



(12) 发明专利申请

(10) 申请公布号 CN 103827858 A

(43) 申请公布日 2014. 05. 28

(21) 申请号 201180073820. 8

(22) 申请日 2011. 09. 30

(30) 优先权数据

1116737. 6 2011. 09. 28 GB

(85) PCT国际申请进入国家阶段日

2014. 03. 28

(86) PCT国际申请的申请数据

PCT/EP2011/067171 2011. 09. 30

(87) PCT国际申请的公布数据

W02013/044987 EN 2013. 04. 04

(71) 申请人 瑞典爱立信有限公司

地址 瑞典斯德哥尔摩

(72) 发明人 E. 弗里曼 A. 阿维德斯森

L. 维斯特伯格

(74) 专利代理机构 中国专利代理(香港)有限公司 72001

代理人 杨美灵 汤春龙

(51) Int. Cl.

G06F 17/30(2006. 01)

H04L 29/08(2006. 01)

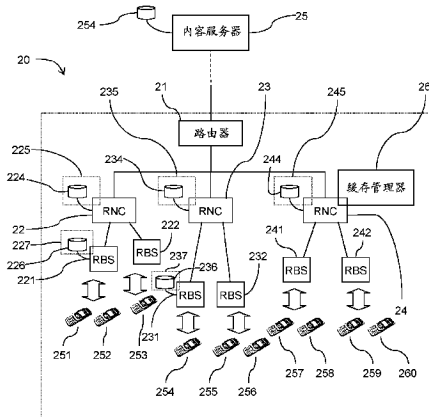
权利要求书3页 说明书12页 附图6页

(54) 发明名称

移动网络中的缓存

(57) 摘要

本文描述了一种用于优化在资源有限网络的缓存域中缓存之间数据对象的分布的方法。在缓存域中的缓存接收对数据对象的用户请求。将通知从收到请求的每个缓存通知到缓存管理器。通知报告用户请求并且识别请求的数据对象。在缓存管理器,整理和存储对象信息,对象信息包括每个请求的数据对象的请求频率和收到请求的缓存的位置。在缓存管理器,基于对象信息识别缓存域内要分布的对象。将在缓存之间分布在那些缓存中存储的数据对象的指示从缓存管理器发送到缓存。根据受欢迎度将对象分类到几个类中,包括具有应分布到缓存域中所有缓存的对象的高受欢迎度类、具有应分布到缓存域中缓存的子集的对象的中等受欢迎度类及具有不应分布的对象低受欢迎度类。



1. 一种用于控制在分组数据网络中缓存域中的缓存(225,235,245,227,237)中存储的内容的分布的缓存管理器(26),所述缓存管理器包括:

对象监视器(55),用于接收来自所述缓存域中缓存的有关在那些缓存请求的数据对象的通知;

对象数据库(52),用于整理和存储有关所述请求的数据对象的对象信息,所述对象信息包括有关已请求所述数据对象的所述缓存和提出所述请求的频率的信息;以及

对象分布器(56),用于基于所述对象信息识别要分布的对象并且指示所述缓存域中的所述缓存在它们之间分布在那些缓存中存储的数据对象;

其中所述对象信息包括用于每个对象的受欢迎度参数,并且所述对象分布器配置成将所述对象分类到至少三个类之一中,所述类包括:

高受欢迎度类,包括应分布到所述缓存域中所有缓存的对象;

中等受欢迎度类,包括应分布到所述缓存域中所述缓存的一个或更多个子集的对象;

以及

低受欢迎度类,包括不应分布的对象。

2. 如权利要求1所述的缓存管理器,还包括:

拓扑监视器(55),用于检索有关所述缓存域的拓扑的拓扑信息;以及

拓扑数据库(51),用于存储所述拓扑信息;

其中所述拓扑信息包括与以下所述一项或更多项有关的数据:

所述缓存域的拓扑结构;

在所述拓扑结构中单元之间链路的带宽限制;

所述缓存的存储能力;

所述拓扑结构中缓存的位置;

在所述拓扑结构中单元之间所述链路承受的当前负载;以及

用于所述缓存域内数据的传输类型;

以及其中所述对象分布器配置成在所述缓存之间应如何分布所述对象的确定中利用所述拓扑信息。

3. 如权利要求2所述的缓存管理器,其中所述拓扑信息包括用于所述缓存域中每个缓存的每日负载简档,所述每日负载简档指示用于该缓存的高负载和低负载的期间,以及其中所述对象分布器配置成将所述中等受欢迎度类中的对象分配到缓存,使得在一天的所有时间,在低负载期间在至少一个缓存上存储每个对象。

4. 如权利要求2或3所述的缓存管理器,其中所述对象分布器配置成指示所述缓存分布所述中等受欢迎度类内的对象,使得更受欢迎对象被分布到较大数量的缓存,并且不那么受欢迎对象被分布到较小数量的缓存。

5. 如权利要求4所述的缓存管理器,其中所述中等受欢迎度类中的所述更受欢迎对象被分布到较靠近用户的许多缓存,并且不那么受欢迎对象被分布到所述网络中较中心的少数缓存。

