

2014년 6월 18일

수신: 고려대학교 산학협력단

참조: 봉수연 선생님

제목: 미국 특허출원 완료 및 양도증 등록 완료 보고의 건 (미국 특허출원 제 14/120,288 호)

귀사관리번호: IP-2014-0065

당소관리번호: OP20140071US

귀사의 무궁한 발전을 기원합니다.

1. 미국 특허출원 완료

당소에 의뢰하신 상기 건의 미국 특허출원을 완료하였기에 그 내역을 알려드립니다. 이에 관하여 현지 대리인으로부터 접수한 서신 및 서류 사본을 동봉하오니 업무에 참조하시기 바랍니다.

2. 양도증 등록완료

상기 건의 발명자로부터 출원인으로서의 발명 양도증이 미국특허청에 등록 완료되었습니다.

3. 국내 특허출원 요약 정보

발명의 명칭	클라우드 컴퓨팅 기반 서버 부하 분산 장치 및 방법		
귀사 담당자	봉수연 선생님		
출원인	고려대학교 산학협력단		
출원번호	10-2013-0054530	출원일자	2013-05-14
발명자	이희조; 알리에브 라샤드; 서동원; 밀번 존		

4. 미국 특허출원 요약 정보

발명의 명칭	DEVICE AND METHOD FOR DISTRIBUTING LOAD OF SERVER BASED ON CLOUD COMPUTING		
출원인	KOREA UNIVERSITY RESEARCH AND BUSINESS FOUNDATION		
출원번호	14/120,288	출원일자	2014-05-14
발명자	Hee Jo LEE; Rashad ALIYEV; Dong Won SEO; John MILBURN		
현지대리인	Pearne & Gordon LLP		

5. 담당자 및 문의처

담당 변리사	이종근	파트너 변리사	02-2051-8200	jglee@mapsip.com
관리 담당자	이영선	선임 매니저	02-2051-8200	yslee@mapsip.com

특허법인 엠에이피에스

동봉물

- 출원지시 서류 사본
- 현지대리인 서신 및 관련서류 사본. 끝.

**DEVICE AND METHOD FOR DISTRIBUTING LOAD OF SERVER BASED ON CLOUD
COMPUTING**

CROSS-REFERNCE TO RELATED APPLICATION

This application claims the benefit of Korean Patent Application No. 10-2013-0054530 filed on May, 14, 2014, the disclosures of which are incorporated herein by reference.

TECHNICAL FIELD

[0001] The embodiments described herein pertain generally to a device or method for defending traffic overload or DDoS attacks, and more particularly, to a device and a method for protecting a server from excessive network traffic utilizing cloud techniques.

BACKGROUND

[0002] A Distributed-Denial-of-Service attack (hereinafter, referred to as "DDoS attack") is one of hacking schemes that attacks a specific site by distributing and arranging a plurality of attackers to thereby simultaneously operate. The DDoS attack implants tools for service attack in a plurality of computers and enables a significantly huge amount of packets that a computer system of a site, an attack target, is incapable of processing to simultaneously flow, thereby degrading performance of a network or paralyzing the computer system.

[0003] Conventionally, a DDoS defense mechanism has been focused on protection of traffic using certain rules of the DDoS attacks. However, newly appearing types of DDoS attacks such as HTTP flood,

Slowloris, RUDY, etc. have traffic patterns similar to normal ones, and, thus, a large amount of malicious traffic can still reach an attack target server even if such rules are applied. Further, if a defense mechanism based on such rules is used, normal traffic concentration such as flash crowds may be misidentified as malicious traffic.

[0004] In this regard, Korean Patent Laid-open Publication No. 10-2012-0066465 (entitled "Method for blocking a denial-of-service attack using an udp flooding") describes a method for blocking DDoS attacks from traffic using certain rules.

SUMMARY

[0005] In view of the foregoing, in order to solve the above-described problem, example embodiments provide a technique capable of continuously providing a service using a cloud replication server even when an overload of normal traffic or a DDoS attack occurs on a target server.

[0006] In accordance with a first aspect, a load distribution device that distributes load of a target server is provided. The load distribution device includes a load detection unit that monitors a load amount of the target server and determines whether the load amount exceeds a predetermined critical value, a server driving unit that drives a replication server when the load amount exceeds the critical value, and a server control unit that distributes part of load to the replication server when the replication server has started to be driven. The replication server is implemented by a cloud computing technique.

