing devices, laptop computing devices, hand-held computing devices, terminal computing devices, mobile devices, wireless devices, various electronic devices and appliances and the like. In an illustrative embodiment, the client computing devices 102 include necessary hardware and software components for establishing communications over a communication network 108, such as a wide area network or local

area network. For example, the client computing devices 102 may be equipped with networking equipment and browser software applications that facilitate communications via the Internet or an intranet.

[0017] Although not illustrated in FIGURE 1, each client computing device 102 utilizes some type of local DNS resolver component, such as a DNS Name server, that generates the DNS queries attributed to the client computing device. In one embodiment, the local DNS resolver component may be provided by an enterprise network to which the client computing device 102 belongs. In another embodiment, the local DNS resolver component may be provided by an Internet Service Provider (ISP) that provides the communication network connection to the client computing device 102.

[0018] The content delivery environment 100 can also include a content provider 104 in communication with the one or more client computing devices 102 via the communication network 108. The content provider 104 illustrated in FIGURE 1 corresponds to a logical association of one or more computing devices associated with a content provider. Specifically, the content provider 104 can include a web server component 110 corresponding to one or more server computing devices for obtaining and processing requests for content (such as Web pages) from the client computing devices 102. The content provider 104 can further include an origin server component 112 and associated storage component 114 corresponding to one or more computing devices for obtaining and processing requests for network resources from the CDN service provider. One skilled in the relevant art will appreciate that the content provider 104 can be associated with various additional computing resources, such as additional computing devices for administration of content and resources, DNS name servers, and the like. For example, although not illustrated in FIGURE 1, the content provider 104 can be associated with one or more DNS name server components that would be authoritative to resolve client computing device DNS queries corresponding to a domain of the content provider.

[0019] With continued reference to FIGURE 1, the content delivery environment 100 can further include a CDN service provider 106 in communication with the one or more client computing devices 102 and the content providers 104 via the communication network 108. The CDN service provider 106 illustrated in FIGURE 1 corresponds to a logical association of one or more computing devices associated with a CDN service provider. Specifically, the CDN service provider 106 can include a number

of Point of Presence (“POP”) locations 116, 122, 128 that correspond to nodes on the communication network 108. Each POP 116, 122, 128 includes a DNS component 118, 124, 130 made up of a number of DNS server computing devices for resolving DNS queries from the client computers 102. Each POP 116, 122, 128 also includes a resource cache component 120, 126, 132 made up of a number of cache server computing devices for storing resources from content providers and transmitting various requested resources to various client computers. The DNS components 118, 124 and 130 and the resource cache components 120, 126 132 may further include additional software and/or hardware components that facilitate communications including, but not limited, load balancing or load sharing software/hardware components.

[0020] In an illustrative embodiment, the DNS component 118, 124, 130 and resource cache component 120, 126, 132 are considered to be logically grouped, regardless of whether the components, or portions of the components, are physically separate. Additionally, although the POPs 116, 122, 128 are illustrated in FIGURE 1 as logically associated with the CDN service provider 106, the POPs will be geographically distributed throughout the communication network 108 in a manner to best serve various demographics of client computing devices 102. Additionally, one skilled in the relevant art will appreciate that the CDN service provider 106 can be associated with various additional computing resources, such additional computing devices for administration of content and resources, and the like.

[0021] One skilled in the relevant art will appreciate that the components and configurations provided in FIGURE 1 are illustrative in nature. Accordingly, additional or alternative components and/or configurations, especially regarding the additional components, systems and subsystems for facilitating communications may be utilized.

[0022] With reference now to FIGURES 2-6, the interaction between various components of the content delivery environment 100 of FIGURE 1 will be illustrated. For purposes of the example, however, the illustration has been simplified such that many of the components utilized to facilitate communications are not shown. One skilled in the relevant art will appreciate that such components can be utilized and that additional interactions would accordingly occur without departing from the spirit and scope of the present disclosure.

[0023] With reference to FIGURE 2, an illustrative interaction for registration of a content provider 104 with the CDN service provider 106 will be described. As

illustrated in FIGURE 2, the CDN content registration process begins with registration of the content provider 104 with the CDN service provider 106. In an illustrative embodiment, the content provider 104 utilizes a registration application program interface ("API") to register with the CDN service provider 106 such that the CDN service provider 106 can provide content on behalf of the content provider 104. The registration API includes the identification of the origin server 112 of the content provider 104 that will provide requested resources to the CDN service provider 106.

[0024] One skilled in the relevant art will appreciate that upon identification of appropriate origin servers 112, the content provider 104 can begin to direct requests for content from client computing devices 102 to the CDN service provider 106. Specifically, in accordance with DNS routing principles, a client computing device request corresponding to a resource identifier would eventually be directed toward a POP 116, 122, 128 associated with the CDN service provider 106. In the event that the resource cache component 120, 126, 132 of a selected POP does not have a copy of a resource requested by a client computing device 102, the resource cache component will request the resource from the origin server 112 previously registered by the content provider 104.

[0025] With continued reference to FIGURE 2, upon receiving the registration API, the CDN service provider 106 obtains and processes the registration information. In an illustrative embodiment, the CDN service provider 106 can then generate additional information that will be used by the client computing devices 102 as part of the content requests. The additional information can include, without limitation, client identifiers, such as client identification codes, content provider identifiers, such as content provider identification codes, executable code for processing resource identifiers, such as script-based instructions, and the like. One skilled in the relevant art will appreciate that various types of additional information may be generated by the CDN service provider 106 and that the additional information may be embodied in any one of a variety of formats.

[0026] The CDN service provider 106 returns an identification of applicable domains for the CDN service provider (unless it has been previously provided) and any additional information to the content provider 104. In turn, the content provider 104 can then process the stored content with content provider specific information. In one example, as illustrated in FIGURE 2, the content provider 104 translates resource

identifiers originally directed toward a domain of the origin server 112 to a domain corresponding to the CDN service provider. The translated URLs are embedded into requested content in a manner such that DNS queries for the translated URLs will resolve to a DNS server corresponding to the CDN service provider 106 and not a DNS server corresponding to the content provider 104. Although the translation process is illustrated in FIGURE 2, in some embodiments, the translation process may be omitted in a manner described in greater detail below.

[0027] Generally, the identification of the resources originally directed to the content provider 104 will be in the form of a resource identifier that can be processed by the client computing device 102, such as through a browser software application. In an illustrative embodiment, the resource identifiers can be in the form of a uniform resource locator (“URL”). Because the resource identifiers are included in the requested content directed to the content **provider**, the resource identifiers can be referred to generally as the “content provider URL.” For purposes of an illustrative example, the content provider URL can identify a domain of the content provider 104 (e.g., contentprovider.com), a name of the resource to be requested (e.g., “resource.xxx”) and a path where the resource will be found (e.g., “path”). In this illustrative example, the content provider URL has the form of:

`http://www.contentprovider.com/path/resource.xxx`

[0028] During an illustrative translation process, the content provider URL is modified such that requests for the resources associated with the translated URLs resolve to a POP associated with the CDN service provider 106. In one embodiment, the translated URL identifies the domain of the CDN service provider 106 (e.g., “cdnprovider.com”), the same name of the resource to be requested (e.g., “resource.xxx”) and the same path where the resource will be found (e.g., “path”). Additionally, the translated URL can include additional processing information (e.g., “additional information”). The translated URL would have the form of:

`http://additional information.cdnprovider.com/path/resources.xxx`

of a DNS server authoritative to the “.cdnprovider” portion of the embedded URL. The issuance of DNS queries corresponding to the “.” and the “.com” portions of a URL are well known and have not been illustrated.

[0032] With reference now to FIGURE 4, in an illustrative embodiment, the successful resolution of the “.cdnprovider” portion of the original URL identifies a network address, such as an IP address, of a DNS server associated with the CDN service provider 106. In one embodiment, the IP address can be a specific network address unique to a DNS server component of a POP. In another embodiment, the IP address can be shared by one or more POPs. In this embodiment, a further DNS query to the shared IP address utilizes a one-to-many network routing schema, such as anycast, such that a specific POP will receive the request as a function of network topology. For example, in an anycast implementation, a DNS query issued by a client computing device 102 to a shared IP address will arrive at a DNS server component logically having the shortest network topology distance, often referred to as network hops, from the client computing device. The network topology distance does not necessarily correspond to geographic distance. However, in some embodiments, the network topology distance can be inferred to be the shortest network distance between a client computing device 102 and a POP.

[0033] With continued reference to FIGURE 4, in either of the above identified embodiments (or any other embodiment), a specific DNS server in the DNS component 118 of a POP 116 receives the DNS query corresponding to the original URL from the client computing device 102. Once one of the DNS servers in the DNS component 118 receives the request, the specific DNS server attempts to resolve the request. In one illustrative embodiment as shown in FIGURE 4, a specific DNS server resolves the DNS query by identifying an IP address of a cache server component that will process the request for the requested resource. As described above and as will be described further below in reference to FIGURE 6, a selected resource cache component can process the request by either providing the requested resource if it is available or attempt to obtain the requested resource from another source, such as a peer cache server computing device or the origin server 112 of the content provider 104.

[0034] In further reference to FIGURE 4, the specific DNS server can utilize a variety of information in selecting a resource cache component. In an illustrative embodiment, and as will be further described below in reference to FIGURES 7 and 8, the DNS server determines a class associated with the requesting client computing device.

For example, the class can correspond to a specific geographic region to which the client computing device belongs or an internet service provider for the client computing device. Such class information can be determined from the client directly (such as information provided by the client computing device or ISP) or indirectly (such as inferred through a client computing device's IP address). Based on the class, the DNS server determines appropriate routing information. Then, for embodiments described specifically in reference to FIGURE 4, the DNS server selects an appropriate resource cache component for providing content associated with the resource request based on the routing information for the determined class of the client computing device. The IP address selected by the DNS server may correspond to a specific caching server in the resource cache. Alternatively, the IP address can correspond to a hardware/software selection component (such as a load balancer). As will also be further described below, the DNS server can further utilize network performance measurements to assist in selecting specific resource cache components for the determined class.