6. 如前面权利要求任一项所述的缓存管理器,配置成控制所述数据对象在所述缓存域中缓存之间的所述分布,以便所述缓存域象大的虚拟缓存一样运行。

7. 如前面权利要求任一项所述的缓存管理器,其中所述网络是移动网络。

8. 如前面权利要求任一项所述的缓存管理器,其中所述缓存管理器是分立实体,并且包括用于与所述网络中单元进行通信的通信系统(59)。

9. 如权利要求8所述的缓存管理器,其中所述缓存管理器与所述缓存域中所述缓存(225,235,245,227,237)之一相关联。

10. 如权利要求1到9任一项所述的缓存管理器,其中所述缓存管理器的功能性分布在所述缓存域中的网络单元之间。

11. 一种用于优化数据对象在资源有限网络的缓存域中缓存(225,235,245,227,237)之间的分布的方法,所述方法包括:

在缓存域中的缓存接收对数据对象的用户请求(S72a);

将报告所述用户请求并且识别所述请求的数据对象的通知从收到请求的所述缓存发送到缓存管理器(26)(S72);

在所述缓存管理器,整理和存储对象信息,所述对象信息包括每个请求的数据对象的请求频率和收到所述请求的所述缓存的位置(S73,S74);

在所述缓存管理器,基于所述对象信息识别所述缓存域内要分布的对象(S75,S76,S77);

将在缓存之间分布在那些缓存中存储的数据对象的指示从所述缓存管理器发送到所述缓存(S78,S79);以及

在所述缓存之间分布所述数据对象(S80);

其中所述对象信息包括用于所述对象的受欢迎度参数,并且所述对象分布器配置成将所述对象分类到至少三个类之一中,所述类包括:

高受欢迎度类,包括应分布到所述缓存域中所有缓存的对象;

中等受欢迎度类,包括应分布到所述缓存域中所述缓存的一个或更多个子集的对象;以及

低受欢迎度类,包括不应分布的对象。

12. 如权利要求11所述的方法,其中使用所述网络在其它情况下将未使用的传送容量,分布所述对象。

13. 如权利要求11或12所述的方法,其中所述缓存管理器(26)在识别要分布的所述对象时还将所述网络的拓扑和/或所述请求的数据对象的优先级和/或对特定数据对象的请求的频率考虑在内。

14. 如权利要求11、12或13所述的方法,其中所述缓存管理器将所述缓存域中每个缓存的每日负载简档考虑在内,以及将所述中等受欢迎度类中的对象分配到缓存,使得在一天的所有时间,在低负载期间在至少一个缓存上存储所述对象,所述每日负载简档指示用于该缓存的高负载和低负载的期间。

15. 如权利要求11到14任一项所述的方法,还包括分布所述中等受欢迎度类内的对象,使得更受欢迎对象被分布到较大数量的缓存,并且不那么受欢迎的对象被分布到较小数量的缓存。

16. 如权利要求15所述的方法,其中所述中等受欢迎度类中的所述更受欢迎对象被分布到较靠近用户的许多缓存,并且不那么受欢迎对象被分布到所述网络中较中心的少数缓存。

17. 如权利要求 11 到 16 任一项所述的方法,其中所述缓存与 RNC、RBS 或其它 RAN 节点相关联或者是其一部分。

18. 如权利要求 11 到 17 任一项所述的方法,其中在非峰值期间内,在所述缓存之间分布所述数据对象 (S70)。

19. 如权利要求 11 到 18 任一项所述的方法,其中所述缓存管理器 (26) 的功能性分布在所述缓存域中的网络单元之间。

20. 如权利要求 11 到 19 任一项所述的方法,还包括:

在每个缓存,形成要分布到其它缓存的对象的分布队列;以及
在所述分布队列的头部,放置最近使用或最常使用的对象。

21. 如权利要求 11 到 20 任一项所述的方法,还包括将包括对象受欢迎度统计的消息从所述缓存发送到所述缓存控制器。

22. 一种计算机程序产品,包括适用于在资源有限网络中的缓存管理器上执行的代码,所述代码可用于促使所述缓存管理器执行以下操作:

从所述缓存管理器控制的缓存域中的缓存检索通知,所述通知包括最近在那些缓存请求的数据对象的细节;

整理和存储有关所述请求的数据对象的对象信息,所述对象信息包括请求的数据对象的位置和受欢迎度;

基于所述对象信息识别要分布的对象;以及

指示所述网络中的所述缓存在它们之间分布在那些缓存中存储的数据对象;

其中所述对象信息包括用于每个对象的受欢迎度参数,并且所述对象到至少三个类之一中,所述类包括:

高受欢迎度类,包括应分布到所述缓存域中所有缓存的对象;

中等受欢迎度类,包括应分布到所述缓存域中所述缓存的一个或更多个子集的对象;

以及

低受欢迎度类,包括不应分布的对象。

23. 如权利要求 22 所述的计算机程序产品,在载体媒体上携带。

移动网络中的缓存

技术领域

[0001] 本发明涉及用于缓存移动分组数据网络中数据的系统。具体地说,本发明涉及适合用于优化资源有限网络中缓存数据的分布的缓存体系结构。

背景技术

[0002] 典型的文件缓存方法包括接收来自文件服务器的文件并且存储整个文件的缓存。以后,在客户端需要文件时,从缓存提供文件而不是从文件服务器提供文件。由于缓存一般是更靠近客户端或者具有比文件服务器更高带宽的服务器,因此,文件迅速地从缓存被提供到客户端。

[0003] 这能够参照图 1 理解,图 1 是具有原服务器 11 和多个缓存 12-16 的网络 10 的示范体系结构的示意图。客户端 17-19 配置成接收来自原服务器 11 和 / 或缓存 12-16 的文件和 / 或流传送数据。