[0007] In accordance with a second aspect, a load distribution method of a load distribution device for distributing load of a target server is provided. The load distribution method includes monitoring a load status of the target server when the target server is driven and a service is provided, activating a replication server when a load amount of the target server exceeds a predetermined critical value, and distributing part of load of the target server to the replication server using a load distribution scheme when the replication server is activated. The replication server is implemented by a cloud computing technique.

[0008] In accordance with the various aspects and example embodiments, performance of an attack target server is not degraded due to a DDoS attack or a traffic overload, and the service provider can keep providing their services.

[0009] Further, in accordance with the various aspects and example embodiments, a false positive, in which a normal user is misidentified as a malicious user during traffic overload, and, thus, a service provided to a target server is stopped, is not generated.

The foregoing summary is illustrative only and is not intended to be in any way limiting. In addition to the illustrative aspects, embodiments, and features described above, further aspects, embodiments, and features will become apparent by reference to the drawings and the following detailed description.

BRIEF DESCRIPTION OF THE DRAWINGS

In the detailed description that follows, embodiments are described as illustrations only since various changes and

modifications will become apparent to those skilled in the art from the following detailed description. The use of the same reference numbers in different figures indicates similar or identical items.

[0010] FIG. 1 is a diagram for describing a filter propagation method as a conventional DDoS defense mechanism.

[0011] FIG. 2 is a diagram for describing an operation of a load distribution device in accordance with an example embodiment.

[0012] FIG. 3 illustrates a detailed configuration of a load distribution device of a target server in accordance with an example embodiment.

[0013] FIG. 4 illustrates an example of system construction of a server control unit to distribute traffic or a load in accordance with an example embodiment.

[0014] FIG. 5 is a flowchart for describing a method of a load distribution device for distributing traffic or a load of a target server in accordance with an example embodiment.

DETAILED DESCRIPTION

[0015] Hereinafter, embodiments of the present disclosure will be described in detail with reference to the accompanying drawings so that the present disclosure may be readily implemented by those skilled in the art. However, it is to be noted that the present disclosure is not limited to the embodiments but can be embodied in various other ways. In drawings, parts irrelevant to the description are omitted for the simplicity of explanation, and like reference numerals denote like parts through the whole document.

[0016] Through the whole document, the term "connected to" or "coupled to" that is used to designate a connection or coupling of one element to another element includes both a case that an element is "directly connected or coupled to" another element and a case that an element is "electronically connected or coupled to" another element via still another element. Further, the term "comprises or includes" and/or "comprising or including" used in the document means that one or more other components, steps, operation and/or existence or addition of elements are not excluded in addition to the described components, steps, operation and/or elements unless context dictates otherwise.

[0017] Through the whole document, the term "step of" does not mean "step for".

[0018] Through the whole document, the term "traffic" means a load given to a certain communication device or system unless context dictates otherwise.

[0019] **FIG. 1** is a diagram for describing a filter propagation method as a conventional DDoS defense mechanism.

[0020] A filter propagation method as a conventional DDoS defense mechanism has been focused on defending an attack target server by installing a firewall or an IDS/IPS (Intrusion Detection/Protection System).

[0021] However, in the case of a new malicious attack in the form similar to a normal state or in the case of normal and temporary traffic concentration, excessive traffic is concentrated on a target server, resulting in a service fault of the server.

[0022] FIG. 2 is a diagram for describing an operation of a load distribution device in accordance with an example embodiment. According to the conventional method, a route from a user to a target server is limited. However, in accordance with the present disclosure, replication servers distributed in several sites provide services instead of the target server, and, thus, a route from a user to the target server is diversified and most of traffic does not reach the target server.

[0023] In order to solve the conventional problem, in the present disclosure, replication servers that perform the same function as the target server are constructed using cloud techniques, and when excessive traffic is concentrated on the target server, the traffic is distributed to the replication servers, and, thus, a service can be continuously provided.

[0024] FIG. 3 illustrates a detailed configuration of a load distribution device of a target server in accordance with an example embodiment.

