

US PATENT & TRADEMARK OFFICE

PATENT APPLICATION FULL TEXT AND IMAGE DATABASE



(1 of 1)

United States Patent Application	20170237684
Kind Code	A1
Smith II; James Thomas ; et al.	August 17, 2017

DECENTRALIZED RESOURCE ALLOCATION

Abstract

The disclosure is directed to routing service requests over a network (50). Service requests may be routed over a network (50) based upon deriving optimized weights for each of a plurality of service providers (130) within a service provider set (135), receiving a plurality of service requests at a broker (110) within the network (50), and routing each of the plurality of service requests from the broker (110) to the plurality of service providers (130) on the network (50). In some implementations, the optimized weights for each of the plurality of service providers (130) may be derived using a non-linear function. In some implementations, the optimized weights for the plurality of service providers (130) associated with a broker (110) may collectively define a weighted distribution. The plurality of service requests may be routed by a broker (110) using its corresponding weighted distribution.

Inventors: **Smith II; James Thomas; (Boulder, CO) ; Shestak; Vladimir Vladimirovich; (Boulder, CO)**

Applicant: **Name** **City** **State** **Country** **Type**

Webscale Networks, Inc. Boulder CO US

Family ID: **48875152**

Appl. No.: **15/275994**

Filed: **September 26, 2016**

Related U.S. Patent Documents

<u>Application Number</u>	<u>Filing Date</u>	<u>Patent Number</u>
13908497	Jun 3, 2013	
15275994		
61654992	Jun 4, 2012	

Current U.S. Class: **709/226**

Current CPC Class: H04L 67/16 20130101; H04L 47/783 20130101; H04L 67/101 20130101; H04L 67/2895 20130101; H04L 67/1004 20130101; H04L 67/1029 20130101; H04L 67/2809 20130101; H04L 67/1034 20130101; H04L 41/5083 20130101; H04L 47/70

14. The method of claim 1, wherein a first portion of said plurality of service requests are routed during said routing step using a first execution of said deriving step, wherein a second portion of said plurality of service requests are routed during said routing step using a second execution of said deriving step, and wherein said first and second portions are completely independent of one another.

15-16. (canceled)

17. The method of claim 1, further comprising: converting said plurality of optimized weights into a cumulative probability mass function; and using said cumulative probability mass function to select said at least one service provider of said plurality of service providers.

18. The method of claim 1, wherein said network is one of a public cloud or a private cloud.

19. (canceled)

20. A method of routing service requests over a network, comprising the steps of: obtaining, by each broker of a plurality of brokers in a network, status information for each service provider of a plurality of service providers within a service provider set on the network, wherein said status information includes a maximum number of connections of said service provider and a number of busy connections of said service provider; deriving, by a processor of each broker of said plurality of brokers, optimized weights for each service provider of said plurality of service providers based on the obtained status information, wherein said optimized weights for said plurality of service providers collectively define a weighted distribution; receiving, at each broker of said plurality of brokers, a plurality of service requests over said network for processing at the service provider set; and routing each service request of said plurality of service requests from each broker of said plurality of brokers to at least one service provider of said service provider set over said network based on said weighted distribution.

21-38. (canceled)

39. A method of routing service requests over a network, comprising the steps of: acquiring, by each broker of a plurality of brokers over a network, status information on each service provider within a service provider set that is associated with said broker, wherein each broker of said plurality of brokers acquires the status information in a manner that is temporally independent from other brokers of said plurality of brokers; deriving, by each broker of said plurality of brokers with said acquired status information, a weighted distribution that comprises an optimized weight for each service provider within said service provider set; receiving, by each broker of said plurality of brokers, a plurality of service requests over said network; and routing, by each broker of said plurality of brokers over said network, said received service requests to one or more of said service providers of the service provider set associated with said broker using said derived weighted distribution.

40-60. (canceled)

61. The method of claim 39, wherein said optimized weight for each service provider is at least partially based on a current demand for a capacity of said service provider.

62. The method of claim 39, wherein each broker of said plurality of broker performs said deriving step autonomously from other brokers of said plurality of brokers.

63. The method of claim 1, wherein each of said deriving, receiving and routing steps is autonomously performed by each broker of said plurality of brokers in said network.