[0035] With reference now to FIGURES 5A-5C, as an alternative to selecting a resource cache component upon receipt of a DNS query as described in reference to FIGURE 4, the CDN service provider 106 can maintain sets of various alternative resource identifiers. The alternative resource identifiers can be provided by the CDN service provider 106 to the client computing device 102 such that a subsequent DNS query on the alternative resource identifier will resolve to a different DNS server component within the CDN service provider's network. In an illustrative embodiment, the alternative resource identifiers are in the form of one or more canonical name ("CNAME") records. In one embodiment, each CNAME record identifies a domain of the CDN service provider 106 (e.g., "cdnprovider.com" or "cdnprovider-1.com"). As will be explained in greater detail below, the domain in the CNAME does not need to be the same domain found in original URL or in a previous CNAME record. Additionally, each CNAME record includes additional information, such as request routing information, (e.g., "request routing information"). An illustrative CNAME record can have the form of:

http://request_routing_information.cdnprovider.com

[0036] In an illustrative embodiment, the CNAME records are generated and provided by the DNS servers to identify a more appropriate DNS server of the CDN

service provider 106. As with selecting an appropriate resource cache component as described above in reference to FIGURE 4, the DNS server receiving the initial DNS query can utilize a variety of information to select a more appropriate DNS server of the CDN service provider 106 to resolve the resource request. In an illustrative embodiment, and as will also be further described below in reference to FIGURES 7 and 8, the DNS server determines a class associated with the requesting client computing device. Again, the class can, for example, correspond to a specific geographic region to which the client computing device belongs or an internet service provider for the client computing device. In any case, the DNS server may obtain class information from the client directly (such as information provided by the client computing device or ISP) or indirectly (such as inferred through a client computing device's IP address). Based on the class, the DNS server determines appropriate routing information. Then, for the embodiments described specifically in reference to FIGURE 5A, the DNS server selects an appropriate alternative DNS server for use in resolving the resource request based on the routing information for the determined class of the client computing device. As will also be further described below, the DNS server can further utilize network performance measurements to select specific alternative DNS servers for the determined class.

[0037] In accordance with an illustrative embodiment, the DNS server maintains a data store that defines CNAME records for various original URLs. If a DNS query corresponding to a particular original URL matches an entry in the data store, the DNS server returns a CNAME record as defined in the data store. In an illustrative embodiment, the data store can include multiple CNAME records corresponding to a particular original URL. The multiple CNAME records would define a set of potential candidates that can be returned to the client computing device. In such an embodiment, the DNS server, either directly or via a network-based service, selects one of the CNAME records defined in the data store as more appropriate routing information based on logic that factors a determined class of the requesting client computing device. It will be appreciated by one skilled in the art and others that the DNS server can implement further additional logic in selecting an appropriate CNAME from a set of possible of CNAMEs. In an illustrative embodiment, each DNS server component 118, 124, 130 maintains the same data stores that define CNAME records, which can be managed centrally by the CDN service provider 106. Alternatively, each DNS server component 118, 124, 130 can have POP specific data stores that define CNAME records, which can be managed

centrally by the CDN service provider 106 or locally at the POP 116, 122, 128. Still further, each DNS server computing device within the DNS server components 118, 124, 130 can utilize shared data stores managed by a respective POP or a local data store specific to an individual DNS server computing device.

[0038] The returned CNAME can also include request routing information that is different from or in addition to the information provided in URL/CNAME of the current DNS query. For example, if the CNAME selection is based on a class associated with the requesting client computing device, a specific class can be identified in the “request_routing_information” portion of the specific CNAME record. A similar approach could be taken to identify service level plans and file management by including a specific identifier in the “request_routing_information” portion of the CNAME record. In another embodiment, request routing information can be found in the identification of a CDN service provider 106 domain different from the domain found in the current URL/CNAME. For example, if the CNAME is based on a regional plan, a specific regional plan domain (e.g., “cdnprovider-region1.com”) could be used in the domain name portion of the specific CNAME record. Any additional request routing information can be prepended to the existing request routing information in the current URL/CNAME such that the previous request routing information would not be lost (e.g., <http://serviceplan.regionalplan.cdnprovider.com>). One skilled in the relevant art will appreciate that additional or alternative techniques and/or combination of techniques may be used to include the additional request routing information in the CNAME record that is selected by the DNS server component 118.

[0039] With continued reference to FIGURE 5A, one skilled in the relevant art will appreciate that the DNS server may select (or otherwise obtain) a CNAME record that is intended to resolve to a more appropriate DNS server of the CDN service provider 106. It may be possible, however, that the same DNS server would also be authoritative for the subsequent DNS query for the CNAME to be provided to the client computing device. For example, a specific DNS server may be authoritative for both a specific regional plan and a service level plan. Thus, returning a CNAME would still result in the DNS query arriving at the same DNS query (which may also be due in part to the client computing device’s geography). In such an embodiment, the DNS server, such as DNS server component 118, may choose to resolve the future DNS query in advance

[0040] With reference now to FIGURE 5B, upon receipt of the CNAME from the DNS server component 118, the client computing device 102 generates a subsequent DNS query corresponding to the CNAME. As previously discussed with regard to FIGURE 4, the DNS query process could first start with DNS queries for the “.” and “com” portions, followed by a query for the “cdnprovider” portion of the CNAME. To the extent, however, that the results of a previous DNS queries can be cached (and remain valid), the client computing device 102 can utilize the cached information and does not need to repeat the entire process. However, at some point, depending on whether the CNAME provided by DNS server component 118 (FIGURE 5A) and the previous URL/CNAME share common CDN service provider domains, the current CNAME DNS query resolves to a different POP provided by the CDN service provider 106. As illustrated in FIGURE 5B, the DNS server component 124 of POP 122 receives the current CNAME based on the different information in the current CNAME previously provided by the DNS server component 118. As previously described, the DNS server component 124 can then determine whether to resolve the DNS query on the CNAME with an IP address of a cache component that will process the content request or whether to provide another alternative resource identifier selected in the manners described above.

[0041] For purposes of illustration, assume that the DNS server component 118 determines that the DNS query corresponding to the current CNAME (provided by DNS server component 116) also corresponds to a CNAME record in its data store. In such an example, the DNS server component 124 would do any necessary processing to select a specific CNAME as described above and return the CNAME to the client computing device. With reference now to FIGURE 5C, the client computing device 102 would now transmit a second subsequent DNS query corresponding to the CNAME provided by DNS server component 124 (FIGURE 5B). In accordance with DNS query processes already described, the DNS query would illustratively be received by the DNS server component 130 of POP 128. Again, the DNS server component 130 can then determine whether to resolve the DNS query on the CNAME with an IP address of a cache component that will process the content request or whether to provide another alternative resource identifier selected in the manners described above. In this example, the DNS server component 130 returns an IP address.

[0042] With continued reference to FIGURE 5C, in an illustrative embodiment, the DNS server components, such as DNS server component 130, can

utilize a variety of information in selecting a resource cache component. In one example, the DNS server component can default to a selection of a resource cache component of the same POP. In another example, the DNS server components can select a resource cache component based on various load balancing or load sharing algorithms. Still further, the DNS server components can utilize network performance metrics or measurements to assign specific resource cache components. Yet further, the DNS server components can select a resource cache component based on routing information for a class of the requesting client computing device as described in reference to FIGURE 4. Again, the IP address selected by a DNS server component may correspond to a specific caching server in the resource cache. Alternatively, the IP address can correspond to a hardware/software selection component (such as a load balancer).

[0043] With reference now to FIGURE 6, in an illustrative example, assume that the DNS server component 130 has selected the resource cache component 132 of POP 128. Upon receipt of the IP address for the resource cache component 132, the client computing device 102 transmits requests for the requested content to the resource cache component 132. The resource cache component 132 processes the request in a manner described above and the requested content is transmitted to the client computing device 102.

[0044] With reference now to FIGURE 7, one embodiment of a routine 700 implemented by the CDN service provider 106 for processing a resource request will be described. One skilled in the relevant art will appreciate that actions/steps outlined for routine 700 may be implemented by one or many computing devices/components that are associated with the CDN service provider 106. Accordingly, routine 700 has been logically associated as being generally performed by the CDN service provider 106, and thus the following illustrative embodiments should not be construed as limiting.

[0045] Routine 700 can apply to embodiments described both in reference to FIGURE 4 and FIGURES 5A-5C. As such, routine 700 will first be described in reference to embodiments corresponding to selecting resource cache components at DNS servers based on routing information for a class of the requesting client computing device, as generally described in reference to FIGURE 4.

[0046] At block 702, one of the DNS server components 118, 124, 130 receives a DNS query corresponding to a resource identifier. As previously discussed, the resource identifier can be a URL that has been embedded in content requested by the

client computing device 102 and previously provided by the content provider 104. The DNS server determines a class of the requesting client associate with the DNS query at block 704. As mentioned above, the class can, for example, correspond to a specific geographic region to which the client computing device belongs or an internet service provider for the client computing device. Such class information can be determined from the client directly (such as information provided by the client computing device or ISP) or indirectly (such as inferred through a client computing device's IP address). In an illustrative embodiment, the determination of class at block 702 can specifically include associating the requesting client computing device to a cluster of other client computing devices based on a variety of criteria. Such criteria can include geographic region and internet service provider data, as mentioned above, in addition to routing path information, networking equipment, client sponsored service level agreements, content provider service level agreements, and the like.

[0047] At a decision block 706, a test is conducted to determine whether the current DNS server is authoritative to resolve the DNS query. In an illustrative embodiment, the DNS server can determine whether it is authoritative to resolve the DNS query if there are no CNAME records corresponding to the received resource identifier. In this illustrative embodiment, there are no CNAME records. Accordingly, the routine continues at block 708 where, in general, the current DNS server determines routing information for the determined class. Specifically, in an illustrative embodiment, the DNS server selects an appropriate resource cache component for providing content associated with the resource request based on routing information for the determined class of the client computing device. The DNS server then provides the IP address of the selected resource cache component to the client computing device.