[0025] A load distribution device 100 includes a load detection unit 110, a server driving unit 120, a sever control unit 130, and a filter 140.

[0026] The load detection unit 110 monitors a load amount of a target server 10 and determines whether or not traffic concentration occurs on the target server 10. In accordance with an example embodiment, the load detection unit 110 can determine whether or not traffic concentration occurs based on whether or not a load amount of the target server 10 exceeds a predetermined critical value. In accordance with the example embodiment, the critical value can be

determined by a service provider in preparation for a DDoS attack that makes it impossible for a system to provide a normal service any more by distributing and arranging multiple attackers to thereby simultaneously make denial of service (DoS) attacks. Further, in order to provide a higher quality service, the service provide may set a critical value to be low such that a load amount over a certain level can be detected.

[0027] The server driving unit 120 drives a replication server 30 for distribution of a load when traffic concentration occurs.

[0028] Further, when the replication server 30 is driven according to an operation of the server driving unit 120, the server control unit 130 distributes traffic or a load to the replication server 30.

[0029] The replication server 30 can be implemented by a cloud computing technique. However, the replication server 30 is not necessarily implemented by the cloud computing technique, but can be configured as a separate internal or external resource. The cloud computing technique is a technique of virtually integrating resources of computers present in different physical locations, and, thus, makes it possible to efficiently use a resource of the replication server 30. In accordance with the present disclosure, the replication server 30 can be constructed using a resource of a server in a virtual space.

[0030] The replication server 30 driven by the server driving unit 120 is classified into three types depending on a construction scheme.

[0031] Firstly, the replication server 30 can be configured by replicating the whole content of the target server 10 into the replication server, which takes a long time to replicate and requires a lot of resources of a storage device, but most stably provides a service to a user.

[0032] Secondly, the replication server can be configured by replicating a specific content frequently requested by a user into the replication server. Such an interest-based replication server can determine whether a content is frequently requested by the user based on the number of user requests for the content. The interest-based replication server requires relatively less resources, but a service provider needs to monitor which content users have been interested in, and update content of the replication server accordingly.

[0033] Finally, a content type-based replication server classifies content into multimedia files, text files, user files, and the like, and then stores the classified content in the replication server. That is, a replication server is in charge of one or more content types. Herein, the content type may refer to a file format, a predefined category or the like of a content.

[0034] The server control unit 130 may use the following method as a method of distributing traffic or a load to the replication server 30.

[0035] A DNS-based load distribution method dynamically uses DNS Round Robin depending on the situation. DNS Round Robin is one of techniques of distributing a service to multiple servers using a DNS (Domain Name System). By way of example, if a server having an IP

address of 1.1.1.1 is in charge of a service regarding www.example.com, when excessive traffic is concentrated, IP addresses of 1.1.1.2, 1.1.1.3, etc. of the replication servers 30 are additionally registered as servers in charge of the corresponding domain, so that traffic of a user can be distributed to the replication servers 30.

[0036] A network switch has a function of delivering a packet having a specific IP range as a source IP address or a packet selected with a certain probability to a specified target. A switch-based load distribution method distributes traffic to the replication server 30 using such a function.

[0037] FIG. 4 illustrates an example of system construction of a server control unit to distribute traffic or a load in accordance with an example embodiment. By way of example, www1 is a web server, and www2 and www3 are replication servers that perform the same function as www1. Traffic toward www1 from users can be distributed to www2 and www3 by DNS Round Robin, a packet delivery function of a switch, or others.

[0038] A network can be implemented in a wired network such a Local Area Network (LAN), a Wide Area Network (WAN), or a Value-Added Network (VAN), or all kinds of wireless network such as mobile radio communication network or a satellite communication network.

[0039] The present disclosure may further include a filter 140. A filter in accordance with an example embodiment is a component configured to process traffic generated by a malicious code among traffic to be distributed to the replication server by the sever control unit 130. The filter 140 is a component configured to

distribute traffic to the replication server when the target server 10 is attacked by a malicious code, and also to perform an extra process regarding the malicious code.

[0040] In accordance with an example embodiment, the server driving unit 130 can inactivate the replication server 30 when traffic concentration is ended, i.e. when a load amount does not exceed a predetermined critical value.

