



(19)中華民國智慧財產局

(12)發明說明書公告本

(11)證書號數：TW I552002 B

(45)公告日：中華民國 105 (2016) 年 10 月 01 日

(21)申請案號：103114547

(22)申請日：中華民國 103 (2014) 年 04 月 22 日

(51)Int. Cl. : G06F15/163 (2006.01)

G06F15/17 (2006.01)

G06F9/50 (2006.01)

(71)申請人：財團法人工業技術研究院(中華民國) INDUSTRIAL TECHNOLOGY RESEARCH INSTITUTE (TW)

新竹縣竹東鎮中興路 4 段 195 號

國立交通大學(中華民國) NATIONAL CHIAO TUNG UNIVERSITY (TW)

新竹市大學路 1001 號

(72)發明人：丁韋智 TING, WEI CHIH (TW)；王濬哲 WANG, JUN ZHE (TW)；陳家旻 CHEN, CHIA MIN (TW)；黃俊龍 HUANG, JUN LONG (TW)

(74)代理人：洪堯順；侯德銘

(56)參考文獻：

CN 102855171A

US 2002/0184575A1

US 2003/0105868A1

US 2009/0276771A1

審查人員：易昶霈

申請專利範圍項數：19 項 圖式數：10 共 47 頁

(54)名稱

公共雲資源動態配置方法及系統

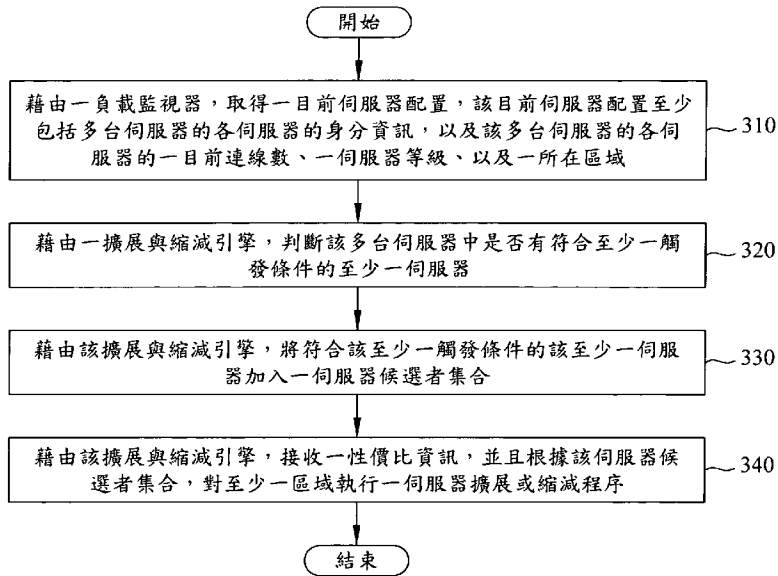
METHOD AND SYSTEM FOR DYNAMIC INSTANCE DEPLOYMENT OF PUBLIC CLOUD

(57)摘要

根據一實施例，一種公共雲資源動態配置方法中，藉由一負載監視器，取得一目前伺服器配置，此目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域；藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合，以及接收一性價比資訊，並且根據此伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序。

According to one exemplary embodiment, a method for dynamic instance deployment of public cloud, uses a load monitor to obtain a current configuration at least including identity information of each of a plurality of servers, and a number of current connections, a server level and a located area of each of the plurality of servers; and uses a scaling engine to determine whether there is at least one server of the plurality of servers satisfies at least one trigger condition, add the at least one server that satisfies the at least one trigger condition into a server candidate set, and receive a performance ratio information to perform a server scaling procedure according to the server candidate set.

指定代表圖：



第三圖

符號簡單說明：

310 . . . 藉由一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域

320 . . . 藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器

330 . . . 藉由該擴展與縮減引擎，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合

340 . . . 藉由該擴展與縮減引擎，接收一性價比資訊，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序

發明摘要

※申請案號：103114541

※申請日：

103. 4. 22

※IPC 分類：

G06F 15/163 (2006.01)

G06F 15/17 (2006.01)

G06F 9/50 (2006.01)

【發明名稱】(中文/英文)

公共雲資源動態配置方法及系統

METHOD AND SYSTEM FOR DYNAMIC INSTANCE
DEPLOYMENT OF PUBLIC CLOUD

【中文】

根據一實施例，一種公共雲資源動態配置方法中，藉由一負載監視器，取得一目前伺服器配置，此目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域；藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合，以及接收一性價比資訊，並且根據此伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序。

【英文】

According to one exemplary embodiment, a method for dynamic instance deployment of public cloud, uses a load

monitor to obtain a current configuration at least including identity information of each of a plurality of servers, and a number of current connections, a server level and a located area of each of the plurality of servers; and uses a scaling engine to determine whether there is at least one server of the plurality of servers satisfies at least one trigger condition, add the at least one server that satisfies the at least one trigger condition into a server candidate set, and receive a performance ratio information to perform a server scaling procedure according to the server candidate set.

【代表圖】

【本案指定代表圖】：第（三）圖。

【本代表圖之符號簡單說明】：

310 藉由一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域

320 藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器

330 藉由該擴展與縮減引擎，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合

340 藉由該擴展與縮減引擎，接收一性價比資訊，並且根

據該伺服器候選者集合，對至少一區域執行一伺服器擴展
或縮減程序

【本案若有化學式時，請揭示最能顯示發明特徵的化學式】：

發明專利說明書

(本說明書格式、順序，請勿任意更動)

【發明名稱】(中文/英文)

公共雲資源動態配置方法及系統

METHOD AND SYSTEM FOR DYNAMIC INSTANCE
DEPLOYMENT OF PUBLIC CLOUD

【技術領域】

本揭露係關於一種公共雲(public cloud)資源動態配置方法及系統。

【先前技術】

網路直播服務如雨後春筍般發展，使用者可以經由網路即時觀賞影片直播，例如遊戲類、娛樂類、新聞類、體育節目類、科技類等。隨著普及的網路直播串流，即時串流服務需要大量且穩定的頻寬。同儕網路的串流影音技術利用網路中各節點間互相分享資料的方法，來增加串流傳輸的效率。在同儕網路中，使用者數目的波動、使用者設備的不良、使用者設備的頻寬的不足、使用者設備距離機房太遠等因素，可能使得即時串流服務網提供的串流品質不穩定。結合伺服器與同儕網路的架構利用分散式伺服器提供穩定的串流輸出來維持使用者的觀看品質。

隨著普及的行動裝置，例如手持式攝影裝置，使用者可

以是串流提供者。不論是播放者或是觀看者，都可以隨時隨地的播放與觀看。此趨勢下，串流平台對於伺服器需求量的負擔將不斷增加，服務業者搭配公共雲業者在公共雲建置分散式伺服器，利用伺服器做為轉繼站(relay)，來符合彈性化的需求。例如，預先評估使用網路直播服務的可容納的最大上線人數，以及事先建立數量足夠的虛擬機器 (virtual Machines, VM) 如雲端伺服器。

● 即使能夠預估網路直播服務的使用者的數量與行為，要滿足如尖峰時段時的使用者的觀看品質，需要建立龐大數量的伺服器來進行待命。在不確定影響範圍的情境下，例如在離峰時段，難以預估使用者數量以及觀看行為的狀況下，需要人員密切注意雲端伺服器的連線情形，也不適合將閒置的伺服器貿然關閉。在轉播工作中，也會發現一些雲端伺服器連線數不多，形同空轉的狀況。此類因伺服器閒置所造成的巨額維運成本也日漸擴大。因此，如何建立自動維運機制才能兼顧使用者觀看品質以及所耗成本最小的彈性伺服器擴充及關閉，已成為一個重要的議題。

雲端伺服器的擴展可以透過垂直擴展 (Vertical scaling) 以及水平擴展 (Horizontal scaling)。垂直擴展是更改伺服器的硬體資源，例如提高中央處理單元 (CPU) / 記憶體 / 頻寬等的等級，而伺服器的數量不變。水平擴展是增減伺服器的數量，

而規格不變，例如透過租賃者預先設定好的範本、伺服器映像檔、或是預設指令腳本，建立許多與標的物同樣規格的虛擬伺服器。目前有些業者需要由租賃者預先將伺服器設為自動擴展(auto-scaling)群組，只有在群組內的伺服器擁有自動擴展功能。有的業者提供服務業者針對不同等級的雲端進行效能評測(benchmarking)。實現方法可採用量測服務的完成時間，來釐清性價比(performance cost ratio)最佳的伺服器等級(instance type)，再藉由訂定政策(policy)實現自動擴展，其政策可基於門檻值觸發、或是固定時間觸發。

現有的伺服器動態增減技術可分成兩類。一類是公共雲業者提供以基礎架構層次(infrastructure-level)為主的反應式(reactive)動態增減，來服務廣大租賃者。此類技術量測目前伺服器的/記憶體/網路使用狀況等，並且有多種指標供租賃者自由選擇。達到門檻值來判斷增減，門檻值可以由使用者(公共雲租賃者)自行設定，或採用預設最佳實務設定。一旦達到其門檻值，透過負載平衡器(load balancer)調配每一伺服器的服務量。另一類是租賃者基於其自身應用的特性，判斷應用層等級(application-level)的服務壓力，透過公共雲業者的編程介面(Application Programming Interface, API)設定企業邏輯，此類大部分是主動式(proactive)技術。技術的參考指標可以是佇列(queue)中待處理資料的數量、平均回應時間、使用者連線數量(number of connections)等。