[0004] 为降低在原服务器 11 上的负载并且节省输送网络 10 中的带宽,一些内容存储在更靠近最终用户 17-19 的缓存 12-16 中。最好是尽可能靠近最终用户推送这些缓存。

[0005] 例如,移动网络体系结构通常包括附接了多个无线电接入网络 (RAN) 的核心网络 (CN),而 RAN 中,终端连接到无线电基站 (RBS)。终端可移动,并且在它们移动时,它们可以无线方式连接到不同 RBS。在 RBS 与网络的剩余部分之间的传送链路通常是被租赁的,并且容量有限。降低对这些链路的需求的一种方式是使用缓存,并且如上所述,缓存应优选尽可能靠近终端以最小化传送需要。

[0006] 在 RAN 中进行缓存(缓存是 RBS 或任何其它 RAN 节点的一部分)有关的一个问题是每个缓存将只接收来自一小群体的终端的业务。通常,缓存存储的内容将是最终用户已请求的内容。如果使用缓存的群体小,则在该缓存中存储的缓存内容小。从统计上来说,对于更小的群体,其他人想从该缓存下载相同内容的概率,即,“缓存命中率”降低。用于大群组的一个大缓存因此从统计上来说比服务于群体的不同更小子集的许多小缓存更有效。

[0007] 解决此问题的一种方式是通过在缓存之间复制内容,从 RBS 中的许多小缓存生成大的“虚拟”缓存。这意味着每个小缓存存储来自大群体的内容(由其它小缓存检测到和缓存),并且这增大了受欢迎内容的本地命中率。缓存的此“预加载”能够视为一种形式的预测缓存,其中,预测是基于在其它缓存中的用户。

[0008] 不可避免的是,这将导致内容跨许多小缓存的重复,并且要付出的代价是所有小缓存一起要求的总存储大小将增大。另外,在缓存之间传送的数据量随着内容的分布而增大。如果缓存之间的带宽受到限制,则另外传送的数据增大了负载,并且可最终造成拥塞。因此,在资源有限网络中能够有助于解决池化和内容分布的方法极其重要。

发明内容

[0009] 本发明的目的是解决至少一些上述缺点。通常,最好是以有效方式在一起形成大的虚拟缓存的小缓存中分布内容。

[0010] 根据本发明的一方面,提供了一种用于控制在称为分组数据网络(可选择性地称为资源有限网络)的缓存域的缓存集中存储的内容的分布的缓存管理器。缓存管理器包括用于接收来自缓存域中缓存的有关在那些缓存请求的数据对象的通知的对象监视器。对象数据库配置成整理和存储有关请求的数据对象的对象信息,对象信息包括有关请求数据对象的缓存和提出请求的频率的信息。对象分布器配置成基于对象信息识别要分布的对象并且指示网络中的缓存在它们之间分布在那些缓存中存储的数据对象。对象信息包括用于每个对象的受欢迎度参数,并且对象分布器配置成将对象分类到至少三个类之一中。类包括具有应分布到缓存域中所有缓存的对象的高受欢迎度类、具有应分布到缓存域中一个或更多个缓存的子集的对象的中等受欢迎度类及具有不应分布的对象的低受欢迎度类。将领会的是,网络可包括一个或更多个缓存域。对象分布器可配置成指示缓存使用网络在其它情况下将未使用的传送容量,分布数据对象。

[0011] 缓存管理器可还包括用于检索有关缓存域的拓扑的拓扑信息的拓扑监视器。也可提供用于存储所述拓扑信息的拓扑数据库。对象分布器可配置成在缓存之间应如何分布对象的确定中利用拓扑信息。更详细地说,对象分布器可决定哪些对象应在哪些缓存存储,以及何时及从何处应进行必需的获取(即,将内容从一个或更多个本地缓存分布到一个或更多个其它本地缓存的适合时间和位置)。

[0012] 拓扑信息可包括与以下一项或更多项有关的数据:缓存域的拓扑结构;在拓扑结构中单元之间链路的带宽限制;缓存的存储能力;拓扑结构中缓存的位置;在拓扑结构中单元之间链路承受的当前负载;以及用于缓存域内数据的传输类型。就缓存的存储能力而言,可能是带有长持久受欢迎度的对象应存储在具有有限次数的写入操作的缓存上(如闪存存储器),并且带有极高受欢迎度的对象不应存储在太少缓存操作盘上。

[0013] 拓扑信息可包括用于缓存域中每个缓存的每日负载简档,每日负载简档指示用于该缓存的高负载和低负载的期间。随后,对象分布器可配置成将中等受欢迎度类中的对象分配到缓存,使得在一天的所有时间,在低负载期间在至少一个缓存上存储每个对象。

[0014] 对象分布器可配置成指示缓存分布中等受欢迎度类内的对象,使得更受欢迎对象被分布到较大数量的缓存,并且不那么受欢迎对象被分布到较小数量的缓存。中等受欢迎度类中的更受欢迎对象可被分布到较靠近用户的许多缓存,并且不那么受欢迎对象被分布到网络中较中心的少数缓存。

[0015] 缓存管理器可配置成控制数据对象在缓存域中缓存之间的分布,以便缓存域象大的虚拟缓存一样运行。

[0016] 缓存管理器可以是分立实体,并且包括用于与网络中单元进行通信的通信系统。它可与缓存域中的缓存之一相关联。备选,缓存管理器的功能性可分布在缓存域中的网络单元之间。