[0048] In an illustrative embodiment, the routing information can be a list of resource cache components that can service the content request for a particular class of client computing devices. The DNS server can use a variety of logic to select a resource cache component from the list. In one embodiment, a probability of selecting each resource cache component on the list can be defined, and the DNS server selects a resource cache component based on these probabilities. Accordingly, in this illustrative embodiment, a DNS server will select a resource cache component on a frequency corresponding to the determined probabilities. For example, the DNS server will most frequently select the resource cache component with the highest probability of selection,

but can also, at times, select a resource cache component with a lower probability of selection. In this case, the probabilities correspond to anticipated performance of the selected computing device. As will be described further below, the CDN service provider 106 can monitor performance of delivering requested resources to clients in a particular class and thereafter update the routing information (e.g., probabilities) accordingly. In another embodiment, the probabilities can correspond to load shedding or other network traffic mitigation. By periodically selecting a non-preferred resource cache component and monitoring its performance for the class, the CDN service provider 106 can thus determine if changes to the routing information for the class are desirable.

[0049] It will be appreciated by one skilled in the relevant art that a number of algorithms or selection logic can be used for selecting a resource cache component to service the resource request from a particular class of client computing devices. For example, in addition to the frequency-based reinforcement algorithm described above, the DNS server may implemented alternative reinforcement learning algorithms. Examples of other reinforcement algorithms include, but are not limited to, algorithms such as State-Action-Reward-State-Action (SARSA), Q-learning, delayed Q-learning, and the like. Additionally, other machine learning approaches, such as support vector machines, neural networks, Bayesian engines, etc. may be utilized in conjunction with a DNS server to select the appropriate resource cache component.

[0050] Next, embodiments in which routing information for a class of the requesting client computing device is used to select an appropriate DNS server for processing the request will be described. In such embodiments, routine 700 similarly commences at block 702 where one of the DNS server components 118, 124, 130 receives a DNS query corresponding to a resource identifier. As described above, the DNS server further determines a class of the requesting client computing device associated with the DNS query at block 704.

[0051] At decision block 706, a test is conducted to determine whether the current DNS server is authoritative to resolve the DNS query. In an illustrative embodiment, the DNS server can determine whether it is authoritative to resolve the DNS query if there are no CNAME records corresponding to the received resource identifier. Alternative or additional methodologies may also be practiced to determine whether the DNS server is authoritative.

[0052] If the current DNS server is authoritative (including a determination that the same DNS server will be authoritative for subsequent DNS queries), the current DNS server resolves the DNS query by returning the IP address of cache server component. In a non-limiting manner, a number of methodologies for selecting an appropriate resource cache component have been previously discussed. Additionally, as described above, the IP address may correspond to a specific cache server of a resource cache component or generally to group of cache servers.

[0053] Alternatively, if at decision block 704 the DNS server is not authoritative, at block 708, the DNS server component selects and transmits an alternative resource identifier. As described above, the DNS server component can utilize a data store to identify a set of potential candidate CNAMEs as a function of the current DNS query. The DNS server then, either directly or via a network-based service, selects one of the CNAME records defined in the data store as more appropriate routing information based on logic that factors a determined class of the requesting client computing device. At block 710, different DNS server components 118, 124, 130 receive a DNS query corresponding to the CNAME. The routine 700 then returns to decision block 704 and continues to repeat as appropriate.

[0054] In an illustrative embodiment, where the DNS server is not authoritative, the routing information can be a set or list of potential candidate CNAMEs which correspond to one or more DNS servers that can service the content request for a particular class of client computing devices. Similar to selecting a cache resource component as described above, the DNS server can use a variety of logic to select a CNAME, or another DNS server, from the list. In one embodiment, a probability of selecting each CNAME in the set can be initially defined in a number of ways, and the DNS server selects a CNAME based on the defined probabilities. Accordingly, in this illustrative embodiment, a DNS server will most frequently select the CNAME with the highest probability of selection, but can also, at times, select a CNAME with a lower probability of selection. In this case, the probabilities correspond to anticipated performance of the corresponding computing device. As will be described further below, the CDN service provider 106 can monitor performance of delivering requested resources to clients in a particular class and thereafter update the probabilities. Again, in further embodiments, the probabilities can correspond to load shedding or other network traffic mitigation. By periodically selecting a non-preferred CNAME and monitoring

performance of the corresponding DNS server for the class, the CDN service provider 106 can thus determine if changes to the routing information for the class are desirable. It will be appreciated by one skilled in the relevant art that a number of algorithms or selection logic can be used for selecting a CNAME/DNS server to service the resource request from a particular class of client computing devices.

[0055] With reference now to FIGURE 8, one embodiment of a request routing routine 800 for updating routing information will be described. One skilled in the relevant art will appreciate that actions/steps outlined for routine 800 may be implemented by one or many computing devices/components that are associated with the CDN service provider 106. Accordingly, routine 800 has been logically associated as being performed by the CDN service provider 106.

[0056] At a block 802, a first DNS server of the CDN service provider 106 receives a DNS query corresponding to a requested resource from a client computing device. As similarly described above in reference to block 704 of FIGURE 7, the DNS server at block 804 determines a class corresponding to the requesting client and associated with the DNS query. Also at block 804, the DNS server determines either DNS or cache routing information based on the determined class of the client computing device as similarly described above. The routine 800 continues at block 806 where network performance criteria associated with delivery of the requested resource is monitored. The network performance criteria can correspond to measurements of network performance for transmitting data from the CDN service provider POPs to the client computing device 102. In one embodiment, network data transfer latencies associated with the delivery of the requested resource are measured by the client computing device 102. Alternatively, the CDN service provider 106, such as through the resource cache component, can measure the performance as part of providing content to a client computing device. Such network performance data can be managed and maintained globally by the CDN service provider and shared with the DNS servers of the CDN or individually by the DNS servers of the CDN service provider. Moreover, network performance criteria can be provided as a batch process from POPs or sent in response to a request from one POP to another.

[0057] With continued reference to FIGURE 8, at a test block 808, a determination is made as to whether an update to the routing information for the identified class is needed based on the performance data. In one embodiment, the update

determination can be made by the CDN service provider globally or by the individual DNS service components or DNS servers. In an illustrative embodiment where individual DNS servers determine whether to update routing information for a class, each DNS server can manage and maintain routing information for the identified class unique to the particular DNS server. In this illustrative embodiment, the performance data can be maintained globally by the CDN service provider and shared with the DNS components and/or DNS servers, with each DNS component and/or DNS server managing how the performance data is used. Accordingly, routing information for a class may vary from one DNS component/server to another.

[0058] Returning to FIGURE 8, if an update is needed, the routing information for the identified class is modified at block 810. In one embodiment, the CDN service provider 106 modifies a list of computing devices (e.g. DNS components/servers and/or resource cache components) for servicing a resource request from a particular class of client computing devices 102. In another embodiment, the CDN service provider and/or specific DNS components/servers can maintain and modify probabilities of selection of particular computing devices for servicing a resource request for a class of client computing devices. For example, if performance data indicates that a DNS server and/or a resource cache component which has a lower probability of selection has performed well, the probability of selection may be increased so that the particular DNS server and/or resource cache component will be selected more frequently for servicing a resource request from a client computing device. After a modification has been made at block 810, or if an update is not needed at block 808, the routine 800 returns to block 802 for further processing as described above.

[0059] It will be appreciated by one skilled in the relevant art that there are a number of ways to modify the routing information associated with requests from a class of client computing devices. It will further be appreciated by one skilled in the relevant art that the timing at which performance is monitored and updates to routing information are made can vary.

[0060] It will be appreciated by those skilled in the art and others that all of the functions described in this disclosure may be embodied in software executed by one or more processors of the disclosed components and mobile communication devices. The software may be persistently stored in any type of non-volatile storage.

[0061] Conditional language, such as, among others, "can," "could," "might," or "may," unless specifically stated otherwise, or otherwise understood within the context as used, is generally intended to convey that certain embodiments include, while other embodiments do not include, certain features, elements and/or steps. Thus, such conditional language is not generally intended to imply that features, elements and/or steps are in any way required for one or more embodiments or that one or more embodiments necessarily include logic for deciding, with or without user input or prompting, whether these features, elements and/or steps are included or are to be performed in any particular embodiment.

[0062] Any process descriptions, elements, or blocks in the flow diagrams described herein and/or depicted in the attached FIGURES should be understood as potentially representing modules, segments, or portions of code which include one or more executable instructions for implementing specific logical functions or steps in the process. Alternate implementations are included within the scope of the embodiments described herein in which elements or functions may be deleted, executed out of order from that shown or discussed, including substantially concurrently or in reverse order, depending on the functionality involved, as would be understood by those skilled in the art. It will further be appreciated that the data and/or components described above may be stored on a computer-readable medium and loaded into memory of the computing device using a drive mechanism associated with a computer readable storing the computer executable components such as a CD-ROM, DVD-ROM, or network interface further, the component and/or data can be included in a single device or distributed in any manner.

[0063] Clause 1. A method comprising:

obtaining a DNS query from a client computing device at a first DNS server, wherein the DNS query corresponds to a requested resource and wherein the first DNS server corresponds to a content delivery network service;

determining a class associated with the client computing device;

determining routing information based on the class of the client computing device;

monitoring performance associated with delivery of the requested resource;

determining whether an update to the routing information for the class is needed based on the delivery performance; and

if so, modifying the routing information for the class.

[0064] Clause 2. The method as recited in Clause 1, wherein modifying the routing information for the class comprises modifying a list of computing devices for servicing the resource request.

[0065] Clause 3. The method as recited in Clause 1, wherein modifying the routing information for the class comprises modifying probabilities of selection of particular computing devices used for servicing the resource request.

[0066] Clause 4. The method as recited in Clause 1, wherein the routing information for the determined class is associated with the first DNS server.