64. The method of claim 63, wherein a first broker of said plurality of brokers derives a first optimized weight for a first service provider of said plurality of service providers, wherein a second broker of said plurality of brokers derives a second optimized weight for the first service provider, and wherein the first and second optimized weights are different.

65. The method of claim 1, wherein each optimized weight of said plurality of optimized weights is at least

discussed below in relation to the status acquiring protocol 200 of FIG. 3). The broker 110 could use status information received from a service provider 130 over the network 50 "as is," or the broker 110 could use status information received from a service provider 130 over the network 50 to derive status information. Another option for a broker 110 to acquire status information on each service provider 130 within its corresponding service provider set 135 is to poll other brokers 110 within the broker set 100 to allow the requesting broker 110 to estimate/calculate the status information for each service provider 130 within its corresponding service provider set 135 (e.g., to identify the extent to which other brokers 110 are transmitting requests to the service providers 130 and in the manner discussed below in relation to the status acquiring protocol 250 of FIG. 4)).

[0037] Each broker 110 will calculate a price (e.g., using at least one processor) for each service provider 130 within its service provider weight set 135 (step 162) using the status information acquired pursuant to step 160. This will be discussed in more detail below. Optimized weights for each service provider 130 (one weight per service provider 130) are derived by each broker 110 (step 164) using the price data acquired pursuant to step 162. This will be discussed in more detail below. In one embodiment, the derivation associated with step 164 utilizes a non-linear function (discussed below). Step 166 of the protocol 150 is directed to assigning the optimized weights from step 164 to what is now an updated or current service provider weight set for the given broker 110.

[0038] With the service provider weight set being current, the protocol 150 of FIG. 2 proceeds to step 154. Step 154 is directed to determining if a given broker 110 has received a request for service over the network 50 from one or more service requesters 120. Each broker 110 may very well continually receive a relatively large volume of service requests from various service requesters 120. Assuming that there are service requests to be distributed by a given broker 110, the broker 110 transmits these service requests over the network 50 to one or more of the service providers 130 in its service provider set 135 using its own optimized weights (steps 160-166). In the illustrated embodiment, the optimized weights for the service providers 130 of a given broker 110 are converted into a cumulative probability mass function (step 156). Service requests from the service requesters 120 are then transmitted by the broker 110 to the various service providers 130 within its service provider set 135 using this cumulative probability mass function (step 158).

[0039] Each broker 110 may execute the protocol 150 of FIG. 2 independently of other brokers 110 within the broker set 100. For instance, the brokers 110 need not be synchronized on any basis for updating their respective optimized weights (e.g., optimized weights for each of the brokers 110 need not be based upon status information from the same point in time). Although two or more brokers 110 in the broker set 100 could acquire status information from their corresponding service provider set 135 that is associated with a common point in time, such is not required by the brokers 110 of the broker set 100. This may be referred to as being "temporally independent status information" on the various service providers 130 for use by the brokers 100 in deriving optimized weights for the service providers 130 in its corresponding service provider set 135.

[0040] FIG. 3 presents one embodiment of a protocol 200 which may be utilized by each broker 110 to acquire and/or request status from each service provider 130 within its corresponding service provider set 135 (e.g., step 160 of protocol 150-FIG. 2). As discussed above, a broker 110 can request status information from each service provider 130 within its corresponding service provider set 135. This may be done by a direct network request (e.g., over the network 50) to the service providers 130. Many service providers (e.g., common web server implementations such as NGinX and Apache) offer a well known mechanism for obtaining the current status for the service provider. For example, an Apache web server provides the well known module `mod_status` that provides the current number of busy and idle threads. Status acquiring protocol 200 is generally directed towards acquiring status information on each service provider 130 by requesting the status information from each service provider 130. The status information may include a capacity of service provider 130 and a flow of service provider 130 (e.g., the current utilization of that service provider 130). The flow/utilization may be a count of the service provider 130 as opposed to a percentage of the service provider utilized. For example, the capacity of the service provider 130 and the flow of the service provider 130 may be acquired in at least two ways and dependent upon whether the service provider 130 is multi-threaded or not. In this regard, step 210 of status acquiring protocol 200 is directed to determining whether each service provider 130 of the service provider set 135 is a multi-threaded service provider.