[0017] 根据本发明的另一方面,提供了一种用于优化在资源有限网络的缓存域中缓存之间数据对象的分布的方法。在缓存域中的缓存接收对数据对象的用户请求。将通知从收到请求的缓存发送到缓存管理器。通知报告用户请求并且识别请求的数据对象。在缓存管理器,整理和存储对象信息,对象信息包括每个请求的数据对象的请求频率和收到请求的缓存的位置。在缓存管理器,基于对象信息识别缓存域内要分布的对象。将指示从缓存管理器发送到缓存以在缓存之间分布在那些缓存中存储的数据对象。可选的是,使用网络在其

它情况下将未使用的传送容量,分布数据对象。对象信息包括用于对象的受欢迎度参数,并且对象分布器配置成将对象分类到至少三个类之一中。类包括具有应分布到缓存域中所有缓存的对象的高受欢迎度类、具有应分布到缓存域中缓存的子集的对象的中等受欢迎度类及具有不应分布的对象的低受欢迎度类。

[0018] 缓存管理器的功能性可分布在缓存域中的网络单元之间。缓存管理器可在识别要分布的对象时将网络的拓扑和/或请求的数据对象的优先级和/或对特定数据对象的请求的频率考虑在内。网络可以是移动网络,并且缓存可与 RNC 或 RBS 相关联或者是其一部分。

[0019] 分布可通过在每个缓存形成要分布到其它缓存的对象的分布队列以及在分布队列的头部放置最近使用或最常使用的对象来实现。

[0020] 包括对象受欢迎度统计的消息可从缓存发送到缓存控制器。

[0021] 根据本发明的又一方面,提供了一种包括适用于在资源有限网络中的缓存管理器上执行的代码的计算机程序产品。代码可用于促使缓存管理器从缓存管理器控制的缓存域中的缓存检索通知,通知包括最近在那些缓存请求的数据对象的细节。代码还可用于促使缓存管理器整理和存储有关请求的数据对象的对象信息,对象信息包括请求的数据对象的位置和受欢迎度。代码还可用于促使缓存管理器基于对象信息识别要分布的对象。代码还可用于促使网络单元指示网络中的缓存在它们之间分布在那些缓存中存储的数据对象。类包括具有应分布到缓存域中所有缓存的对象的高受欢迎度类、具有应分布到缓存域中一个或更多个缓存的子集的对象的中等受欢迎度类及具有不应分布的对象的低受欢迎度类。

[0022] 本发明也提供在载体媒体上携带的如上所述的计算机程序产品。

附图说明

[0023] 现在将仅通过示例,并参照附图描述一些优选实施例,其中:

图 1 是具有原服务器和多个缓存的网络的示意图;

图 2 是带有配置成充当缓存服务器的无线电基站和无线电网络控制器的无线电接入网络的示意图;

图 3 是示出用于服务于用户请求的过程的流程图;

图 4 是示出缓存群组的不同峰值负载期间的示意图;

图 5 是示出用于缓存管理器的信息模型的示意图;

图 6 是配置成充当缓存服务器的 RNC 的示意图;以及

图 7 是示出图 2 的网络中缓存的数据的管理的示意图。

具体实施方式

[0024] 如在上面背景中所述,设计缓存系统时,一般存在在中心放置缓存与本地放置缓存之间的选择。在两种情况下,传送方面的增益均很小,要么因为放置点在网络中位置高(就中心放置缓存而言),要么因为命中率变得太低(就本地缓存而言)。

[0025] 因此,最好是组合中心缓存的命中率和本地缓存的传送优点。这能够通过使用确保本地缓存具有接近或等于中心缓存的高命中率的“缓存管理器”而得以实现。它通过根据全局统计预加载本地缓存(因此命中率高)而得以完成并且确保此预加载在非峰值时间进行(以便不妨碍带宽节省)。

[0026] 因此,实际上从许多小的缓存生成大的虚拟缓存,其中,将传送的限制考虑在内。这能够通过缓存之间分布内容的“缓存均衡化”来实现。每个小缓存存储大群体请求的内容(由其它小缓存检测到和缓存),并且这增大了受欢迎内容的本地命中率。在小缓存之间的对象数据分布能够视为预测缓存的一种形式,其中,预测是基于通过其它缓存接收数据的用户的行为。

[0027] 要付出的代价是在缓存之间传送的数据量将增大。如果带宽在缓存之间受到限制,则另外传送的数据将增大负载,并且最终造成拥塞。因此,在资源有限网络中能够有助于解决内容分布的方法极其重要。

[0028] 业务在典型的 24 小时期间内有所不同,并且通常在晚上时间有大量的空闲容量。此外,甚至在峰值期间,也存在空闲容量可用的期间。因此,可能在峰值时间期间(例如,在白天期间)分布少量的高优先级数据对象,并且在非峰值时间内(例如,在晚上)分布更大量的不那么重要的数据对象。在此上下文中,数据对象将被理解为表示终端请求和 / 或缓存存储的任何数据集。

[0029] 图 2 是包括连接到三个 RNC 22、23、24 的路由器 21 的网络 20 的示意图,每个 RNC 连接到 RBS 221、222 ;231、232 ;241、242。每个 RBS 221、222、231、232、241、242 连接到一个或多个终端 251-260。每个 RNC 22、23、24 和一些 RBS221、231 与可构建或不可构建到 RNC 或 RBS 本身中的缓存存储单元 224、234、244、226、236 相关联。RNC 22、23、24 和路由器 21 通过传送链路互连。RNC (和与缓存存储单元相关联的 RBS 221、231 的那些 RBS)能够充当缓存服务器 - 即,每个包括表现得象嵌入式代理或使用深度分组检查 (DPI) 检测和引导诸如 HTTP-GET 等最终用户请求的逻辑实体。为方便起见,缓存服务器(例如,RNC 22)和缓存存储单元(例如,224)中的缓存功能性的组合将在本文档中称为“缓存”。无论缓存存储单元是否物理构建到缓存服务器中,这都适用。这意味着实际上有与每个 RNC 22、23、24 和两个 RBS 221、231 相关联的缓存 225、235、245、227、237。