[0067] Clause 5. The method as recited in Clause 4, wherein the routing information for the determined class for the first DNS server is different from routing information for the determined class for a second DNS server of the content delivery network service provider.

[0068] Clause 6. The method as recited in Clause 1, wherein the class corresponds to an internet service provider.

[0069] Clause 7. The method as recited in Clause 1, wherein the class corresponds to a designated geographic location.

[0070] Clause 8. The method as recited in Clause 1, wherein determining routing information comprises selecting a cache component for providing content associated with the resource request based on routing information for the determined class of the client computing device.

[0071] Clause 9. The method as recited in Clause 1, wherein determining routing information comprises:

determining whether the first DNS server is authoritative based on routing information for the determined class of the client computing device; and

if not, selecting a second DNS server for processing the resource request based on the routing information for the determined class of the client computing device.

[0072] Clause 10. A method comprising:

obtaining a DNS query from a client computing device at a first DNS server, wherein the DNS query corresponds to a requested resource and wherein the first DNS server corresponds to a content delivery network service provider;

associating the client computing device with a cluster of other client computing devices; and

routing the DNS query based on routing information for the associated cluster of the client computing device.

[0073] Clause 11. The method as recited in Clause 10 further comprising: monitoring performance associated with delivery of the requested resource; determining whether an update to the routing information for the cluster is needed based on the delivery performance; and

if so, modifying the routing information for the cluster.

[0074] Clause 12. The method as recited in Clause 11, wherein modifying the routing information for the cluster comprises modifying a list of computing devices for servicing the resource request.

[0075] Clause 13. The method as recited in Clause 11, wherein modifying the routing information for the cluster comprises modifying probabilities of selection of particular computing devices used for servicing the resource request.

[0076] Clause 14. The method as recited in Clause 11, wherein the routing information for the cluster is associated with the first DNS server.

[0077] Clause 15. The method as recited in Clause 14, wherein the routing information for the cluster for the first DNS server is different from routing information for the cluster for a second DNS server of the content delivery network service provider.

[0078] Clause 16. The method as recited in Clause 10, wherein routing the DNS query comprises selecting a cache component for providing content associated with the resource request based on routing information for the cluster of the client computing device.

[0079] Clause 17. The method as recited in Clause 16, wherein selecting the cache component comprises selecting a cache component based on a frequency defined by probabilities of selection.

[0080] Clause 18. The method as recited in Clause 17 further comprising: monitoring performance associated with delivery of the requested resource to the cache component having the lower probability of selection; and

modifying the probability of selection of the cache component having the lower probability of selection for the cluster based on the delivery performance.

[0081] Clause 19. The method as recited in Clause 16, wherein selecting the cache component comprises periodically selecting a non-preferred cache component for providing content associated with the resource request for the cluster.

[0082] Clause 20. The method as recited in Clause 20 further comprising:
monitoring performance associated with delivery of the requested resource to the non-preferred cache component; and
modifying the routing information for the cluster based on the delivery performance of the requested resource to the non-preferred cache component.

[0083] Clause 21. The method of Clause 10, wherein routing the DNS query comprises:
determining whether the first DNS server is authoritative based on routing information for the cluster of the client computing device; and
if not, selecting a second DNS server for processing the resource request based on the routing information for the cluster of the client computing device.

[0084] Clause 22. The method as recited in Clause 22 further comprising:
monitoring performance associated with delivery of the requested resource using the second DNS server; and
modifying the routing information for the cluster based on the delivery performance of the requested resource using the second DNS server.

[0085] Clause 23. The method as recited in Clause 23, wherein modifying the routing information for the cluster comprises modifying a list of DNS servers for processing the resource request for the cluster.

[0086] Clause 24. The method as recited in Clause 23, wherein modifying the routing information for the cluster comprises modifying probabilities of selection of particular DNS servers used for processing the resource request for the cluster.

[0087] Clause 25. A system comprising:
a first network point of presence associated with a content delivery network service provider, wherein the first network point of presence includes a DNS server component that receives a DNS query from a client computing device, wherein the DNS query corresponds to a requested resource, and wherein the DNS server in the first network point of presence is operable to:

associate the client computing device with a cluster of other client computing devices; and
route the DNS query based on routing information for the cluster.

[0088] Clause 26. The system as recited in Clause 26, wherein the DNS server in the first network point of presence is further operable to:

monitor performance associated with delivery of the requested resource;
determine whether an update to the routing information for the cluster is
needed based on the delivery performance; and

if so, modify the routing information for the cluster.

[0089] Clause 27. The system as recited in Clause 26, wherein the routing information for the cluster determined at the DNS server in the first network point of presence is different from routing information for the cluster for a second DNS server in a second network point of presence associated with the content delivery network service provider.

[0090] Clause 28. The system as recited in Clause 26, wherein the DNS server in the first network point of presence is further operable to select a cache component for providing content associated with the resource request based on routing information for the cluster of the client computing device.

[0091] Clause 29. The system as recited in Clause 26, wherein the DNS server in the first network point of presence is further operable to:

determine whether the first DNS server is authoritative based on routing information for the cluster of the client computing device; and
if not, select a second DNS server in a second network point of presence for processing the resource request based on the routing information for the cluster of the client computing device.

[0092] Clause 30. A system comprising:

a first network point of presence associated with a content delivery network service provider, wherein the first network point of presence includes a DNS server component that receives a DNS query from a client computing device, wherein the DNS query corresponds to a requested resource associated with a first resource identifier, and wherein the DNS server component in the first network point of presence is operable to:

associate the client computing device with a cluster of other client computing devices;

determine that the first resource identifier is associated with an alternative resource identifier based on routing information for the cluster; and

transmit an alternative resource identifier to the client computing device, wherein the alternative resource identifier includes information for causing a DNS query to resolve to a domain corresponding to the content delivery service provider.

[0093] Clause 31. The system as recited in Clause 31, wherein the alternative resource identifier corresponds to a canonical name record identifier.

[0094] Clause 32. The system as recited in Clause 31 further comprising:
a second network point of presence associated with a content delivery network service provider, wherein the second network point of presence includes a DNS server component that receives a subsequent DNS query from the client computing device, wherein the subsequent DNS query corresponds to a requested resource associated with the alternative resource identifier, and wherein the DNS server component in the second network point of presence is operable to:

associate the client computing device with a cluster of other client computing devices; and

route the DNS query based on routing information for the cluster.

[0095] Clause 33. The system as recited in Clause 31 further comprising:
a second network point of presence associated with a content delivery network service provider, wherein the second network point of presence includes a DNS server component that receives a subsequent DNS query from a client computing device, wherein the subsequent DNS query corresponds to a requested resource associated with the alternative resource identifier, and wherein the DNS server component in the second network point of presence is operable to:

resolve the subsequent DNS query to identify a cache component for providing content associated with the original resource request; and

transmit information identifying the identified cache component to the client computing device.

[0096] Accordingly, general purpose computing devices may be configured to implement the processes, algorithms and methodology of the present disclosure with the processing and/or execution of the various data and/or components described above.

[0097] It should be emphasized that many variations and modifications may be made to the above-described embodiments, the elements of which are to be understood as being among other acceptable examples. All such modifications and variations are

intended to be included herein within the scope of this disclosure and protected by the following claims.

WHAT IS CLAIMED IS:

1. A method comprising:
 - obtaining a DNS query from a client computing device at a first DNS server, wherein the DNS query corresponds to a requested resource and wherein the first DNS server corresponds to a content delivery network service;
 - determining a class associated with the client computing device;
 - determining routing information based on the class of the client computing device;
 - monitoring performance associated with delivery of the requested resource;
 - determining whether an update to the routing information for the class is needed based on the delivery performance; and
 - if so, modifying the routing information for the class.
2. The method as recited in Claim 1, wherein modifying the routing information for the class comprises modifying a list of computing devices for servicing the resource request.
3. The method as recited in Claim 1, wherein modifying the routing information for the class comprises modifying probabilities of selection of particular computing devices used for servicing the resource request.
4. The method as recited in Claim 1, wherein the class corresponds to a designated geographic location.
5. The method as recited in Claim 1, wherein determining routing information comprises selecting a cache component for providing content associated with the resource request based on routing information for the determined class of the client computing device.
6. A method comprising:
 - obtaining a DNS query from a client computing device at a first DNS server, wherein the DNS query corresponds to a requested resource and wherein the first DNS server corresponds to a content delivery network service provider;
 - associating the client computing device with a cluster of other client computing devices; and
 - routing the DNS query based on routing information for the associated cluster of the client computing device.

7. The method as recited in Claim 6 further comprising:
monitoring performance associated with delivery of the requested resource;
determining whether an update to the routing information for the cluster is
needed based on the delivery performance; and
if so, modifying the routing information for the cluster.

8. The method as recited in Claim 6, wherein routing the DNS query comprises
selecting a cache component for providing content associated with the resource request
based on routing information for the cluster of the client computing device.

9. The method as recited in Claim 8, wherein selecting the cache component
comprises periodically selecting a non-preferred cache component for providing content
associated with the resource request for the cluster.

10. The method of Claim 6, wherein routing the DNS query comprises:
determining whether the first DNS server is authoritative based on routing
information for the cluster of the client computing device; and

if not, selecting a second DNS server for processing the resource request based
on the routing information for the cluster of the client computing device.

11. A system comprising:

a first network point of presence associated with a content delivery network
service provider, wherein the first network point of presence includes a DNS server
component that receives a DNS query from a client computing device, wherein the DNS
query corresponds to a requested resource, and wherein the DNS server in the first
network point of presence is operable to:

associate the client computing device with a cluster of other client computing
devices; and

route the DNS query based on routing information for the cluster.

12. The system as recited in Claim 11, wherein the DNS server in the first
network point of presence is further operable to:

monitor performance associated with delivery of the requested resource;

determine whether an update to the routing information for the cluster is needed
based on the delivery performance; and

if so, modify the routing information for the cluster.