[0041] FIG. 5 is a flowchart for describing a method of a load distribution device for distributing traffic or a load of a target server in accordance with an example embodiment.

[0042] When a target server is being driven and a service is provided, a load distribution device monitors a load status of the target server (S410).

[0043] Then, in the case of normal traffic concentration such as flash crowds referring to a phenomenon in which after a DDoS attack or a some interesting event or announcement occurs, the number of people accessing a relevant site suddenly increases, a replication server is activated (S420).

[0044] In accordance with an example embodiment, whether or not traffic concentration occurs can be determined based on whether or not a load amount of the target server exceeds a predetermined critical value.

[0045] The activated replication server can be classified into three types: a whole replication server; an interest-based replication server; and a content type-based replication server, depending on a construction scheme.

[0046] Firstly, a replication server can be configured by replicating the whole content of a target server into the replication server, which takes a long time to replicate and requires a lot of resources of a storage device, but most stably provides a service to a user.

[0047] Secondly, a replication server can be configured by replicating a specific content frequently requested by a user into the replication server. Such an interest-based replication server can determine whether a content is frequently requested by the user based on the number of user requests for the content. The interest-based replication server requires relatively less resources, but a service provider needs to monitor which content users have been interested in, and update content of the replication server accordingly.

[0048] If a replication server is configured as an interest-based replication server, before the replication server is activated (S420), a step of checking whether the user-requested content has been replicated into the interest-based replication server may be further included in order to redistribute the load caused by the request of the user for the interest-based content into the replication server.

[0049] Finally, a content type-based replication server classifies content into multimedia files, text files, user files, and the like, and then stores the classified content in the replication server. That is, a replication server is in charge of one or more content types. Herein, a content type may be a file format, a predefined category or the like of a content.

[0050] If a replication server is configured as a content type-based replication server, before the replication server is activated (S420), a step of checking the type of the user-requested content may be further included, and in a step of distributing a load (S430) to be described later, a load of the target server can be distributed in order to redistribute a load to each replication server depending on a type of user content.

[0051] The replication server 30 can be implemented by a cloud computing technique. However, the replication server 30 is not necessarily implemented by the cloud computing technique, but can be configured as a separate internal or external resource. The cloud computing technique is a technique of virtually integrating resources of computers present in different physical locations, and, thus, makes it possible to efficiently use a resource of the replication server 30. In accordance with the present disclosure, the replication server 30 can be constructed using a resource of a server in a virtual space.

[0052] Then, a load is distributed using a load distribution scheme (S430). The following method may be used as a method of distributing a load.

[0053] A DNS-based load distribution method dynamically uses DNS Round Robin and a client characteristic-based method depending on the situation. According to the DNS Round Robin, it is possible to distribute a service to multiple servers using a DNS (Domain Name System). Using the client characteristic-based method, it is possible to distribute clients to multiple servers based on their characteristics.

[0054] A network switch has a function of delivering a packet having a specific IP range as a source IP address or a packet selected with a certain probability to a specified target. A switch-based load distribution method distributes traffic to a replication server using such a function.

[0055] Then, in accordance with an example embodiment, a process of known malicious traffic can be determined using a filter (S440).

[0056] Thereafter, in accordance with an example embodiment, a load status of the target server is continuously monitored, and when traffic overload on the target server is ended, the replication server is inactivated (S450).

[0057] According to the load distribution device or the load distribution method in accordance of the present disclosure, performances of an attacked target server is not degraded due to a DDoS attack or traffic overload, and the service provider can provide their services without service fault. Further, a false positive, in which a normal user is misidentified as a malicious user during traffic overload, and, thus, a service provided to a target server is stopped, is not generated.

[0058] For reference, each of components illustrated in **FIG. 3** in accordance with an example embodiment may imply software or hardware such as a field programmable gate array (FPGA) or an application specific integrated circuit (ASIC), and they carry out a predetermined function.

[0059] However, the components are not limited to the software or the hardware, and each of the components may be stored in an

addressable storage medium or may be configured to implement one or more processors.

[0060] Accordingly, the components may include, for example, software, object-oriented software, classes, tasks, processes, functions, attributes, procedures, sub-routines, segments of program codes, drivers, firmware, micro codes, circuits, data, database, data structures, tables, arrays, variables and the like.