[0030] 虽然图 2 示出带有由 RNC 22、23、24 和 RBS 221、231 操作的缓存 225、235、245、227、237 的树状拓扑,但将领会的是,这只是示例。如将变得明白的一样,用于本文中所述分布式缓存数据的机制可适用于许多其它拓扑和在分布式缓存 225、235、245、227、237 之间存在有限带宽的情形。

[0031] 如上所示,与此类缓存有关的问题涉及以下实际情况:由于配置为缓存的每个 RNC 22、23、24 和 RBS 221、231 只供应数据到有限数量的终端,因此,来自每个 RNC 22、23、24 和 RBS 221、231 的业务量小。来自本地缓存业务的“命中率”小。解决方案是通过将由群组中任何小缓存 225、235、245、227、237 在本地缓存的内容分布到该群组中的其它小缓存,增大命中率。随后,群组中的缓存将包含在其它缓存请求的内容,并且这将增大命中率。

[0032] 因此,如果内容由终端之一 251 请求,则此请求通过终端 251 连接到的 RBS 221 和 RNC 22 传递。内容数据从网络(例如,在核心网络中,未示出)中上游的内容服务器 25 (与内容存储单元 254 相关联)通过 RNC 22 和 RBS 221 输送到终端 251。除输送内容数据到终端 251 外,RNC 22 和 RBS 221 的每个也在其相关联缓存存储单元 224、226 中保存内容数据。随后,下次附接到该 RBS 的终端 251 或 252 请求该内容时,RBS 221 便能够从缓存存储单元 226 提取它。附接到相同 RNC 22 的另一终端 253 请求该内容时,能够从缓存存储单元 224 (或由 RBS 221 从缓存存储单元 226)提取它并将它供应到请求终端。另外,内容数据被复

制到群组中的其它缓存 235、245、237 (由 RNC 23、24 和其它 RBS 231 保持),以便如果被请求,则它也能够被供应到任何其它终端 254-260。

[0033] 与分布有关的问题是链路带宽有限,使得更新能够造成拥塞,特别是如果它们是在忙时间期间进行。因此,更新应受到控制以便优化带宽的使用。此优化可包括以下所述的一些或全部:

●使用在其它情况下未利用的传送容量(例如,在非峰值时间期间)。

[0034] ●优先处理可能受欢迎的数据。例如,能够在带宽变得可用时以及基于数据对象受欢迎度,根据排队原则,排队和提供要分布的数据对象。

[0035] ●如果传送带宽或存储容量是限制因素,则限制极少使用的数据对象的分布。可不分布一些数据。如果要求,则能够从已保存它的缓存提取它。

[0036] 此外,情况可以是经 RBS 221 对与该 RBS 相关联的缓存 227 中或与 RBS 221 的 RNC 22 相关联的缓存 225 中未存储但在附近缓存 235、245、237 中存储的数据对象提出请求。在此情况下,可将数据对象标记为高优先级并且立即从附近 RNC 235、245、237 获取它,而不是等待带宽变得可用。

[0037] 为控制数据在缓存 225、235、245、227、237 之间的分布,集中式缓存管理器单元 26 操作性耦合到网络 20 中的节点之一。在图 2 中,它示为与 RNC 之一 24 相关联,但将领会的是,它也能够与任何 RNC 22、23、24 或例如路由器 21 等任何另一网络单元或 RBS 之一相关联。缓存管理器单元 26 控制在分布式缓存集之间信息的扩散。将领会的是,也可能设计提供相同功能性的分布式缓存管理器解决方案。重要的特征是缓存管理器(无论它是如图 2 所示单个单元 26 还是分布式系统)能够优化数据在缓存 225、235、245、227、237 之间的分布。

[0038] 特定管理器控制的缓存集可被视为“缓存域”(CD)。CD 可相互排斥或重叠。管理器可放置在带有其域的拓扑中中心位置的站点或者与其相关联,并且它可以是该站点的 RBS (或任何另一缓存实体)的一部分。CD 能够对应于 RAN,并且可以只是 RAN 的一部分,或者包含多于一个 RAN。

[0039] 缓存管理器 26 知道在任何 RNC 22、23、24 (和与缓存相关联的 RBS 221、231) 请求的所有数据对象的受欢迎度、在缓存 225、235、245、227、237 中的所有可能出现(occurrence) 和在网络中的带宽限制,并且它配置成根据以上所述使用此信息管理缓存更新。另外,它配置成检测和重新引导来自用户的“失败的”数据请求,使得从相同 CD 中优先于内容服务器 25 的其它缓存获取与用户直接相邻的缓存不能输送的请求的数据对象。这样,缓存管理器 26 能够在图 2 中路由器 21 与更高层之间的跳提供(相当大的)传送增益,并且因此同样更有效地使用链路。将注意的是,在连接由租赁传送或微波链路提供的情况下,此方面特别有吸引力。