有一技術提供緊密整合的自動化管理，包括跨雲自動化管理，讓使用者設定各種範本、巨集、腳本等，觀察指標可以排入一陣列，對於增減的邏輯則由租賃者自行判斷。有一技術提出主動式的人工神經網絡訓練的二維矩陣，判斷是否增減伺服器。有一技術認為網頁文件存取有其固定的導覽路線，要找出當中壓力最重的路線進行伺服器擴展。有一技術解決兩層式應用服務，此技術透過一鏈結系統(linkage system)去觀察第一層的反應效能，以決定第二層是否開始擴展(scale-up)。有一技術根據目前虛擬機器(VMs)的總體流量狀態，控制負載平衡器調配負載至其他伺服器。有些技術指出可以根據計費週期來關閉機器。

有一技術考慮違反服務層級協議(Service Level Agreement, SLA)付出的代價與節省經費兩者之間的最佳平衡點。此技術用在多層(multi-tier)的應用，並且基於應用的容量做擴展以及預測系統所需的容量，同時考慮成本模型(cost model)與資源模型(resource model)，所有的要求(requests)都會經由閘道器與負載平衡器。大部份的虛擬機器(VM)具有相同的一般資源配置，其中一部分的虛擬機器具有較低的資源配置。當應用的容量需要擴展(scale up)時，將較低配置的虛擬機器垂直擴展至一般資源配置。當應用的容量需要縮減(scale down)時，進行垂直擴展或水平擴展至較低的資源配

置。

在上述現行的伺服器動態增減技術中，有的技術未評估關閉伺服器後，對於服務提供商的衝擊。有的技術只根據前一台伺服器的狀態，從一群機器中任意選一台關閉。有的技術無法透過負載平衡器來完全控制用戶向誰取得資料。有的技術未充分利用公共雲的特性於節省費用，例如未充分利用不同資料中心的位置與價格並不相同、公共雲的租用計費週期不足 1 小時仍以 1 小時計算、串流服務商可以利用多個公共雲服務商的雲端伺服器等特性。因此，如何建立公共雲的自動維運機制來兼顧服務品質以及所耗成本最小的彈性伺服器擴充與縮減，是值得研究的議題。

【發明內容】

本揭露的實施例可提供一種公共雲資源動態配置方法及系統。

本揭露的一實施例是關於一種公共雲資源動態配置方法。此方法可包含：藉由一負載監視器(Load Monitor)，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身份資訊(Identity Information)，以及該多台伺服器的各伺服器的一目前連線數(current number of connections)、一伺服器等級(level)、以及一所在區域(located

area);藉由一擴展與縮減引擎(Scaling Engine),判斷該多台伺服器中是否有符合至少一觸發條件(trigger condition)的至少一伺服器;藉由該擴展與縮減引擎,將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合(server candidate set);以及藉由該擴展與縮減引擎,接收一性價比資訊,並且根據該伺服器候選者集合,對至少一區域執行一伺服器擴展或縮減程序。

● 本揭露的另一實施例是關於一種公共雲資源動態配置系統。此系統包含一負載監視器、以及一擴展與縮減引擎。此負載監視器取得一目前伺服器配置,該目前伺服器配置至少包括多台伺服器的各台伺服器的身份資訊,以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域。此擴展與縮減引擎判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器;將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合;以及接收一性價比資訊,並且根據該伺服器候選者集合,對至少一區域執行一伺服器擴展或縮減程序。

茲配合下列圖示、實施例之詳細說明及申請專利範圍,將上述及本發明之其他優點詳述於後。

【圖式簡單說明】

第一圖是根據本揭露的一實施例，定義公共雲的租賃費用率的一範例。

第二圖是根據本揭露的一實施例，說明伺服器縮減的觸發時機的一示意圖。

第三圖是根據本揭露的一實施例，說明一種公共雲資源動態配置方法。

第四 A 圖是根據本揭露的一實施例，說明一種公共雲資源動態配置系統。

第四 B 圖是根據本揭露的一實施例，說明第四 A 圖之系統的一應用情境的範例。

第四 C 圖是根據本揭露的一實施例，說明以封包的往返時間來劃分區域的一範例。

第五 A 圖是根據本揭露的一實施例，說明一區域的各伺服器等級對應的每條連線的單位價格的資訊的一範例。

第五 B 圖是根據本揭露的一實施例，說明一區域的各伺服器等級對應的最大連線數的資訊的一範例。

第六圖是根據本揭露的一實施例，說明至少一區域的各區域內的伺服器擴展或縮減的運作流程。

第七圖是根據本揭露的一實施例，說明如何計算一目標配置的運作。

第八 A 圖與第八 B 圖是根據本揭露的一實施例，舉一範例說明一區域內的伺服器擴展或縮減，其中，第八 A 圖是調整前，該區域內各伺服器的狀態資訊；第八 B 圖是調整後，

該區域內各伺服器的狀態資訊。

第九圖將是根據本揭露的一實施例，說明跨區域的伺服器縮減的運作流程。

第十圖是根據本揭露的一實施例，說明 t 值的選擇、與跨區域百分比、節省費用比，之間的關係。

【實施方式】

以下，參考伴隨的圖式，詳細說明依據本揭露的實施例，俾使本領域者易於瞭解。所述之發明創意可以採用多種變化的實施方式，當不能只限定於這些實施例。本揭露省略本領域者已熟知部分(well-known part)的描述，並且相同的參考號於本揭露中代表相同的元件。

依據本揭露的實施例，提供一種公共雲資源動態配置方法及系統。其技術蒐集目前服務在一或多個公共雲所有伺服器的配置狀態，考量對租賃者(向公共雲業者租賃機器者)的服務在公共雲上進行效能測量，從而了解如各等級之伺服器的連線數、以及所在區域等，而一公共雲有至少一伺服器。第一圖是根據本揭露的一實施例，定義公共雲的租賃費用率的一範例。在第一圖的範例中，可依伺服器等級(instance type)定義五種等級(即小、中、大、超大、CPU 增強，分別記為等級 S、等級 M、等級 L、等級 XL、等級 CC2.8XL)的租賃費用率。例如，等級 S 的租賃費用率為每小時 0.060 元，等

級 M 的租賃費用率為每小時 0.120 元，等級 L 的租賃費用率為每小時 0.240 元，等級 XL 的性價比為每小時 0.480 元，等級 CC2.8XL 的性價比為每小時 1.920 元。

租賃者根據這些伺服器的連線數，可計算各等級的伺服器的性價比。租賃者可根據其服務的需求，設定至少一觸發條件，依據本揭露的一實施例，符合觸發條件的伺服器可被加入於一伺服器候選者集合；當符合該觸發條件的情況發生時，可根據輸入的性價比資訊、以及該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序。

依據本揭露的實施例，此至少一觸發條件可被設定為有一伺服器的一或多種運行狀態已達到一門檻值時觸發、以一排程方式於一整點時觸發、有一伺服器已達到距離一計費週期的一結尾的一時間區間內時觸發、一固定時段週期性地觸發，之前述一種或一種以上的觸發條件任意組合。例如，此至少一觸發條件可設定有一伺服器的 CPU、記憶體、頻寬等的所謂的閒置率或資源利用率已達到門檻值時觸發，或是以排程方式於整點觸發，或是有一伺服器接近一計費週期的結尾時觸發，或是每分鐘觸發。而閒置率一般可定義為數值 1 減去資源利用率。

在本揭露中，根據一實施範例，性價比的定義是平均每

條連線所需的單位價格(unit price)。第五 A 圖是根據本揭露的一實施例，定義性價比的一應用範例。在第五 A 圖的範例中，可依伺服器等級(instance type)定義五種等級(即小、中、大、超大、CPU 增強，分別記為等級 S、等級 M、等級 L、等級 XL、等級 CC2.8XL)的性價比，其每條連線的單位價格。例如，等級 S 的性價比為每小時 0.0012 元，等級 M 的性價比為每小時 0.0010 元，等級 L 的性價比為每小時 0.0008 元，等級 XL 的性價比為每小時 0.0006 元，等級 CC2.8XL 的性價比為每小時 0.0024 元。在第五 B 圖的範例中，其中等級 S 的最大連線數為 50 台伺服器，等級 M 的最大連線數為 120 台伺服器，等級 L 的最大連線量為 300 台伺服器，等級 XL 的最大連線數為 800 台伺服器，等級 CC2.8XL 的最大連線數為 800 台伺服器。伺服器例如可以是虛擬機器、主機等的其中一種或一種以上的組合。對於租賃者，各等級的伺服器的性價比需要做效能評測，性價比越高越好。