[0061] The components and functions thereof can be combined with each other or can be divided.

[0062] The illustrative embodiments can be embodied in a storage medium including instruction codes executable by a computer or processor such as a program module executed by the computer or processor. A data structure in accordance with the illustrative embodiments can be stored in the storage medium executable by the computer or processor. A computer readable medium can be any usable medium which can be accessed by the computer and includes all volatile/non-volatile and removable/non-removable media. Further, the computer readable medium may include all computer storage and communication media. The computer storage medium includes all volatile/non-volatile and removable/non-removable media embodied by a certain method or technology for storing information such as computer readable instruction code, a data structure, a program module or other data. The communication medium typically includes the computer readable instruction code, the data structure, the program module, or other data of a modulated data signal such as a carrier wave, or other transmission mechanism, and includes information transmission mediums.

[0063] The load distribution device and method in accordance with the present disclosure can be implemented by a computer-readable code in a computer-readable storage medium. The computer-readable storage medium includes all kinds of storage media in which computer-readable data are stored and may include, for example, a ROM (Read Only Memory), a RAM (Random Access Memory), a magnetic tape, a magnetic disc, a flash memory, an optical data storage device, etc. Further, the computer-readable storage medium can be distributed in a computer system connected via a computer communication network and can be stored and executed as a code that is readable in a distribution manner.

[0064] The device and method of the present disclosure has been explained in relation to a specific embodiment, but its components or a part or all of its operation can be embodied by using a computer system having general-purpose hardware architecture can be applied.

[0065] The above description of the present disclosure is provided for the purpose of illustration, and it would be understood by those skilled in the art that various changes and modifications may be made without changing technical conception and essential features of the present disclosure. Thus, it is clear that the above-described embodiments are illustrative in all aspects and do not limit the present disclosure. For example, each component described to be of a single type can be implemented in a distributed manner. Likewise, components described to be distributed can be implemented in a combined manner.

[0066] The scope of the present disclosure is defined by the following claims rather than by the detailed description of the embodiment. It shall be understood that all modifications and embodiments conceived from the meaning and scope of the claims and their equivalents are included in the scope of the present disclosure.

WE CLAIM

1. A load distribution device that distributes a load of a target server, the load distribution device comprising:

a load detection unit that monitors a load amount of the target server and determines whether the load amount exceeds a predetermined critical value;

a server driving unit that drives a replication server when the load amount exceeds the critical value; and

a server control unit that distributes part of load to the replication server when the replication server has started to be driven,

wherein the replication server is implemented by a cloud computing technique.

2. The load distribution device of Claim 1,

wherein the server driving unit drives a whole replication server into which whole content of the target server has been replicated.

3. The load distribution device of Claim 1,

wherein the server driving unit drives an interest-based replication server into which part of content frequently requested by a user more than certain number of times has been replicated from the target server.

4. The load distribution device of Claim 1,

wherein the server driving unit drives a content type replication server into which part of content classified by content type has been replicated from the target server.

5. The load distribution device of Claim 1,

wherein the server control unit distributes a load of the target server by a DNS distribution method in which part of load is distributed using a DNS, or using a switch-based load distribution method in which packets selected with a certain probability is delivered to a specified target.

6. The load distribution device of Claim 1, further comprising:

a filter that processes traffic generated by a malicious code, among traffic to be distributed to the replication server by the server control unit.

7. The load distribution device of Claim 1,

wherein the server driving unit inactivates the replication server when the load amount does not exceed the critical value.

8. A load distribution method of a load distribution device for distributing a load of a target server, the load distribution method comprising:

monitoring a load status of the target server when the target server is driven and a service is provided;

activating a replication server when a load amount of the target server exceeds a predetermined critical value; and

distributing part of load of the target server to the replication server using a load distribution scheme when the replication server is activated,

wherein the replication server is implemented by a cloud computing technique.

9. The load distribution method of Claim 8,

wherein the activating of the replication server includes activating a whole replication server into which whole content of the target server has been replicated.

10. The load distribution method of Claim 8,

wherein the activating of the replication server includes activating an interest-based replication server into which part of content frequently requested by a user more than certain number of times has been replicated from the target server.