[0040] 换言之,缓存管理器 26 通过在有空闲带宽的时间期间在小的本地缓存 225、235、245 之间分布内容,增大小的本地缓存 225、235、245 的“有效”命中率。通过分析来自不同缓存 225、235、245、227、237 的请求,确定要推送的内容。在峰值时间期间已在一个或多个 RBS 上请求(且因此在对应缓存存储单元 224、234、244、226、236 缓存)的内容将在非峰值时间期间被分布到其它缓存(且因此在这些 RBS 上请求时可供输送)。这样,一些业务从峰值时间(提出请求时的时间)移到非峰值时间(预填充缓存时的时间)。不同缓存存储单元

224、234、244、226、236 的内容因此将在峰值时间期间有差异(这是因为在不同 RNC 22、23、24 的用户请求不同数据对象),并且在某种程度上在非峰值时间期间汇聚(这是因为通过复制任何缓存中的内容到所有缓存,缓存可变得均衡)。如下面将更详细解释的一样,不是所有内容必需对所有缓存是均衡的。

[0041] 缓存管理器知道在何处找到数据对象,并且决定在每个缓存上存储的内容。因此,可能可决定根本不存储数据对象,或者存储数据对象,但只在一些缓存存储而在其它缓存不存储。在哪些缓存中存储的选择可取决于带宽、盘空间、普遍受欢迎度及在某些用户群组中的受欢迎度。下面提供有关如何做出此选择的更多细节。

[0042] 用于服务于用户请求的典型过程在图 3 中示出并且如下所述继续:

S1 移动终端 251 的用户发出对某段内容的请求(例如,HTTP-GET)。

[0043] S2 在 RBS 221 (或其它网络单元)的站点由(可构建或不可构建到 RBS 221 中的)缓存 227 的代理 /DPI 逻辑单元 53 (参见图 5) 截接请求。

[0044] S3 逻辑单元 53 询问与 RBS 22 相关联的缓存存储单元 226 以了解内容是否存储在该缓存存储单元 226 中。

[0045] S4 如果是,这称为“主要缓存命中”:从缓存存储单元 226 检索内容,并且进行到终端 251 的本地输送。在此情况下,在图 2 中的所有链路上实现了带宽节省。

[0046] S5 如果存在第一缓存缺失,即,如果请求的内容在缓存存储单元 226 中不可用,则请求将被转发到“真实的”内容服务器 25,到更高层代理或到缓存管理器 26。

[0047] S6 请求在缓存管理器 26 被截接(在其到内容服务器 26 或更高层代理的路径上),或者到达缓存管理器 26 (作为选择的目的地)。

[0048] S7 缓存管理器 26 分析请求并且检查在 CD 中的任何其它本地缓存存储单元中是否找到请求的内容。

[0049] S8 如果是(即,存在次要缓存命中),则将请求转发到包含内容的任何缓存 235。

[0050] S9 从缓存存储单元 234 检索内容,并且将其输送到用户。

[0051] S10 第一缓存 (RBS 221) 可截接输送,并且在其相关联缓存存储单元 226 中本地存储内容以用于将来的请求。在此情况下,在图 2 中路由器 21 上方的所有链路上实现了带宽节省。将领会的是,这能够在路径中的任何缓存发生,如带有相关联存储单元 224 的 RNC 22。

[0052] S11 如果存在次要缓存命中(即,如果请求的内容在 CD 中的任何其它缓存存储单元中不可用),或者如果根据朝向其它本地缓存的链路负载条件指示,则将对内容的请求转发到内容服务器 25 或某一更高层代理。备选是在以管理器作为最后解决办法的每个更高层缓存截接请求。更高层缓存可具有内容,或者知道在它们下面将找到内容的缓存(以管理器知道它们的相同方式)。至少在第一情况下,继续到管理器没有意义。

[0053] S12 将内容从内容服务器 25 或更高层代理输送到用户。

[0054] S13 第一缓存 (RBS 221) 或路径中的任何另一缓存 (RNC 22) 可截接输送,并且在其相关联缓存存储单元 226 中本地存储内容以用于将来的请求。

[0055] 这样,将为本地缓存 225、235、245、227、237 填充已由本地用户请求的内容。从不同来源输送内容:(i) 从本地缓存本身(如果该 RNC 或 RBS 或站点的其它用户已请求它,或者在非峰值时间期间已将它推送到该站点,非峰值时间能够是空闲带宽可用的任何时间

期), (ii) 从 CD 中的另一本地缓存(如果 CD 中的另一用户已请求它), 或者 (iii) 从更远的来源(如果 CD 中的任何用户尚未请求它)。从本地缓存获取内容时获得最大传送增益, 从另一本地缓存获取内容时获得相当大的增益, 以及从远处来源获取内容时获得小增益或未获得增益(虽然在此阶段缓存数据可在将来产生增益)。

[0056] 因此, 明显的是, 如果本地缓存存储单元 224、234、244、226、236 的至少一些包含整个 CD 中其它用户已请求的内容, 并且因此, 本地用户将请求它的概率高, 则此过程最有效。如果情况是如此, 则主要缓存命中 (S4) 将更频繁发生, 从而降低执行步骤 S7-S13 的需要及传送链路和带宽的其相关联使用。换言之, 这增大了能够获得最高传送增益 (i) 的请求的部分。