如之前所述，當判斷出有已符合至少一觸發條件的伺服器時，可根據輸入的性價比資訊，以及伺服器候選者集合進行至少一區域的擴展或縮減程序。擴展伺服器的範例，譬如可以在某一區域增加一台高性價比的伺服器、或是增加一台等級最小的伺服器、或是增加一台等級最大的伺服器、或是增加一台各等級中最大連線數最大的伺服器，然後等待下一次的觸發。縮減伺服器的範例，譬如可將資源利用率較低的

伺服器關閉，或是將低性價比的伺服器關閉，讓使用者分散到其他高性價比的伺服器去。

當使用者隨時間的經過而逐漸減少，閒置的伺服器將因而增加。根據本揭露一實施例，可將低性價比的伺服器關閉，讓使用者分散到其他高性價比的伺服器去，以節省多餘的伺服器的成本花費。擴展或縮減伺服器的觸發的時間點，譬如可以採用如 CPU、記憶體、頻寬等的閒置率已達到門檻值（例如，以 CPU 的閒置率(idle rate)為 80%與 20%分別作為上限門檻值與下限門檻值）時觸發，或是以排程方式於整點觸發，或是有任何一台伺服器接近計費週期結尾時觸發，或是每分鐘觸發。觸發時可以考慮將目前所有的伺服器全部列入伺服器候選者集合、或是考慮將該伺服器是否已接近其計費週期的結尾才列入伺服器候選者集合。第二圖是根據本揭露的一實施例，說明伺服器縮減的觸發時機的一示意圖，其中一伺服器的一計費週期如標號 210 所示。

在第二圖中，係考慮將一或多台使用中已接近其計費週期(billing cycle)結尾的伺服器列入要被關閉的候選者(reducing candidate)集合，其實施方式例如可設定一門檻值 t ，並且將離計費週期 t 分鐘內即將完成一計費週期的一或多台伺服器列入伺服器候選者集合。在第二圖的範例中，根據此門檻值 t ，伺服器 A、伺服器 C、以及伺服器 D 都是接

近其計費週期結尾的伺服器候選者。因此，伺服器 A、伺服器 C、以及伺服器 D 也可以觸發伺服器縮減(server reduction)。也就是說，根據本揭露的實施例，可採用條件式觸發來產生伺服器擴展或縮減程序。

第三圖是根據本揭露的一實施例，說明一種公共雲資源動態配置方法。參考第三圖，此方法可包含：藉由一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身份資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域(步驟 310)；藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器(步驟 320)；藉由該擴展與縮減引擎，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合(步驟 330)；以及藉由該擴展與縮減引擎，接收一性價比資訊，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序(步驟 340)。挑選自該目前伺服器配置中的該至少一伺服器的該伺服器候選者集合，其中也包括了每一伺服器的身份資訊、一目前連線數、一伺服器等級、以及一所在區域等資訊。

依此，根據本揭露的一實施例，一種公共雲資源動態配置系統 400 可如第四 A 圖所示。系統 400 可包含一負載監視器 410、以及一擴展與縮減引擎 420。此負載監視器 410 取

得一目前伺服器配置 412，該目前伺服器配置至少包括多台伺服器的各伺服器的身份資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域。此擴展與縮減引擎 420 判斷該至少一伺服器中是否有符合至少一觸發條件的至少一伺服器；將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合 422；以及接收一性價比資訊 424，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序 426。挑選自該目前伺服器配置中的該至少一的伺服器候選者集合，此其中也包括了每一伺服器的身份資訊、一目前連線數、一伺服器等級、以及一所在區域等資訊。

第四 B 圖是根據本揭露的一實施例，說明第四 A 圖之系統的一應用情境的範例。在第四 B 圖的範例中，負載監視器 410 可取得一或多個公共雲上的一目前伺服器配置，此目前伺服器配置例如是位於多個不同區域(例如新加坡、日本、美國、巴西、...)的多台伺服器的目前狀態資訊，此狀態資訊包括至少此多台伺服器的每一伺服器的身分資訊、目前連線數、伺服器等級、以及所在區域等的狀態資訊。身分資訊可為例如是一伺服器代號，用以區分不同的伺服器。擴展與縮減引擎 420 從負載監視器 410 取得這些狀態資訊，當此多台伺服器中有符合觸發條件者(例如位於新加坡的伺服器)，擴展與縮減引擎 420 可對位於此區域(新加坡)的伺服器

可藉由，但不限定是發出一或多個伺服器擴展或縮減指令 (scaling commands)⁴³⁰，以執行伺服器擴展或縮減程序⁴²⁶，將成本效益較低的伺服器關閉，令使用者分散到其他成本效益較高的伺服器去。其中縮減指令例如是「aws ec2 terminate-instances」。其中擴展指令例如是「aws ec2 run-instances」、「aws ec2 terminate-instances」、「aws ec2 modify-instance-attribute」這三種的其中之一或二或三種的任意組合。根據本揭露的實施例，公共雲資源動態配置系統⁴⁰⁰可在單一公共雲上運行，也可以跨越在多個公共雲上運行。

本揭露所謂的「區域 (area)」，可以是以地理位置 (geographical location) 來劃分的區域、或是以封包的往返時間 (Round Trip Time, RTT) 來劃分的區域。第四 C 圖是根據本揭露的一實施例，說明以封包的往返時間來劃分區域的一範例。在第四 C 圖的範例中，有六個在不同所在位置的雲端中心 (記為雲端中心 431~雲端中心 436)，其中雲端中心 431~雲端中心 433 的各雲端中心的封包的往返時間皆小於等於 120 毫秒 (即 $RTT \leq 120ms$)，而雲端中心 434~雲端中心 436 的各雲端中心的封包的往返時間皆小於等於 500 毫秒且大於等於 120 毫秒 (即 $120ms < RTT \leq 500ms$)，依此，雲端中心 431~雲端中心 433 被劃分在區域 441，而雲端中心 434~雲端中心 436 被劃分在區域 442。

根據本揭露的實施例，性價比資訊至少包含該至少一區域的各區域的各伺服器等級對應的每條連線的單位價格的資訊、以及該至少一區域的各區域的各伺服器等級對應的最大連線數的資訊。第五 A 圖是根據本揭露的一實施例，說明一區域的各伺服器等級對應的每條連線的單位價格的資訊的一範例。第五 A 圖的範例說明了並非越高等級的伺服器的單位成本越便宜，可由租賃者自行進行各等級的效能評測，例如租用最貴的叢集 CPU 等級的伺服器可能對於多媒體的應用毫無幫助，其性價比會非常低。一般而言，因為頻寬的關係會在較高伺服器等級如 L、XL 等級得到較高的性價比。某些服務消耗記憶體非常大，此時可以選針對記憶體優化的伺服器等級的性價比較高。第五 B 圖是根據本揭露的一實施例，說明該區域的各伺服器等級對應的最大連線數的資訊的一範例。

根據本揭露的一實施例，伺服器擴展或縮減程序可以分為兩階段，第一階段是區域內(inter-area)的伺服器擴展或縮減，第二階段是跨區域(intra-area)的伺服器縮減。也就是說，當有符合至少一觸發條件的伺服器時，先對該至少一區域的各區域內執行一伺服器擴展或縮減後，再執行一跨區域的伺服器縮減。根據本揭露的實施例，此兩階段的伺服器擴展或縮減程序，第一階段在不造成跨區域連線的前提下，先把

所有區域的每一區域內各自的伺服器運行成本縮減到最低，以減少大部分的跨區域連線，讓大部分的使用者都能經由同區域的伺服器提供連線，第二階段的伺服器縮減可能造成少部分的使用者必須由跨區域的伺服器提供連線。此伺服器擴展或縮減程序從而能夠在節省伺服器成本以及滿足使用者品質(減少跨區域連線)上達成平衡。

第六圖是根據本揭露的一實施例，說明至少一區域的各區域內的伺服器擴展或縮減的運作流程。參考第六圖，擴展與縮減引擎 420 接收一性價比資訊，此性價比資訊至少包含該至少一區域的各區域內各伺服器等級各自對應的每條連線的單位價格的資訊、以及該至少一區域的各區域內各伺服器等級各自對應的最大連線數的資訊（步驟 610）；根據此性價比資訊，計算一目標配置，從而產生該至少一區域的各區域內各伺服器等級各自對應的一伺服器數量（步驟 620）；以及發出一或多個伺服器擴展或縮減指令，調整該至少一區域的各區域中各伺服器等級對應的伺服器數量至該目標配置中各伺服器等級各自對應的伺服器數量（步驟 630）。當需要從多個相同等級的伺服器中關閉其中至少一伺服器時，可優先考量，但不限定是，關閉該多個相同等級的伺服器中目前連線數最少的伺服器。

第七圖是根據本揭露的一實施例，說明如何計算一區域

的一目標配置的運作。參考第七圖，擴展與縮減引擎 420 將該伺服器候選者集合中該區域中所有伺服器的目前連線數的總合做爲一未分派連線數（步驟 710）；並且根據該區域內各伺服器等級各自對應的每條連線的單位價格、該區域內各伺服器等級各自對應的最大連線數、以及該未分派連線數，分配該區域內各伺服器等級各自對應的一目標伺服器數量（步驟 720）。一伺服器等級對應的每條連線單位價格越低，其性價比越高。計算一伺服器等級對應的該目標伺服器數量有多種方式，以下的公式是其中的一個範例。