11. The load distribution method of Claim 8,

wherein the activating of the replication server includes activating a content type replication server into which part of content classified by content type has been replicated from the target server, and

the distributing of the load includes distributing part of load of the target server depending on the content type.

12. The load distribution method of Claim 8, wherein the distributing of the load includes distributing part of load by a DNS distribution method in which the part of load is distributed using a DNS, or using a switch-based load distribution method in which packets selected with a certain probability is delivered to a specified target.

13. The load distribution method of Claim 8, further comprising:

filtering traffic generated by a malicious code, among traffic to be distributed to the replication server.

14. The load distribution method of Claim 8, further comprising:

inactivating the replication server when the load amount of the target server does not exceed the critical value.

ABSTRACT

A load distribution device that distributes load of a target server is provided. The load distribution device includes a load detection unit that monitors a load amount of the target server and determines whether the load amount exceeds a predetermined critical value, a server driving unit that drives a replication server when the load amount exceeds the critical value, and a server control unit that distributes part of load to the replication server when the replication server has started to be driven. The replication server is implemented by a cloud computing technique.

FIG. 1
(PRIOR ART)

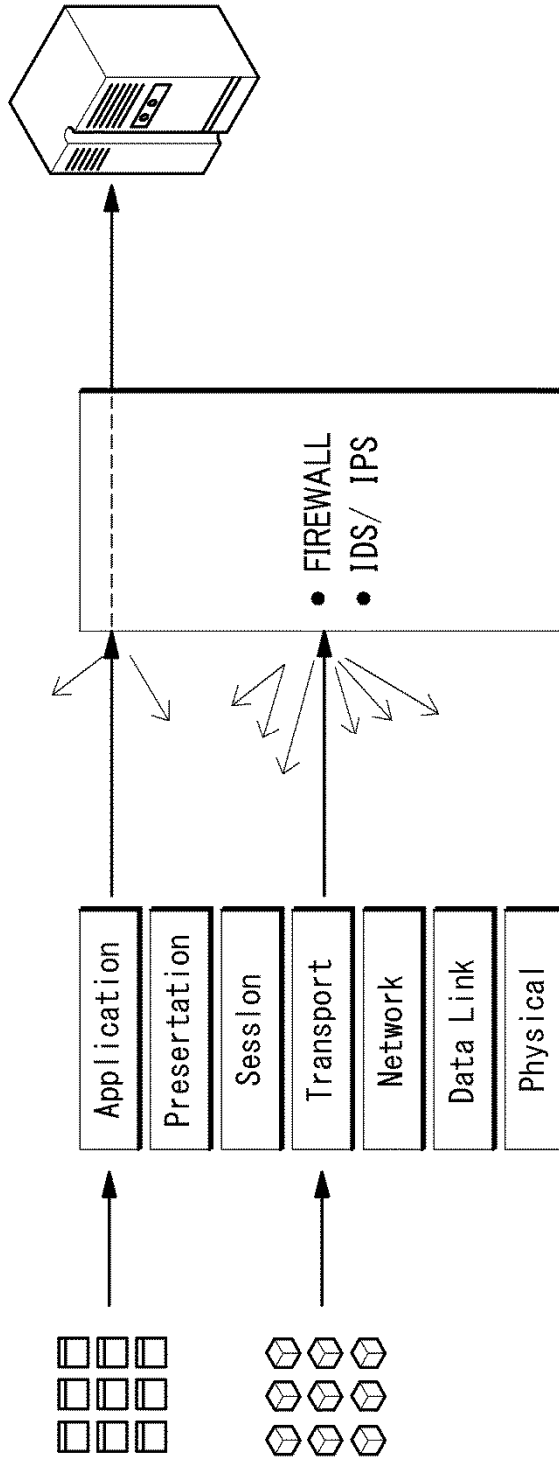
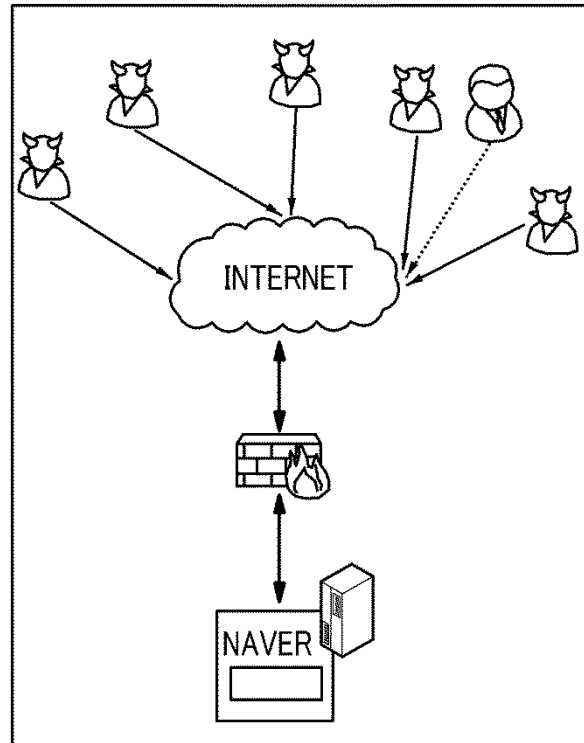