[0057] 为实现此目的, 缓存管理器 26 收集有关此类请求的信息(主要在峰值时间期间发生), 并且随后使用此信息预填充缓存(主要在非峰值时间期间)。

[0058] 例如, 在 CD 中(例如, 由终端 251) 初次请求数据对象 X 时, 从内容服务器 25 或更高层代理获取 X, 并且在路由请求通过的缓存 225 和 227 中存储它。如果不久以后相同终端 251 或另一终端 252 经相同 RBS 221 请求 X, 则从缓存 227 获得它, 或另一终端 253 经相同 RNC 22 请求 X, 则从缓存 225 获取它并将它输送到用户。如果随后经相同 CD 中的不同 RNC 23 请求 X, 则从与第一 RNC 22 相关联的缓存存储单元 224 获取它, 并且也将它存储在第二 RNC 23 的缓存存储单元 234 中。随后, 在某个非峰值时间, 缓存管理器 26 确保 X 被分布到一些或所有缓存 235、245、237, 以便由发生请求的缓存服务于对 X 的随后请求。

[0059] 将领会的是, RNC 之间的传输链路可具有不同的容量。中心缓存控制器 26 因此能够决定分布是否应立即启动, 或者它是否应延迟到更低业务期间或晚上时间。另一选择是使用“低于尽力而为型”优先级类或比常见 TCP 更快回退的 TCP, 以便赋予分布业务低后台优先级并且不影响高优先级业务。

[0060] 为实现此操作, 缓存控制器 26 维护用于其域内被请求的每个对象的计数器。每次收到对该对象的新请求时, 增大计数器, 直至达到阈值, 在该点启动对所有缓存的均衡化。

[0061] 实际上, 内容可在受欢迎度方面有所有不同。受欢迎的内容应被广泛分布到 CD 中的所有(或几乎所有)缓存, 但不那么受欢迎的内容应只被分布到少数缓存。为监视受欢迎度, 中心控制器 26 跟踪来自移动网络或能够被视为形成“池”的其域中缓存覆盖的区域中 UE 251-260 的对内容的所有请求。它也跟踪访问频率以确定对象的受欢迎度。

[0062] 对象的受欢迎度用于确定应如何缓存对象。能够认为分布有三个基本类别, 但将领会的是, 其它子类别也可有用。三个基本分布类别是:

- (i) 分布到域中所有缓存的内容(“均衡化”),
- (ii) 根本不缓存的内容, 或者
- (iii) 分布到缓存的一个或更多个子集(“池化”)的内容。

[0063] 在最后的情况下, 对象的受欢迎度也能够用于确定应存储该对象的缓存的数量和位置。另外, 确定应从缓存删除哪些对象时, 应使用类似的统计。

[0064] 如上所示, 内容受欢迎度最初可分类为“高”、“中等”或“低”。另外:

- 均衡化适用于优选在所有缓存中存储其副本的高受欢迎度的内容,
- 池化适用于优选在一些缓存中存储其副本的中等受欢迎度的内容, 以及
- 根本无动作适用于不存储其副本的低受欢迎度的内容。

[0065] 简单的池化策略是只在输送对象到请求该对象的最终用户经过的网络单元(RNC, RBS)的一些或所有缓存中存储请求的对象。然而,此简单的方案可能导致通过高负载链路的输送、在缓存之间请求的不均匀分布和/或在长距离内的输送。

[0066] 池化对象的决定涉及选择应保持所述内容的缓存。如上所述,此类选择应优选计及至少负载(优选负载更低的链路)、拓扑(优选更短的距离)和容量。

[0067] 链路负载的问题包括将对象发送到缓存(由于池化它的决定的原因),但具体而言,从缓存输送对象(由于用户请求对象的原因)。第一方面类似于均衡化,并且能够在一定程度上受到控制,这是因为缓存控制器 26 能够确保在非峰值期间内将对象发送到缓存。第二方面不同,表现在将对象从缓存输送到用户可能多次执行和在由用户而不是由系统确定的时间执行。

[0068] 如前面所述,业务随时间有规律地变化。然而,前面未考虑的是变化模式在不同节点之间不同的实际情况。例如,对于更小区域和/或每节点的更少用户,变化可更大。通过根据缓存峰值时间将缓存编组,并且在每个群组中一个或更多个缓存上存储对象,这些差别能够用于避免用于被池化对象的峰值负载。随后,从当前未遇到峰值负载的群组中的任何缓存获取对象。

[0069] 这能够通过参照图 4 理解,图 4 示出用于三个缓存群组 41、42、43 的对照时间的平均负载模式,每个群组被分类到上午(M)、中午(N)和下午(A)之一。“上午”群组 41 的峰值负载 44 发生在 5.00 与 11.00 之间,“中午”群组 42 的峰值负载 45 发生在 10.00 与 16.00 之间,以及“下午”群组 43 的峰值负载 46 发生在 15.00 与 21.00 之间。例如,这能够通过居住在(其中大多数缓存可能形成“下午”群组的一部分)和商业区(其中大多数缓存可能形成中午群组的一部分)中网络单元之间的差别而设想到。如果对象存储在每个此类群组的至少一个成员上,则始终将有包含该对象的不可能经受峰值负载的至少一个缓存。将理解的是,此示例可易于扩展到多于三个群组。

[0070] 识别此类群组以允许从“便宜”资源池化,这是缓存控制器 26 中池化布置机制的任务。

[0071] 另外,如上所述,从缓存获取对象的“成本”(例如,延迟和带宽)可不但取决于时间,而且也取决于拓扑。因此,优选的是可完全根据被池化对象的受欢迎度以不同方式处理它们。例如,可有利地将其受欢迎度几乎足以分类为“均衡化”的对象发送到许多缓存(每个群组中的许多节点),而可有利地将其受欢迎度几乎足以分类为“无动作”的对象保持在一个位置(每个群组中的一个节点)。将注意的是,在此布置中,对受欢迎对象的请求将被分布在比对不那么受欢迎对象的请求更多的节点,并且这样在不同节点上的负载变得更均匀分布。