一伺服器等級對應的目標伺服器數量

= 該未分派連線數 / 該伺服器等級對應的最大連線數；

以及，更新該未分派連線數如下：

該未分派連線數

= 該未分派連線數 Mod 該伺服器等級對應的最大連線數；

其中，Mod 是一模數運算。

在步驟 720 中，有多種實施方式可分配該區域內各伺服器等級各自對應的該目標伺服器數量。例如根據一實施例，可由該區域內多台伺服器等級對應的一最低的單位價格至一最高的單位價格高，依序地分配該區域內各伺服器等級各自對應的該目標伺服器數量。假設將距離一計費週期(60 分鐘)結束 t 分鐘內的伺服器加入一伺服器候選者集合，或將

所有伺服器皆加入關閉的伺服器候選者集合(即 $t=60$)。則一區域內的伺服器擴展或縮減程序可運作如下。加總該伺服器候選者集合內所有伺服器的連線數做爲一未分派連線數。依序從性價比高(伺服器等級對應的每條連線單位價格最低)的伺服器等級開始分配連線數。例如，XL 等級的伺服器其性價比最高並且假設最多可以支援 800 條連線，則先分配 $\lfloor \text{未分派連線數}/800 \rfloor$ 台 XL 等級的伺服器。分配後，將該未分派連線數更新爲 $\lfloor \text{未分派連線數} \bmod 800 \rfloor$ 。當更新後的未分派連線數尚未歸零時，再繼續分配下一等級伺服器的目標伺服器數量，直到該未分派連線數成爲零。若該未分派連線數小於該伺服器等級對應的最大連線數，該目標伺服器數量加 1。欲積極節費的租賃者可調整公式爲放棄該未分派連線數，使用該目標伺服器數量。有多種實施方式可在此進行微調，仍不違背由性價比高的伺服器開始分配之精神。此時已完成一區域的目標配置(包含該區域內各伺服器等級對應的伺服器數量)。根據該目標配置與該區域內目前的伺服器配置數量上的差異進行調整，此時可能會增加或減少各種等級的伺服器。當需要增加伺服器時，可直接增加;當需要關閉伺服器時，可採用，但不限定於，一最小編輯距離(minimum edit distance; Levenshtein)爲原則來進行伺服器數量的調整，其依據爲目前使用該伺服器的連線數。舉例來說，若有兩台同樣是 XL 等級的伺服器要關閉其中一台伺服器，此時可選擇目前連線數較少的那台伺服器。

依據上述的實施例，第八 A 圖與第八 B 圖舉一範例說明一區域內的伺服器擴展或縮減程序，其中，假設一伺服器候選者集合中一區域內總共有 1628 使用者之連線。第八 A 圖是調整前，該區域內各伺服器的狀態資訊。租賃者經過效能評測後，認為 XL 等級伺服器之性價比較高，優先將連線數分派給 XL 等級的伺服器，並且根據上述目標配置的運作流程及求得目標伺服器數量的公式範例，計算出該區域內的目標配置是 2 台 XL 等級的伺服器、以及 1 台 S 等級的伺服器。

根據此目標配置與該區域內目前的伺服器配置數量上的差異，因此，應關閉一台 XL 等級的伺服器、一台 L 等級的伺服器、以及一台 S 等級的伺服器。縮減伺服器時，可考慮同等級伺服器中具有最小編輯距離者，例如，目前的 XL 等級的伺服器共有三台可選，可從中選擇關閉目前連線人數最低的 XL 等級的伺服器，因而關閉伺服器代號為 i-PSRHEDNF 的 XL 等級的伺服器(XL 等級的伺服器中目前連線數最低者)、伺服器代號為 i-PHAQQQYT 的 L 等級的伺服器、以及伺服器代號為 i-KGMUCWEE 的 S 等級的伺服器(S 等級的伺服器中目前連線數最低者)，如第八 B 圖所示之調整後，該區域內各伺服器的狀態資訊，其中刪除線表示關閉該伺服器。

根據本揭露的一實施例，第二階段的跨區域的伺服器縮減係依據伺服器候選者集合 422 中所有伺服器的閒置率或資源利用率以進行縮減，譬如可依照這些伺服器的閒置率由高至低排序或資源利用率由低至高排序，依序進行縮減。一伺服器的資源利用率計算方法，其中的一個範例如以下的公式：

資源利用率 = 該伺服器的目前連線數與該伺服器對應的伺服器等級所對應的最大連線數的比值。

第九圖將是根據本揭露的一實施例，說明跨區域的伺服器縮減的運作流程。

參考第九圖，擴展與縮減引擎 420 計算一服務容量與一目前總連線數，其中服務容量=該伺服器候選者集合中所有伺服器的伺服器等級對應的最大連線數的總合，目前總連線數=該伺服器候選者集合中所有伺服器的目前連線數的總合(步驟 910);依照該伺服器候選者集合中所有伺服器的閒置率由高至低排序(步驟 920);然後，從閒置率最高的一伺服器開始依次進行判定，當該服務容量與該伺服器的伺服器等級對應的最大連線數相減後的差大於等於該目前總連線數時，擴展與縮減引擎 420 判定關閉該伺服器(步驟 930)。當該服務容量與該伺服器的伺服器等級對應的最大連線數相減後的

差小於該目前總連線數時，擴展與縮減引擎 420 判定不關閉該伺服器(步驟 940)。直到該伺服器候選者集合中不再有伺服器可以被關閉。

也就是說，跨區域的伺服器縮減可依據該伺服器候選者集合中所有伺服器的伺服器等級對應的最大連線數的總合、該伺服器候選者集合中所有伺服器的目前連線數的總合、以及各伺服器等級所對應的最大連線數，判定是否關閉該伺服器。

根據本揭露實施例的公共雲資源動態配置技術，於第二階段經過跨區域縮減後才會產生跨區域連線，若租賃者不希望產生任何跨區域連線，可以設定擴展與縮減引擎 420 不執行跨區域的伺服器縮減階段，但是獲得較差的節費效果。第十圖是根據本揭露的一實施例，說明 t 值的選擇、與跨區域百分比、節省費用比，之間的關係。其中，橫軸代表 t 值(單位:分鐘)，縱軸代表百分比。曲線 1010 代表觸發時不考慮 t 值而將所有伺服器全部列入伺服器候選者集合的一種原始方法所產生的跨區域百分比，曲線 1020 代表只將距離計費週期結尾 t 分鐘內的伺服器列入伺服器候選者集合的跨區域百分比，曲線 1030 代表該原始方法的節省費用比，曲線 1040 代表考慮 t 值時的節省費用比。

參考第十圖，從曲線 1040 可以看出， t 值的選擇越高，跨區域的伺服器縮減所產生的節省費用效果越強；其代價是所產生的跨區域連線數也越高。若 t 值設為 60 分鐘表示所有伺服器都被列入考慮關閉的伺服器候選者集合即等同於該原始方法。假如 t 值選擇為 5 分鐘，則節省費用效果很差，若 t 值增加為 10 分鐘，則節省費用效果很明顯提升近 1 倍。當 t 值選擇為 35 分鐘以上開始出現節省費用的邊際效益遞減。

綜上所述，依據本揭露的實施例提供一種公共雲資源動態配置方法及系統。其技術利用一負載監視器，取得公共雲上的一目前伺服器配置，提供給一擴展與縮減引擎。此擴展與縮減引擎採用條件觸發式產生伺服器縮減事件，並且可動態調整各等級伺服器的目標伺服器數量，從而降低伺服器的運行成本並維持租賃者的服務品質。此技術可在單一公共雲上運行，也可以跨越在多個公共雲上運行。

以上所述者僅為依據本揭露的實施範例，當不能依此限定本揭露實施之範圍。即大凡發明申請專利範圍所作之均等變化與修飾，皆應仍屬本揭露專利涵蓋之範圍。

【符號說明】

S、M、L、XL、CC2.8XL 伺服器等級

t 門檻值

210 一計費週期

A、C、D 候選者伺服器

310 藉由一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域

320 藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器

330 藉由該擴展與縮減引擎，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合

340 藉由該擴展與縮減引擎，接收一性價比資訊，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序

400 公共雲資源動態配置系統

410 負載監視器

420 擴展與縮減引擎

422 伺服器候選者集合

424 性價比資訊

426 伺服器擴展或縮減程序

412 目前伺服器配置

430 伺服器擴展或縮減指令

610 接收性價比資訊，此性價比資訊至少包含該至少一區域的各區域內各伺服器等級各自對應的每條連線的單位價格的資訊、以及該區域內各伺服器等級各自對應的最大連線數的資訊

620 根據此性價比資訊，計算一目標配置，從而產生該至少一區域的各區域內各伺服器等級各自對應的一伺服器數量

630 發出一或多個伺服器擴展或縮減指令，調整該至少一區域的各區域中各伺服器等級對應的伺服器數量至該目標配置中各伺服器等級各自對應的伺服器數量

710 將該伺服器候選者集合中該區域中所有伺服器的目前連線數的總合做為一未分派連線數

720 根據該區域內各伺服器等級各自對應的每條連線的單位價格、該區域內各伺服器等級各自對應的最大連線數、以及該未分派連線數，分配該區域內各伺服器等級各自對應的一目標伺服器數量