FIG. 2

Existing solution



Cloud-based solution

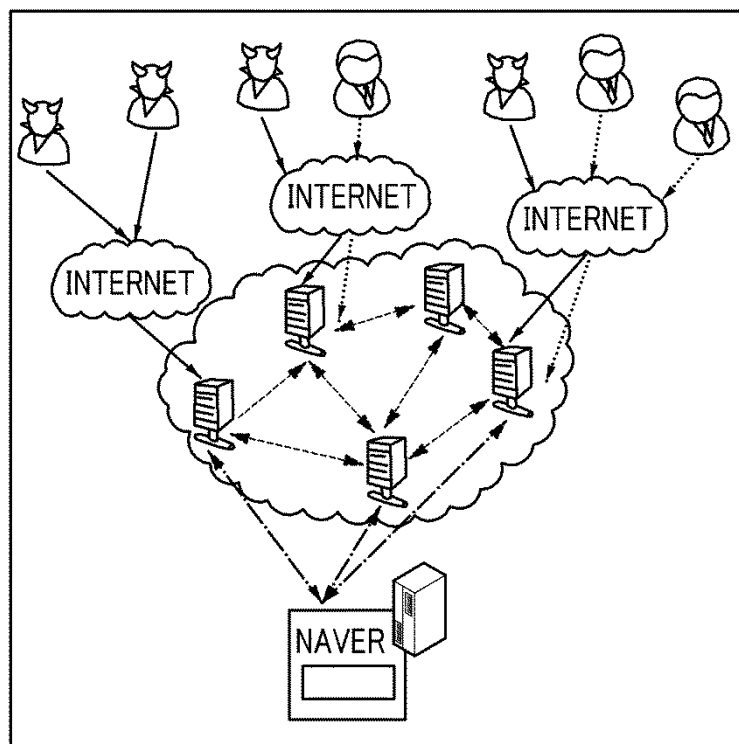


FIG. 3