13. A system comprising:

a first network point of presence associated with a content delivery network service provider, wherein the first network point of presence includes a DNS server component that receives a DNS query from a client computing device, wherein the DNS query corresponds to a requested resource associated with a first resource identifier, and wherein the DNS server component in the first network point of presence is operable to:

associate the client computing device with a cluster of other client computing devices;

determine that the first resource identifier is associated with an alternative resource identifier based on routing information for the cluster; and

transmit an alternative resource identifier to the client computing device, wherein the alternative resource identifier includes information for causing a DNS query to resolve to a domain corresponding to the content delivery service provider.

14. The system as recited in Claim 13, wherein the alternative resource identifier corresponds to a canonical name record identifier.

15. The system as recited in Claim 13 further comprising:

a second network point of presence associated with a content delivery network service provider, wherein the second network point of presence includes a DNS server component that receives a subsequent DNS query from the client computing device, wherein the subsequent DNS query corresponds to a requested resource associated with the alternative resource identifier, and wherein the DNS server component in the second network point of presence is operable to:

associate the client computing device with a cluster of other client computing devices; and

route the DNS query based on routing information for the cluster.

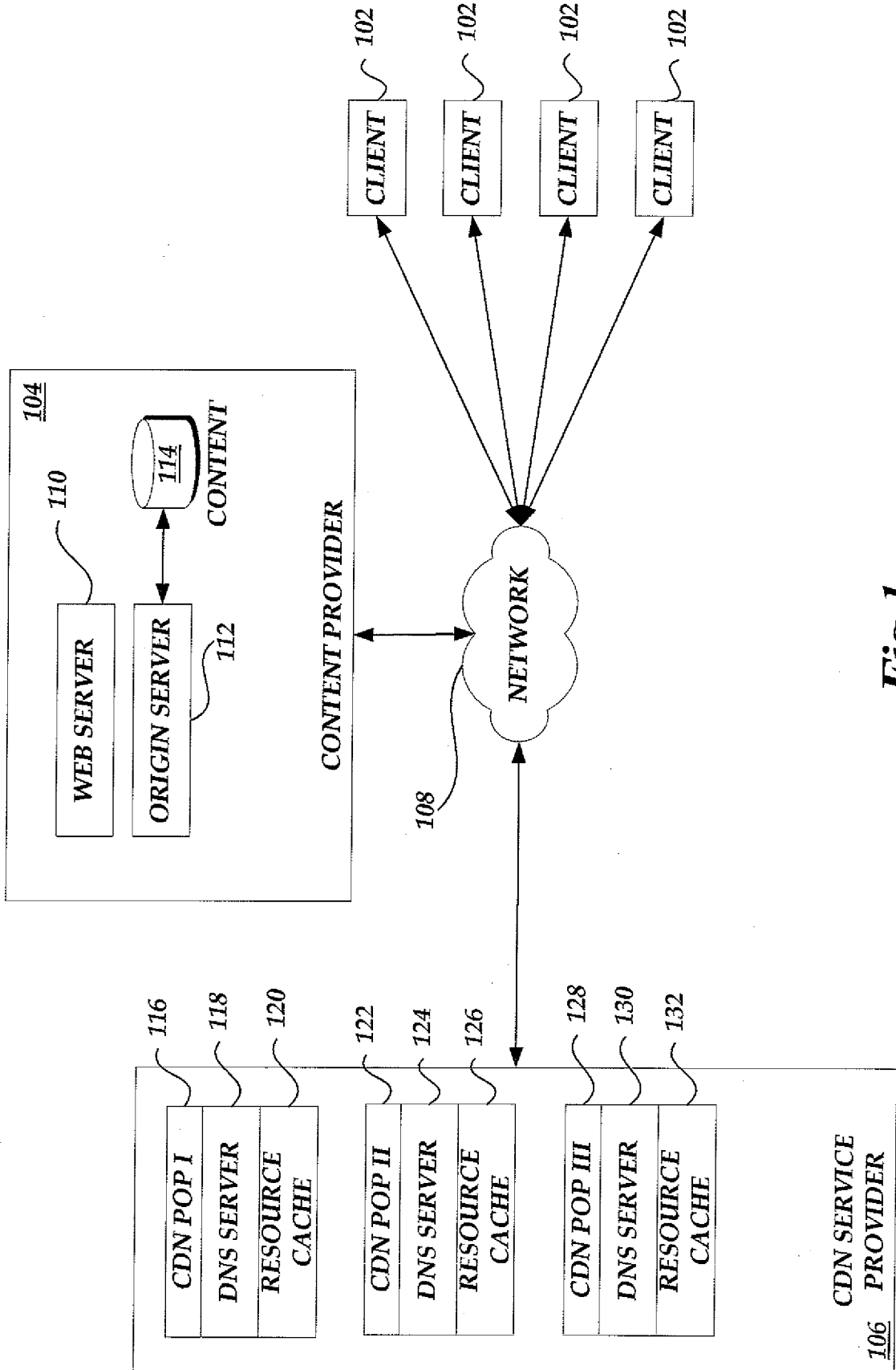


Fig. 1.

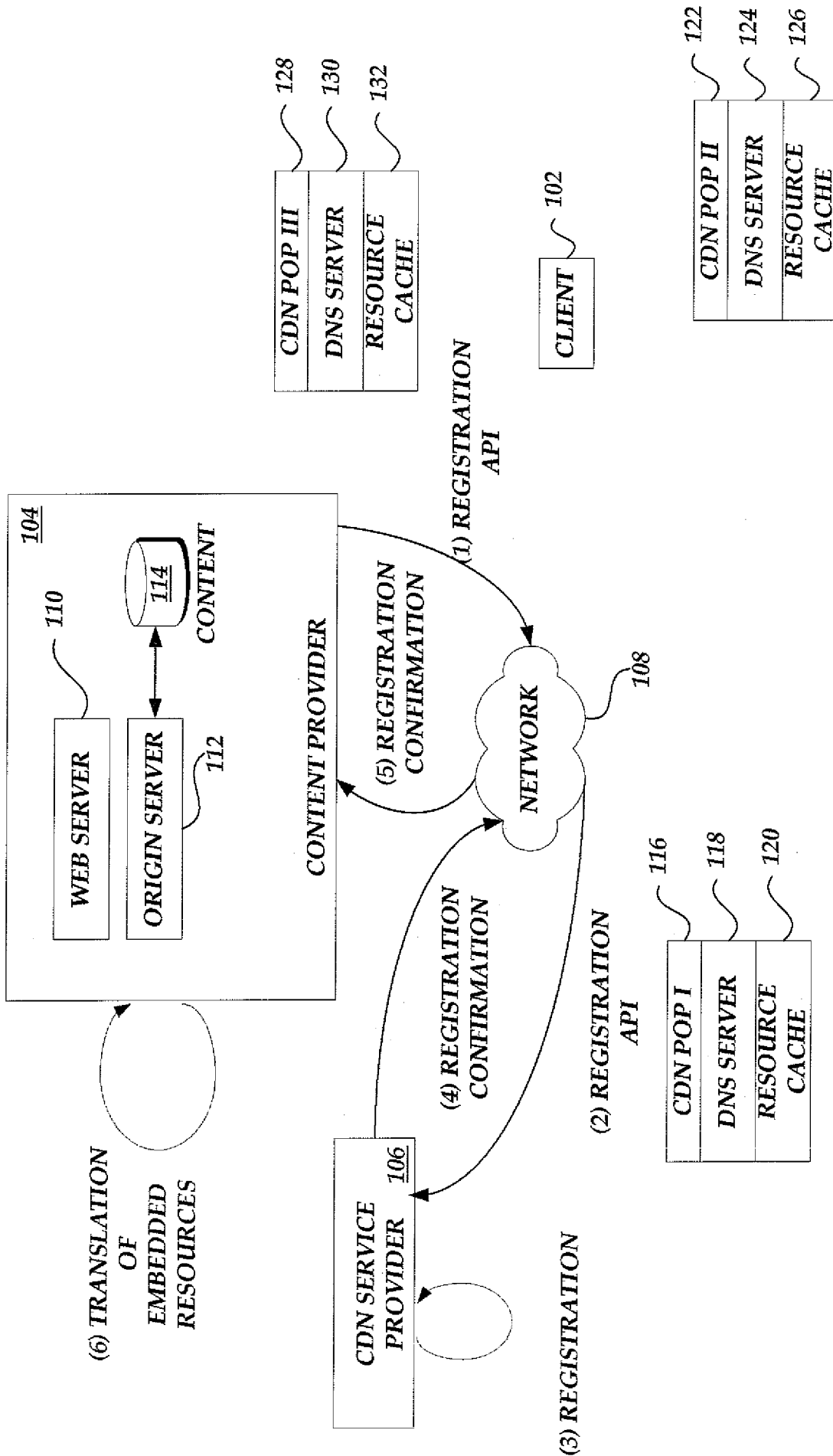


Fig. 2.