910 計算一服務容量與一目前總連線數，其中服務容量=該伺服器候選者集合中所有伺服器的伺服器等級對應的最大連線數的總合，目前總連線數=該伺服器候選者集合中所有伺服器的目前連線數的總合

920 依照該伺服器候選者集合中所有伺服器的閒置率由高至低排序

930 從閒置率最高的伺服器開始，當該服務容量與該伺服器的伺服器等級對應的最大連線數相減後的差大於等於該目前總連線數時，判定關閉該伺服器

940 當該服務容量與該伺服器的伺服器等級對應的最大連線數相減後的差小於該目前總連線數時，判定不關閉該伺服器

1010 曲線，代表原始方法所產生的跨區域百分比

1020 曲線，代表考慮 t 值的跨區域百分比

1030 曲線，代表原始方法的節省費用比

● 1040 曲線，代表考慮 t 值的節省費用比

申請專利範圍

1. 一種公共雲資源動態配置方法，包含：

藉由一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域；

藉由一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器；

藉由該擴展與縮減引擎，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合；以及

藉由該擴展與縮減引擎，接收一性價比資訊，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序。

2. 如申請專利範圍第 1 項所述之方法，其中該性價比資訊至少包括該至少一區域的各區域內各伺服器等級各自對應的每條連線的單位價格的資訊、以及該至少一區域的各區域內各伺服器等級各自對應的最大連線數的資訊。

3. 如申請專利範圍第 1 項所述之方法，其中執行該伺服器擴展或縮減程序是先對該至少一區域的各區域內執行一伺服器擴展或縮減後，再執行一跨區域的伺服器縮減。

4. 如申請專利範圍第 1 項所述之方法，其中該至少一觸發條件被設定為有一伺服器的一或多種運行狀態已達到一門檻值時觸發、以一排程方式於一整點時觸發、有一伺

服务器已達到距離一計費週期的一結尾的一時間區間內時觸發、一固定時段週期性地觸發，之前述一種或一種以上的觸發條件任意組合。

5. 如申請專利範圍第 2 項所述之方法，其中該方法還包括：根據該性價比資訊，計算一目標配置，從而產生該至少一區域的各區域內各伺服器等級各自對應的一伺服器數量；以及

發出一或多個伺服器擴展或縮減指令，調整目前該至少一區域的各區域中各伺服器等級對應的伺服器數量至該目標配置中各伺服器等級各自對應的伺服器數量。

6. 如申請專利範圍第 5 項所述之方法，其中計算該目標配置還包括：

將該伺服器候選者集合中該至少一區域的各區域中所有伺服器的目前連線數的總合做為一未分派連線數；以及

根據該至少一區域的各區域內各伺服器等級各自對應的每條連線的單位價格、該至少一區域的各區域內各伺服器等級各自對應的最大連線數、以及該未分派連線數，分配該至少一區域的各區域內各伺服器等級各自對應的一目標伺服器數量。

7. 如申請專利範圍第 6 項所述之方法，其中該方法係由該至少一區域的各區域內各伺服器等級各自對應的一最低的每條連線的單位價格至一最高的每條連線的單位價格，依序地分配該至少一區域的各區域內各伺服器等級

各自對應的該目標伺服器數量。

8. 如申請專利範圍第 1 項所述之方法，其中當需要從多個相同等級的伺服器中關閉其中至少一伺服器時，被關閉的該至少一伺服器係該多個相同等級的伺服器中目前連線數最少的伺服器。
9. 如申請專利範圍第 3 項所述之方法，其中該跨區域的伺服器縮減係將該伺服器候選者集合中所有伺服器依據該些伺服器的一閒置率或一資源利用率以進行縮減。
10. 如申請專利範圍第 9 項所述之方法，其中該閒置率是數值 1 減去該資源利用率，該資源利用率是該伺服器的一目前連線數與該伺服器對應的一伺服器等級所對應的一最大連線數的比值。
11. 如申請專利範圍第 3 項所述之方法，其中該跨區域的伺服器縮減係依據該伺服器候選者集合中所有伺服器的伺服器等級各自對應的最大連線數的總合、該伺服器候選者集合中所有伺服器的目前連線數的總合、以及一伺服器的伺服器等級所對應的最大連線數，判定是否關閉該伺服器。
12. 一種公共雲資源動態配置系統，包含：
 - 一負載監視器，取得一目前伺服器配置，該目前伺服器配置至少包括多台伺服器的各伺服器的身分資訊，以及該多台伺服器的各伺服器的一目前連線數、一伺服器等級、以及一所在區域；以及

一擴展與縮減引擎，判斷該多台伺服器中是否有符合至少一觸發條件的至少一伺服器，將符合該至少一觸發條件的該至少一伺服器加入一伺服器候選者集合；以及接收一性價比資訊，並且根據該伺服器候選者集合，對至少一區域執行一伺服器擴展或縮減程序。

13. 如申請專利範圍第 12 項所述之系統，其中當該至少一伺服器中有符合該至少一觸發條件的該至少一伺服器時，該擴展與縮減引擎對位於該至少一區域的該至少一伺服器發出一或多個伺服器擴展或縮減指令，以執行該伺服器擴展或縮減程序。
14. 如申請專利範圍第 12 項所述之系統，其中該伺服器擴展或縮減程序分為兩階段，其中第一階段是區域內的伺服器擴展或縮減，第二階段是跨區域的伺服器縮減。
15. 如申請專利範圍第 12 項所述之系統，其中該至少一觸發條件被設定為有一伺服器的一或多種運行狀態已達到一門檻值時觸發、以一排程方式於一整點時觸發、有一伺服器已達到距離一計費週期的一結尾的一時間區間內時觸發、一固定時段週期性地觸發，之前述一種或一種以上的觸發條件任意組合。
16. 如申請專利範圍第 12 項所述之系統，其中該擴展與縮減引擎從該負載監視器取得該目前伺服器配置的資訊。
17. 如申請專利範圍第 12 項所述之系統，其中該性價比資訊至少包括該至少一區域的各區域內各伺服器等級各

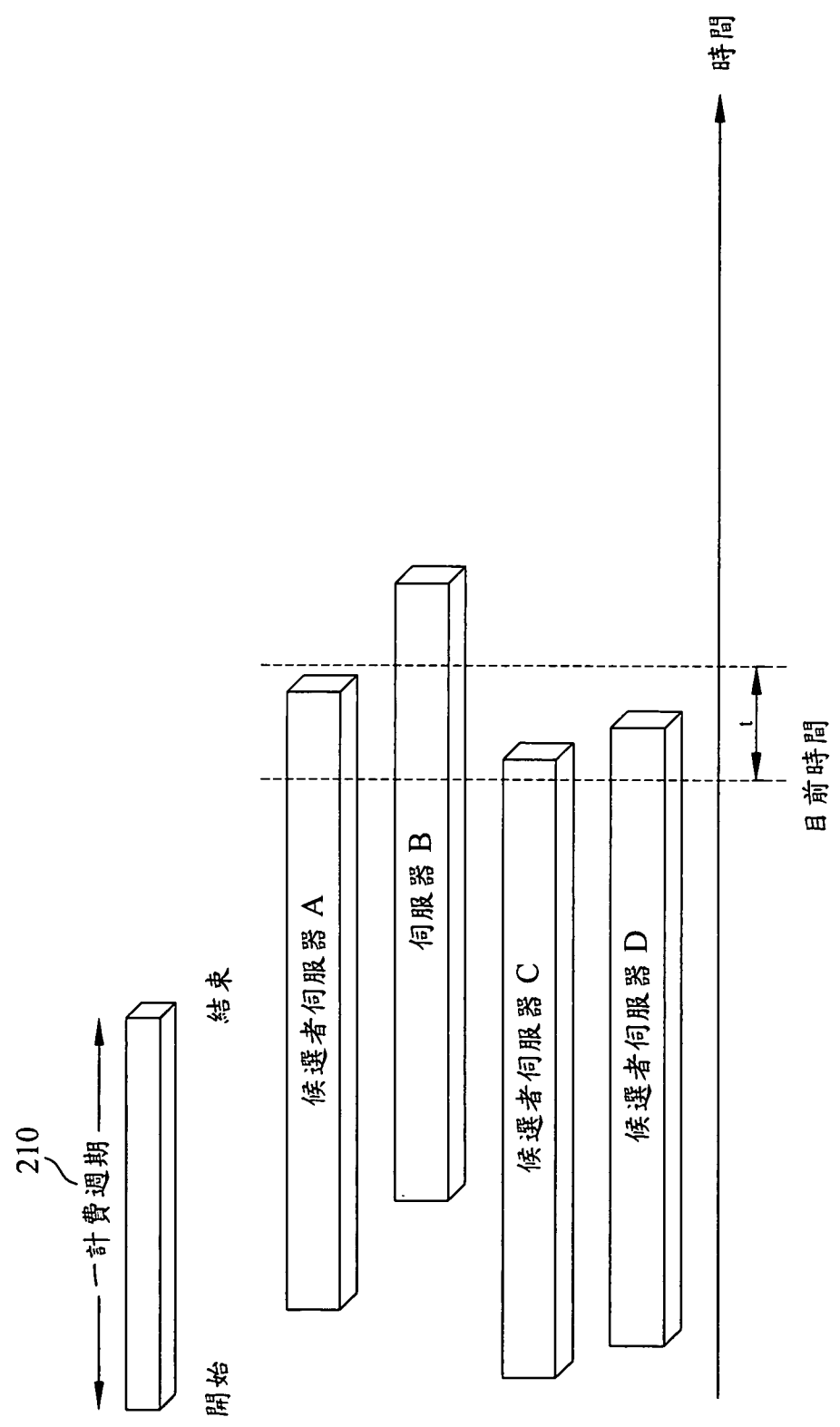
自對應的每條連線的單位價格的資訊、以及該至少一區域的各區域內各伺服器等級各自對應的最大連線數的資訊。