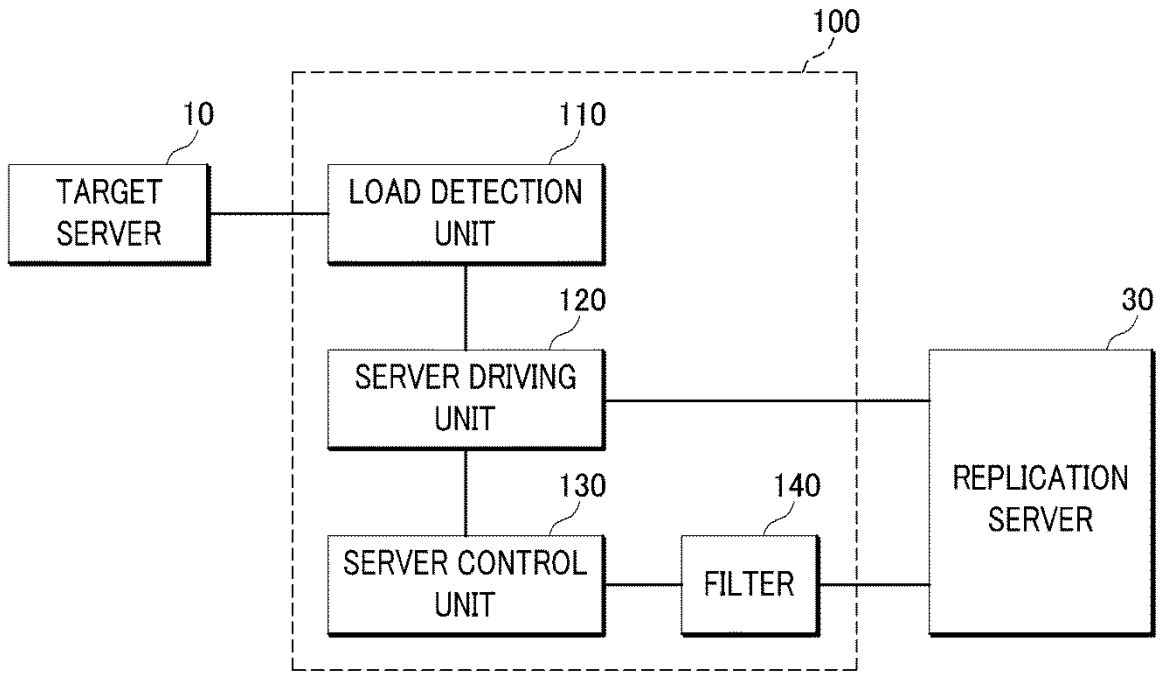


FIG. 4

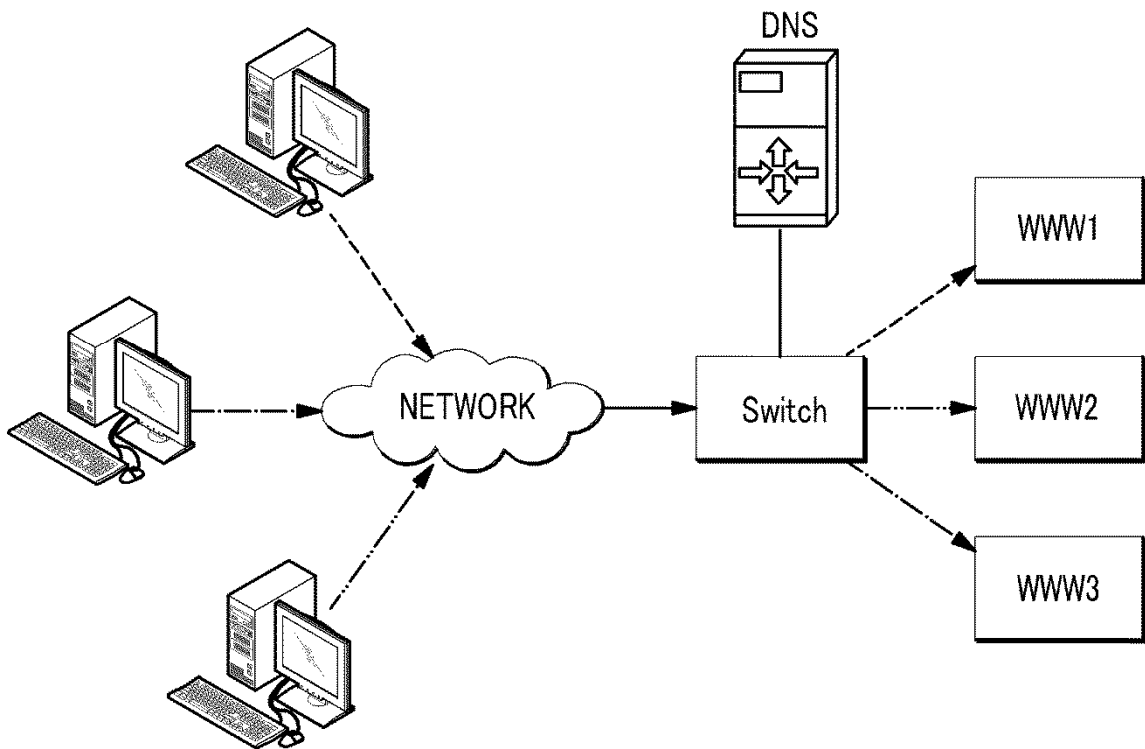


FIG. 5