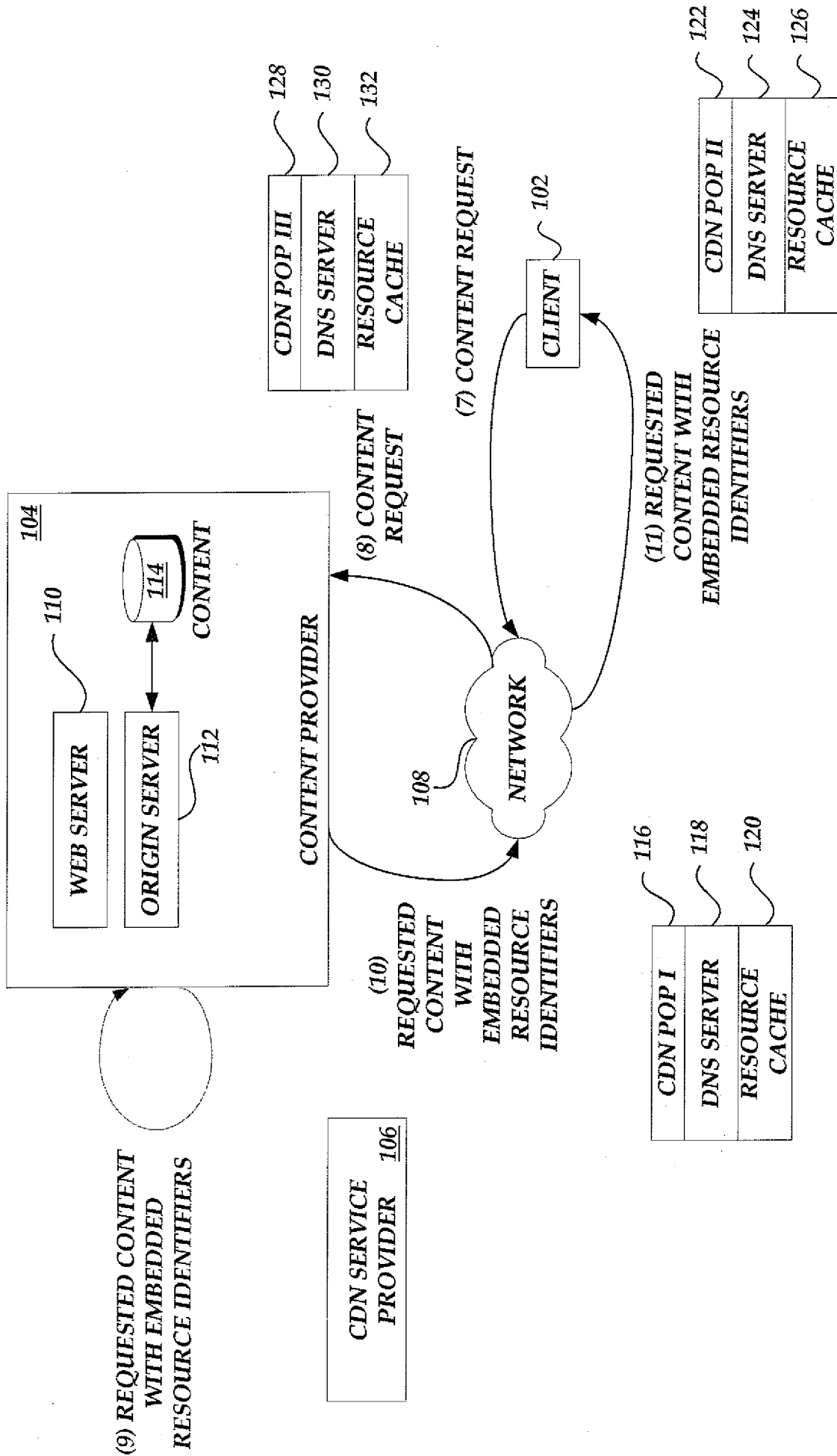


Fig. 3.

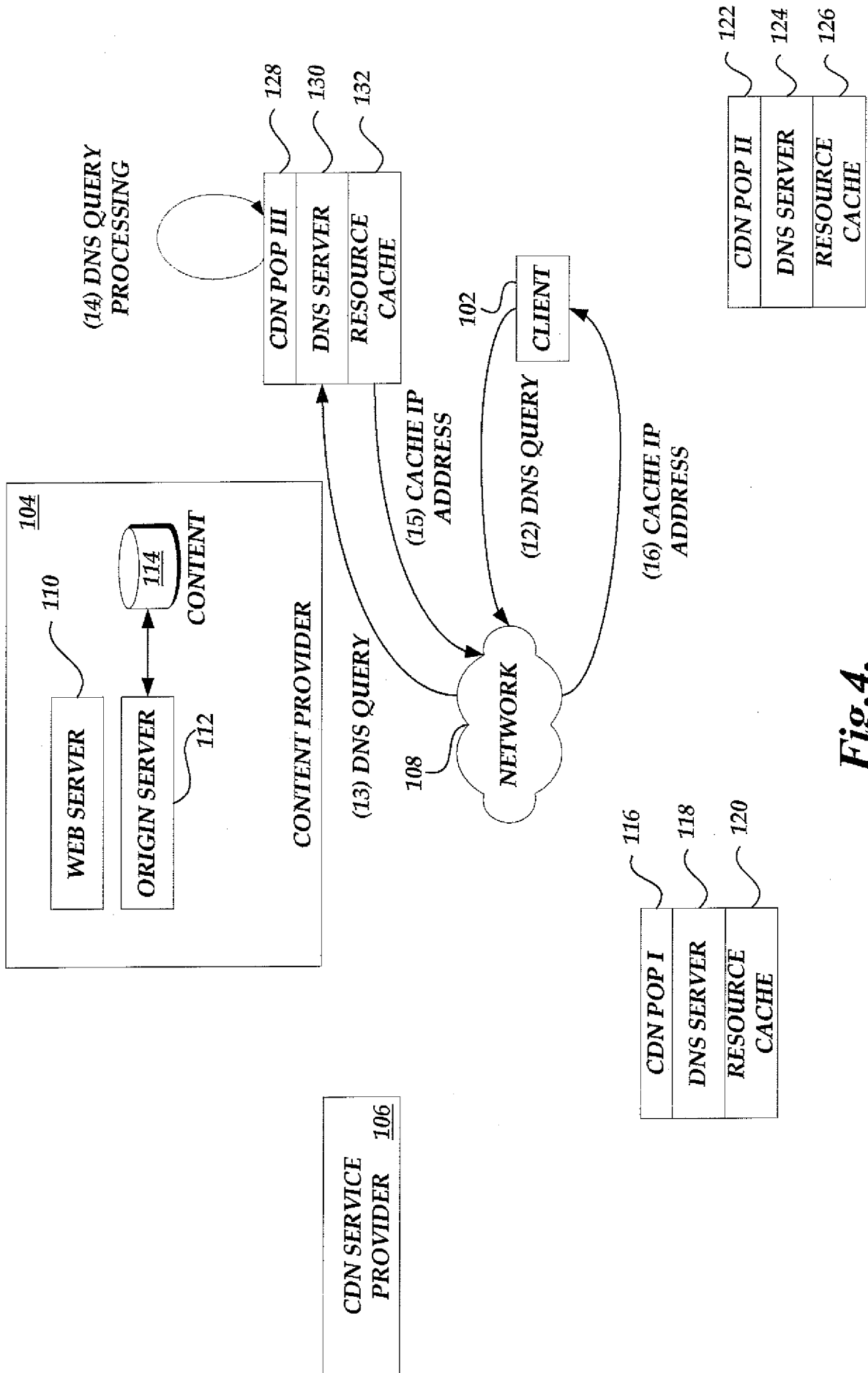


Fig.4.

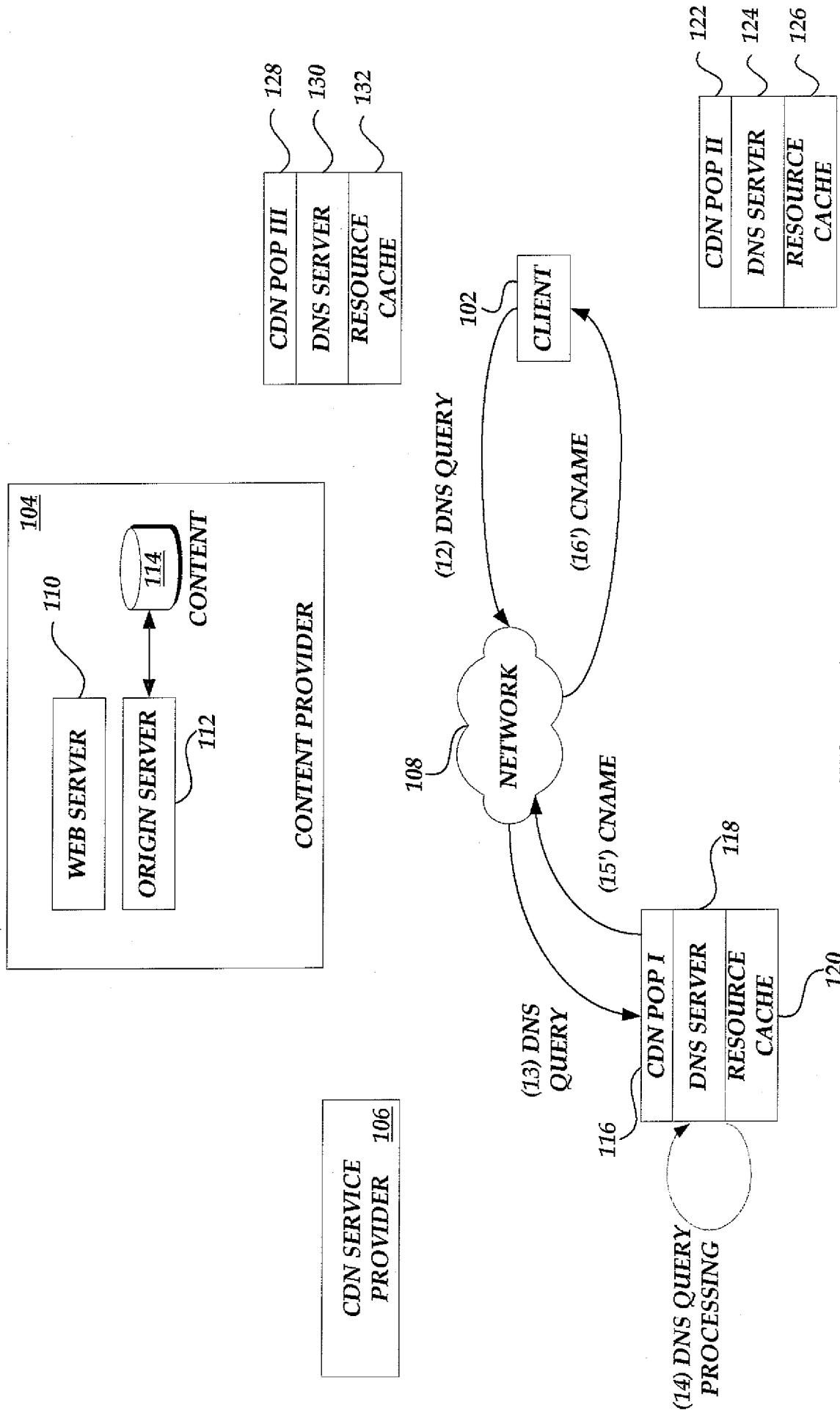


Fig.5A.