18. 如申請專利範圍第 12 項所述之系統，其中該至少一伺服器是至少一虛擬機器以及至少一主機，的其中一種或一種以上的組合。
19. 如申請專利範圍第 12 項所述之系統，其中該系統係在一或多個公共雲上運行。

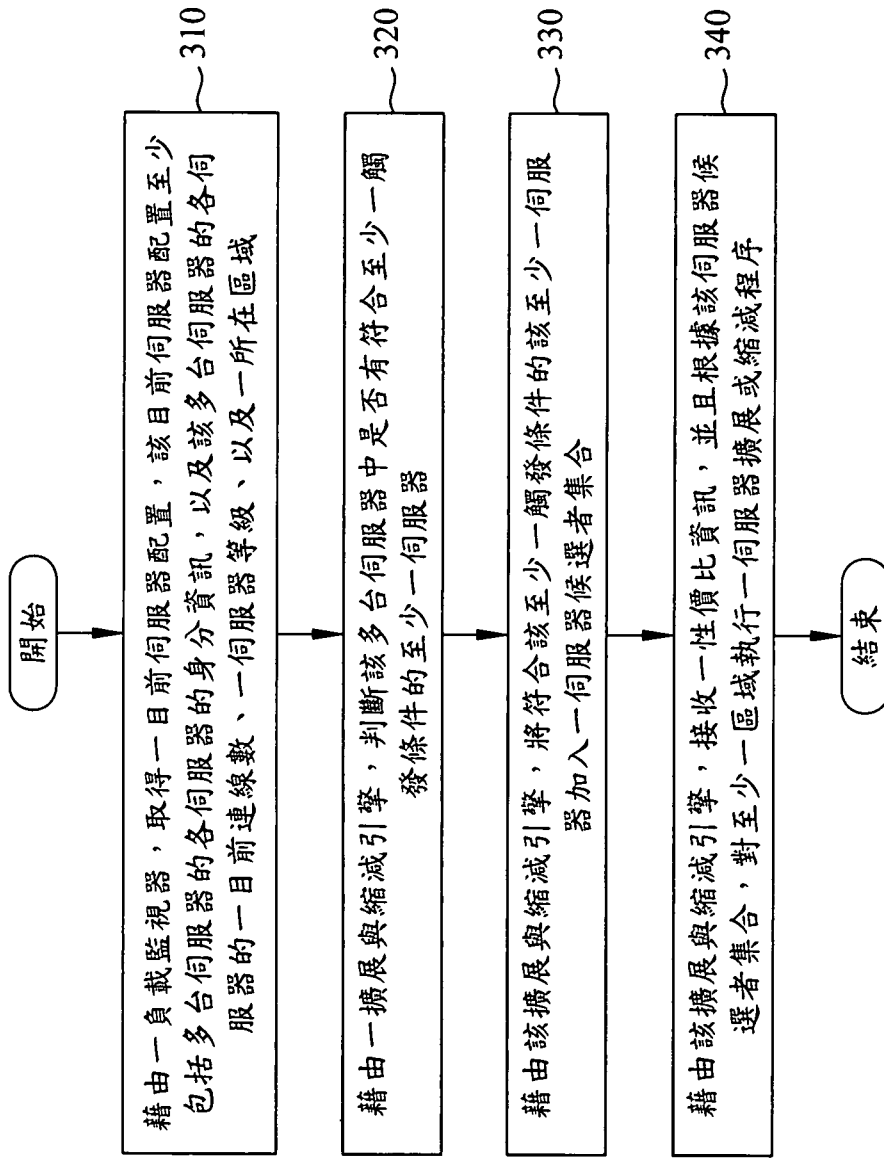
圖式

等級S	\$0.060每小時
等級M	\$0.120每小時
等級L	\$0.240每小時
等級XL	\$0.480每小時
等級CC2.8XL	\$1.920每小時

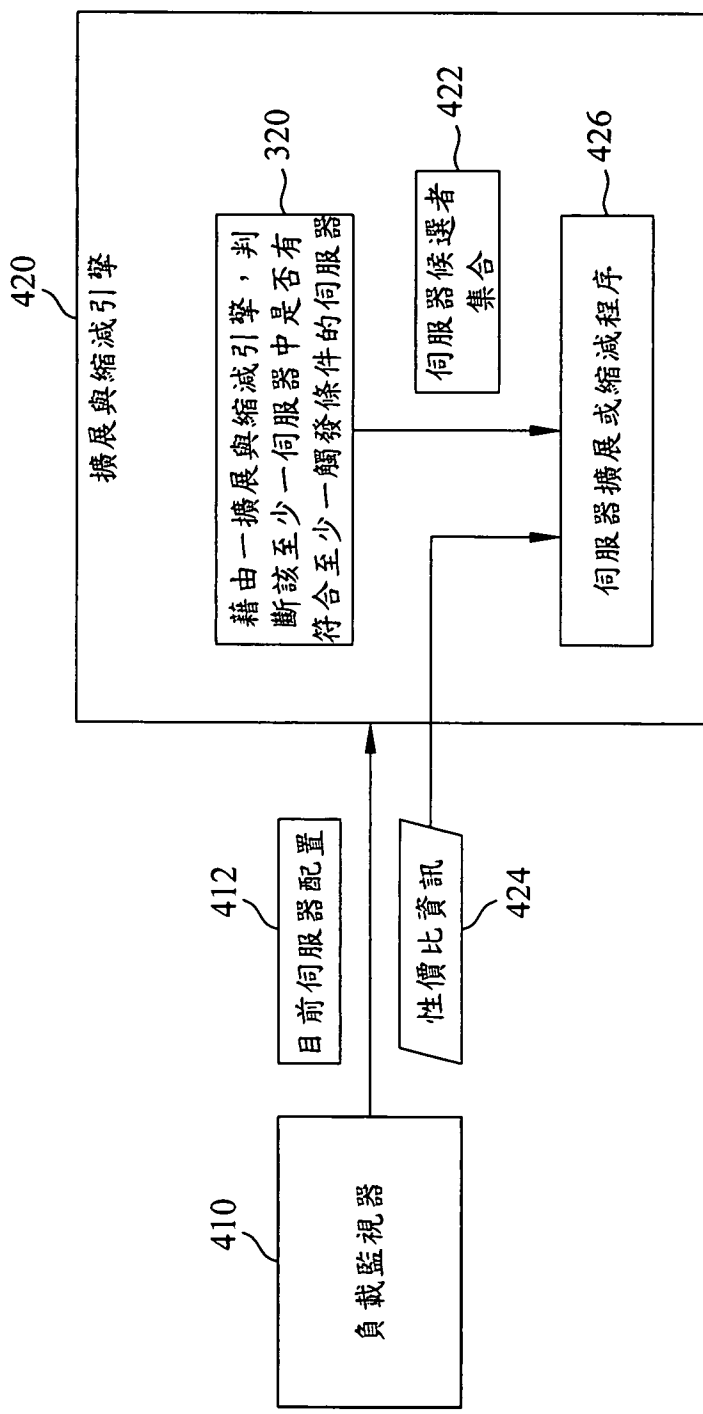
第一圖



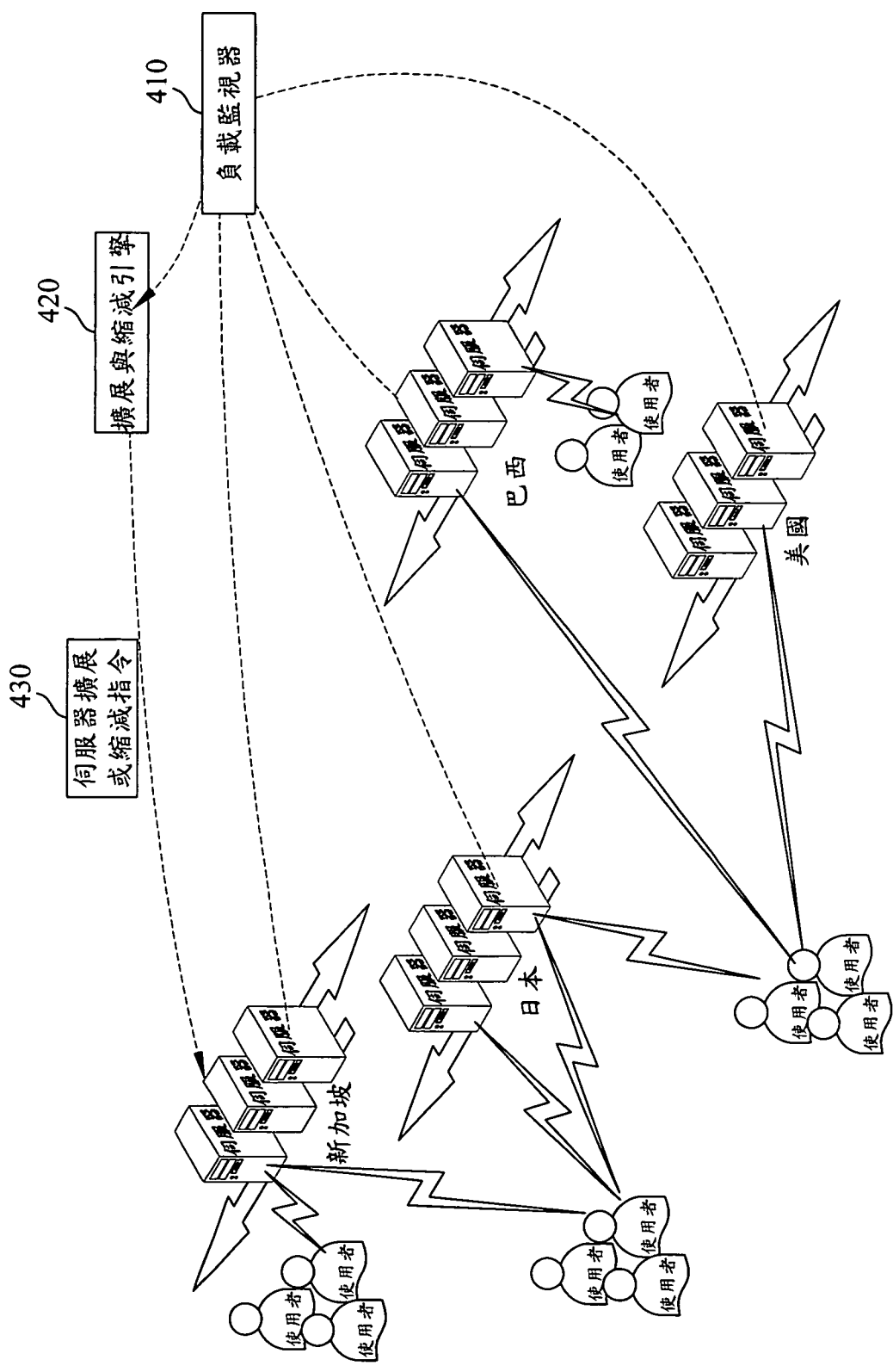
第二圖



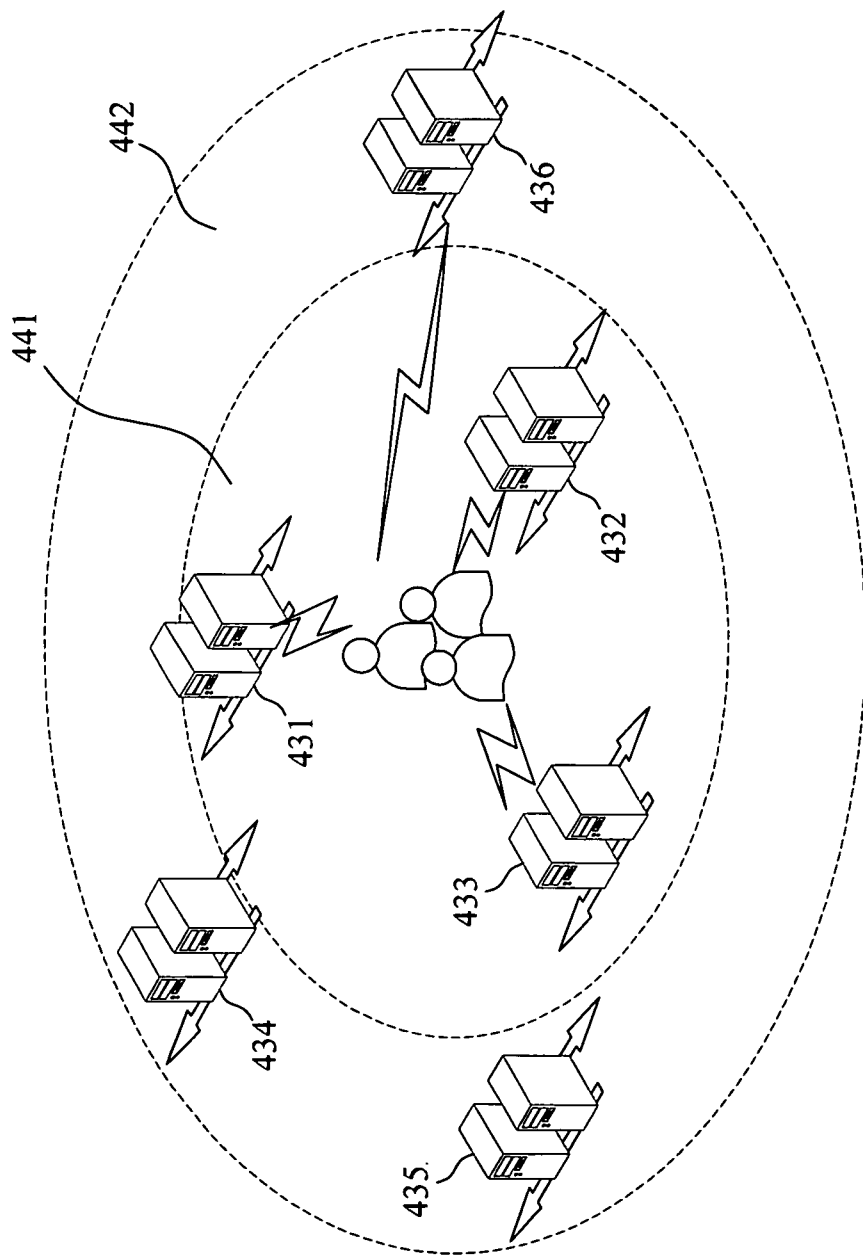
第三圖



第四A圖



第四B圖



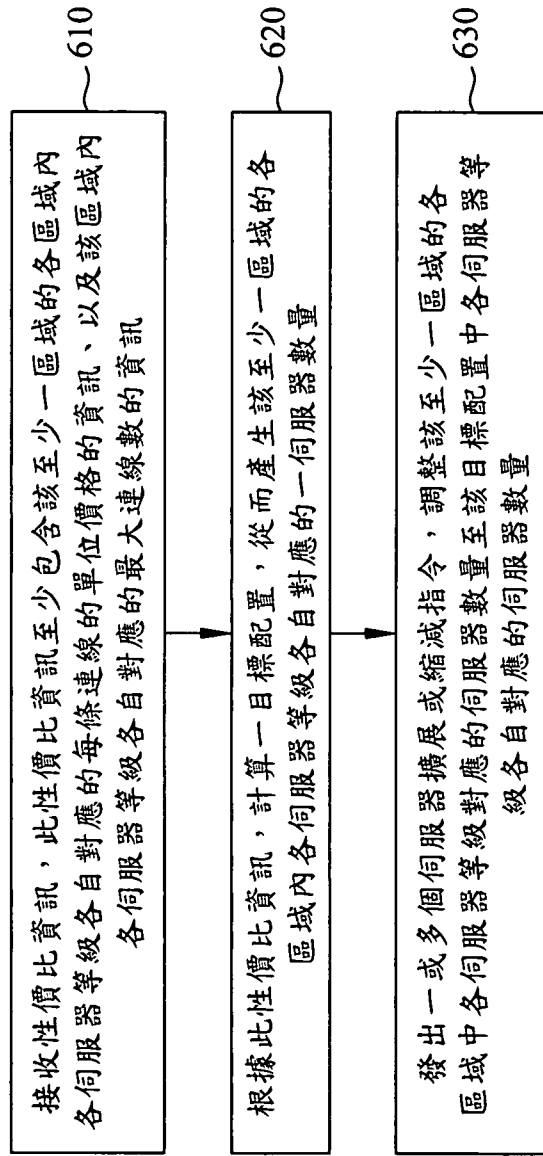
第四C圖

伺服器等級	每條連線的單位價格(每小時)
S	\$0.0012
M	\$0.0010
L	\$0.0008
XL	\$0.0006
CC2.8XL	\$0.0024