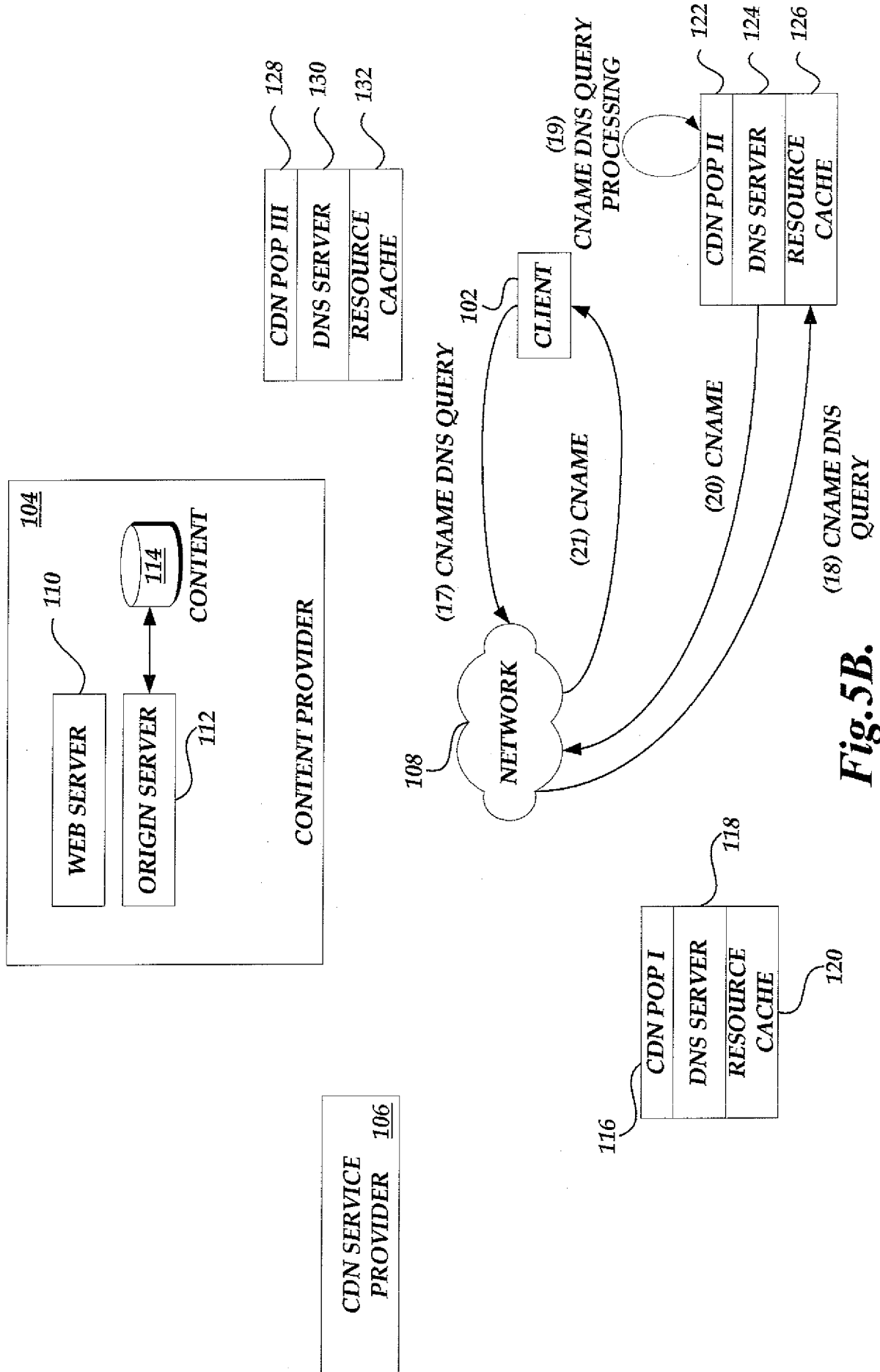


Fig. 5B.

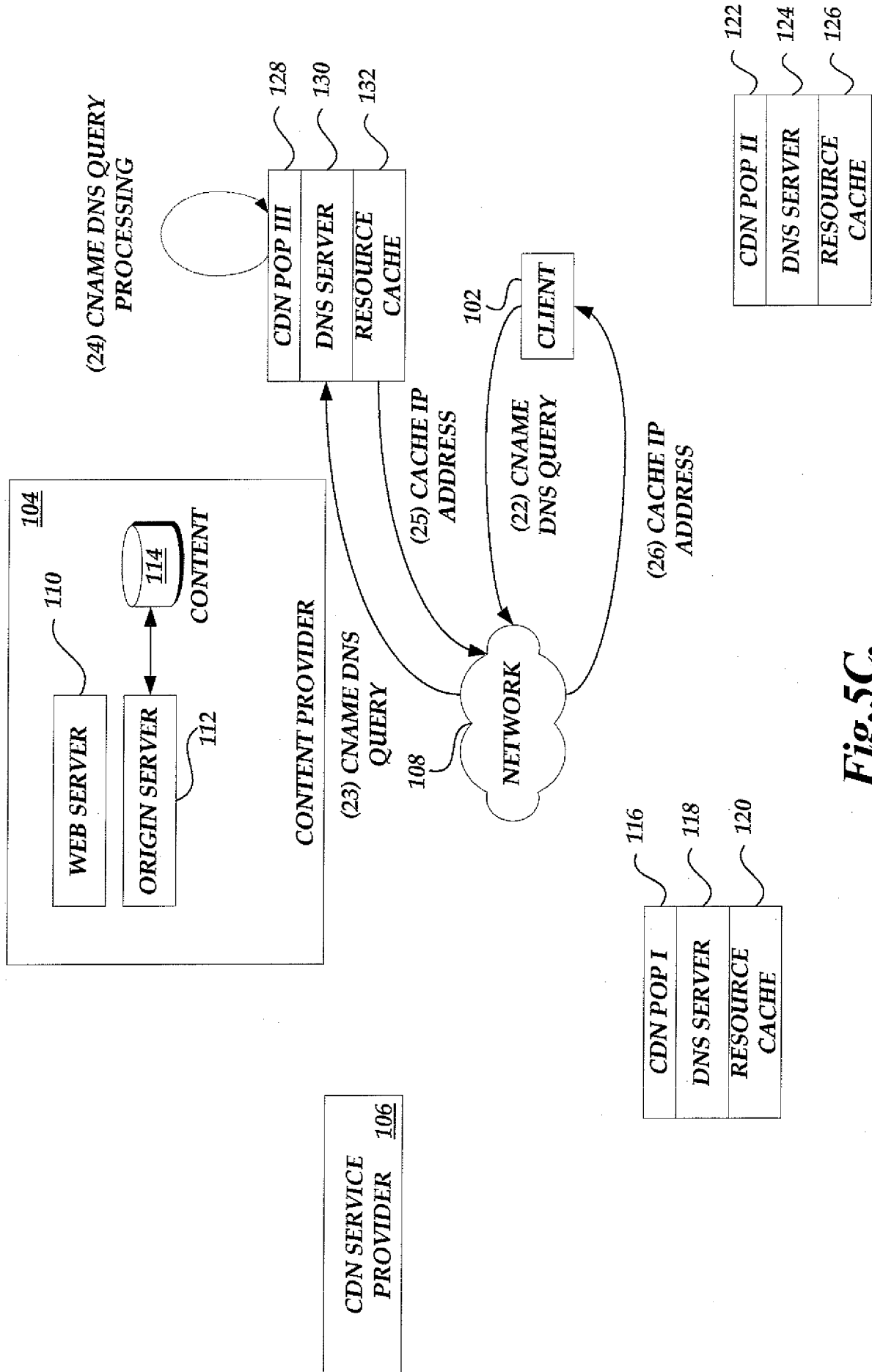


Fig. 5C.