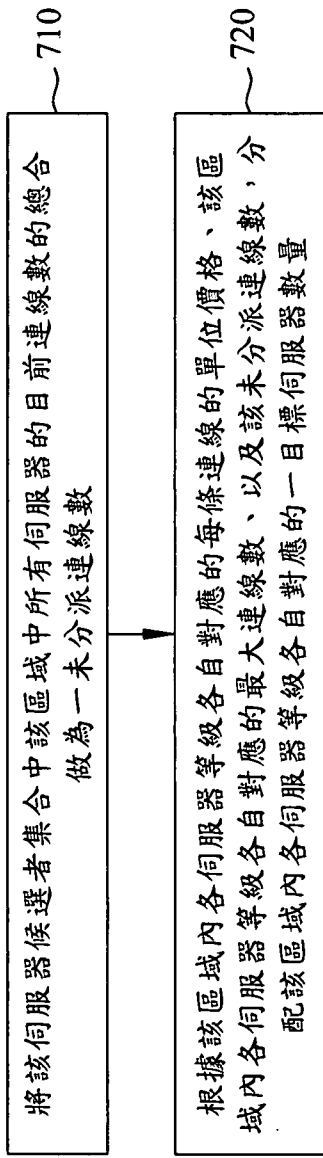
第五A圖

伺服器等級	最大連線數
S	50
M	120
L	300
XL	800
CC2.8XL	800

第五B圖



第六圖



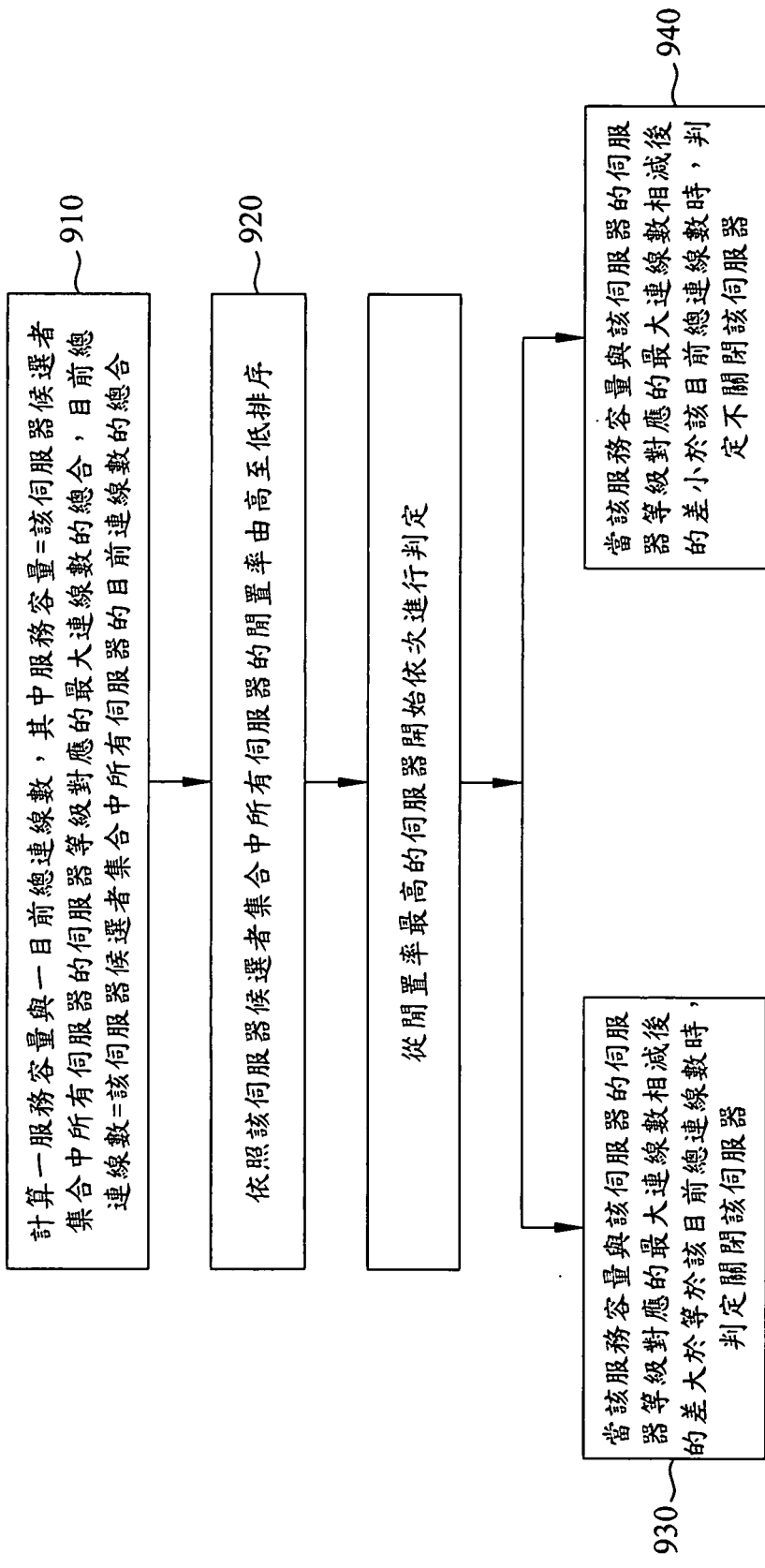
第七圖

伺服器代號	伺服器等級	最大連線數	目前連線數	閒置率	區域代號
i-FOHLBLOQ	S	50	46	0.08	2
i-KGMUCWEE	S	50	30	0.4	2
i-PHAQQQYT	L	300	1	1	2
i-PSRHEDNF	XL	800	134	0.83	2
i-NVVPRYUI	XL	800	773	0.03	2
i-HHIGBFQV	XL	800	644	0.2	2

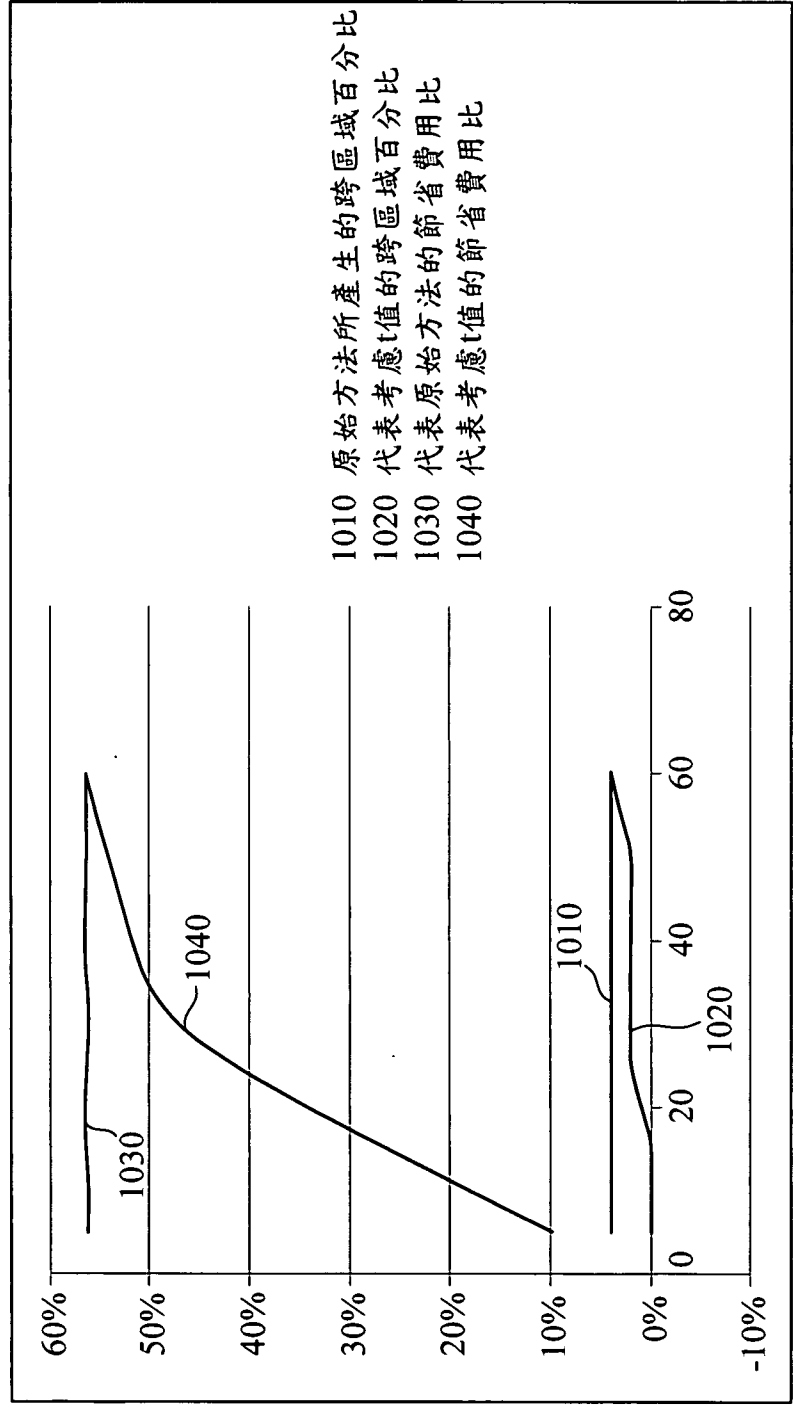
第八A圖

伺服器代號	伺服器等級	最大連線數	目前連線數	閒置率	區域代號
i-FOHLBLOQ	S	50	46	0.08	2
i-KGMUCWEE	S	50	30	0.4	2
i-PHAQQQYT	L	300	1	1	2
i-PSRHEDNF	XL	800	134	0.83	2
i-NVVPRYUI	XL	800	773	0.03	2
i-HHIGBFQV	XL	800	644	0.2	2

第八B圖



第九圖



第十圖