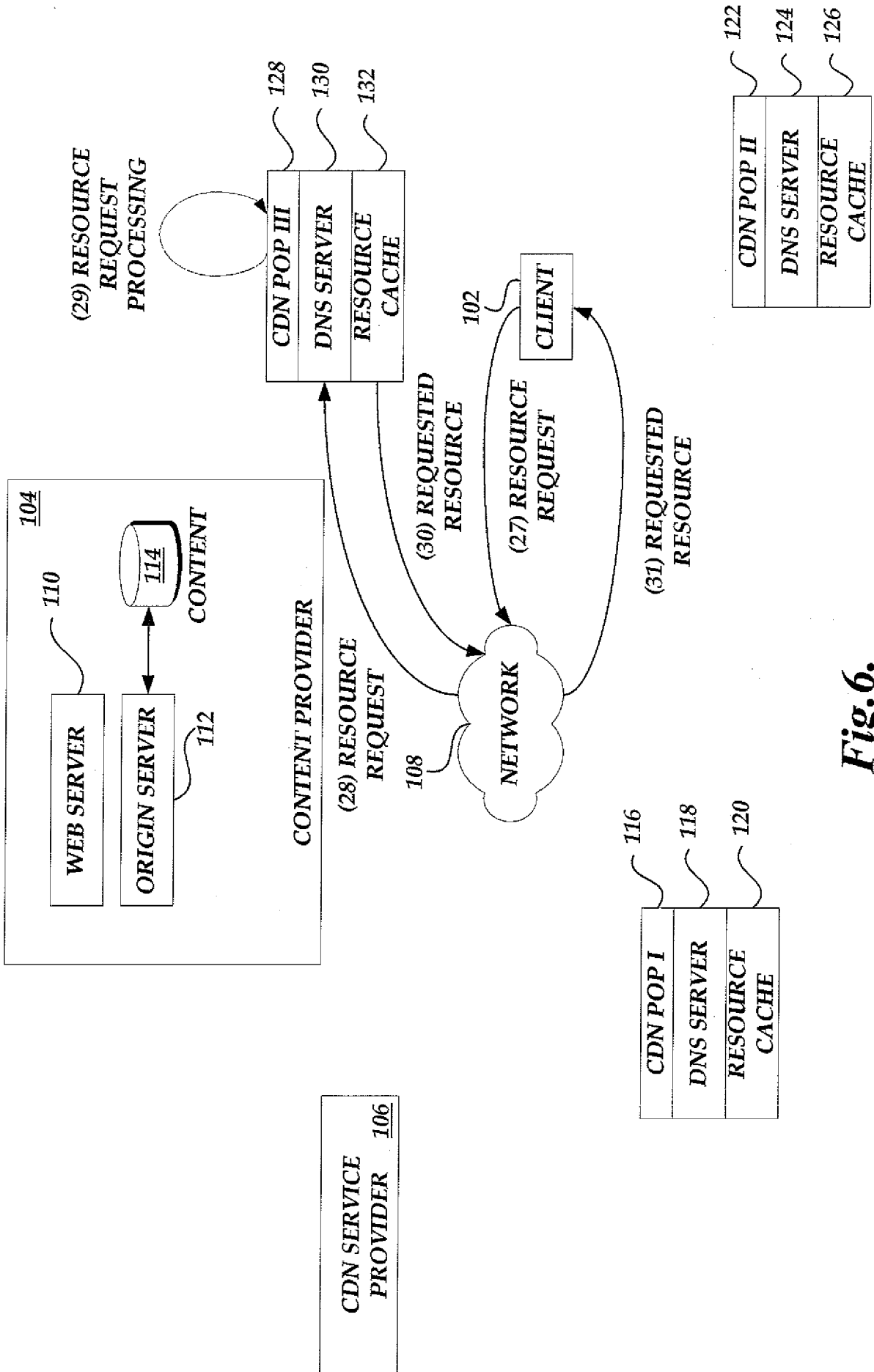


Fig. 6.

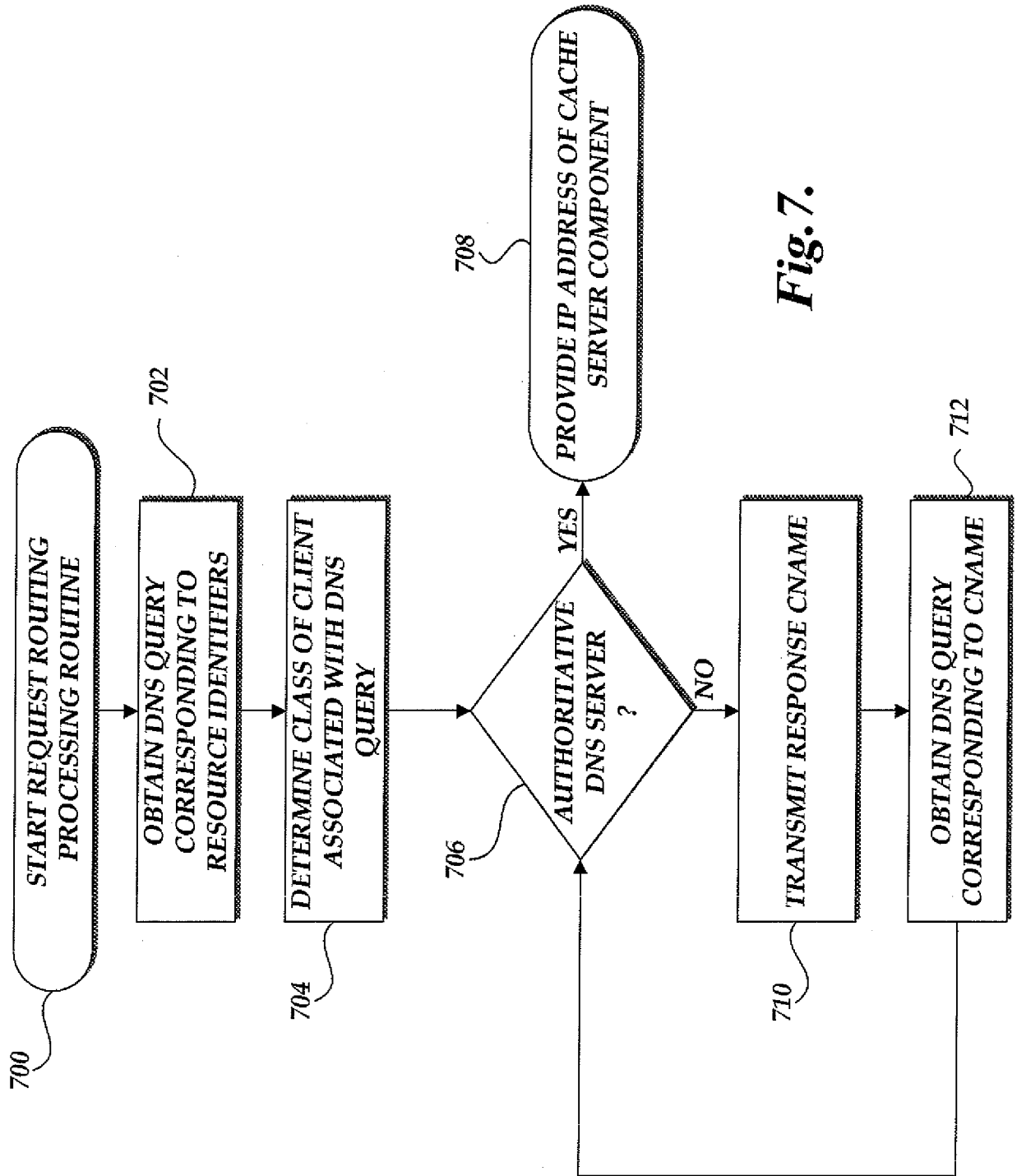


Fig. 7.

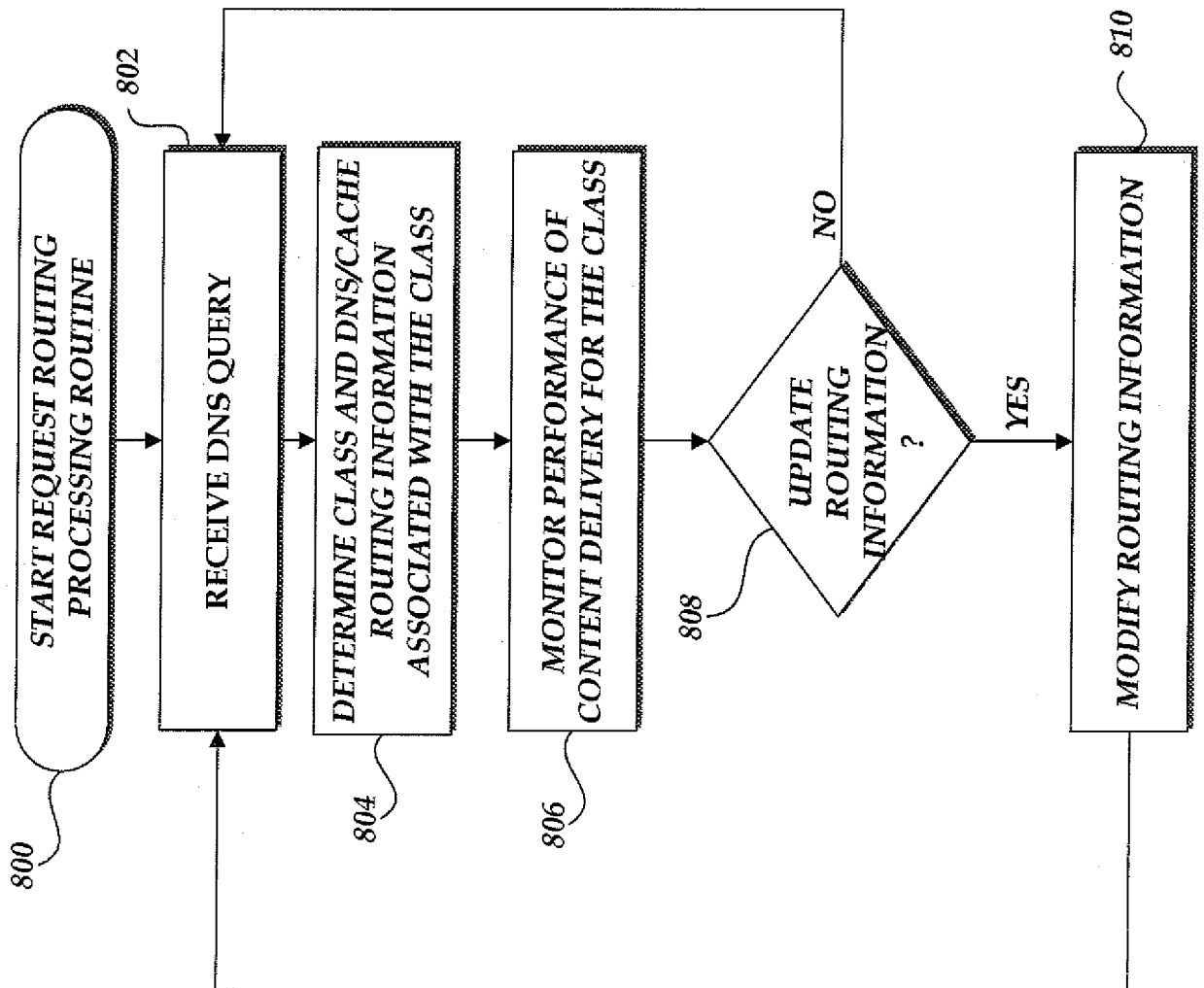


Fig. 8.