[0072] 此外,存储对象的位置越少,此位置位于拓扑中心将越是优选,且反之亦然。这是因为包含所述对象的节点越多,能够从拓扑“方便的”节点服务于请求的机会就越高。

[0073] 因此,识别被池化对象的适合数量的实例和应存储对象的位置也是池化布置机制的任务。一个实现可基于群组的数量。例如,关于两个群组,在每个峰值时间群组中在多个“远程”位置存储第一受欢迎度群组中的最(受欢迎)和最不(受欢迎)对象,而在每个峰值负载群组中在少数“中心”位置存储第二受欢迎度群组中的对象。明显此示例自然扩展到多于两个群组。

[0074] 换而言之,在树状拓扑中,应在每个群组中的许多“叶”缓存存储更受欢迎的被池化对象,并且因此确保至少一个对象将靠近每个访问点。不那么受欢迎的对象存储在更少量的缓存,但那些选择的缓存是“根”缓存,以便同样地,至少一个对象较靠近每个访问点。

[0075] 池化布置机制应按其传送链路的峰值时间将缓存编组,将被池化对象的受欢迎度分类,并且匹配每个这样的类和每个峰值时间群组中节点的子集,使得不那么受欢迎的对象往往存储在更少但更中心的位置,而更受欢迎的对象往往存储在更多但更不中心的位置。

[0076] 除缓存的每日负载简档和拓扑外,在确定分布时将缓存本身的存储技术考虑在内也是有用的。例如,闪存存储器受以写入操作的次数衡量的其寿命限制:闪存存储器只能够被重写有限的次数。因此,在此类缓存上存储有“长期”受欢迎度的对象将是有益的,但优选避免(在可能的情况下)存储受欢迎度快速减弱的将需要经常替换的对象。相比之下,盘存储装置具有更长的寿命,但受带宽限制,因此,最好是避免在相同盘上存储太多受欢迎的对象,这是因为不可能以要求的速率将这些对象输送到许多用户。

[0077] 因此,如上所述,应根据适合的准则选择被池化对象分布到的缓存,这些准则可包括:

- 在一段时间内的业务负载模式;
- 在传送基础设施内的地理分布;
- 缓存存储技术的限制。

[0078] 为实现分布,可采用在缓存删除算法中经常使用的方案的反向方案。在满缓存中要删除哪些对象的问题的普通解决方案是删除最近最少使用(least-recently-used,LRU)或最不常使用(LFU)对象。在控制对象的分布时,此方案能够反向进行,从而先分布最近使用(most-recently-used,MRU)或最常使用(MFU)。例如,我们可在缓存使用“分布队列”(例如,每目的地一个队列),其中(a)存储了在此缓存中发现并且应分布到其它缓存的对象,以及(b)使用空闲带宽从中传送当前MRU或MFU对象(即,作为相对于是前台的用户请求的后台业务)。

[0079] 图5是示出用于缓存管理器26的适合信息模型的示意图。管理器包括以下信息库:

- 可存储与以下所述有关的信息的拓扑数据库51:
 - CD的拓扑结构:(例如,树型、星形等)
 - 拓扑中链路的带宽限制。

[0080] ● 拓扑结构中缓存的位置。

[0081] ● 当前链路负载。拓扑可支持多个业务类型(例如,交互式业务和谈话业务)。链路负载监视在网络单元之间链路上的负载,以避免数据对象的分布造成拥塞。

[0082] ● 传输类型:多播或单播。

[0083] - 可存储与以下所述有关的信息的对象数据库52:

- 数据对象的位置和能力,即,其中存储了不同数据对象的缓存(如果有)。

[0084] ● 包括诸如对数据对象的先前请求等历史信息的数据对象的受欢迎度。

[0085] ● 缓存的数据对象,即,最近缓存的数据对象列表。

[0086] - 可包括以下所述的对象分布状态数据库53:

●对象分布状态:这应是提供有关对象应被分布的广泛程度的信息的参数(例如,“应均衡化”、“应被池化到许多缓存”、“应被池化到少数缓存”、“不应分布”等)。

[0087] ●优先级:数据对象分布优先级。

[0088] ●挂钟时间:用于确定应开始分布的时间点的计时信息。

[0089] 管理器 26 中的功能性可分成以下块:

- 拓扑监视器 54:检索或收集有关分布式缓存 225、235、245、227、237 的拓扑的信息的功能。

[0090] - 对象监视器 55:检索或收集来自 CD 内本地缓存 225、235、245、227、237 的有关最近缓存的数据对象的信息的功能。信息例如可通过缓存的定期轮询,通过预订来自缓存的信息,或者通过监视和分析在缓存与服务器之间的业务而获得。本领域技术人员将明白其它可能性。

[0091] - 对象分布器 56:识别应重新分布的数据对象并且执行重新分布的功能。这些动作可以是连续的过程,或者它们可以某个频率,在某些时间和 / 或在某些负载条件发生。

[0092] - 对象优先化器 57:使用例如受欢迎度统计估计数据对象的多个请求的概率以便区分已识别要重新分布的不同数据对象的优先级的功能。

[0093] - 资源管理器 58:确保重新分布不使网络过载的功能。这例如能够通过限制在单个特定非峰值时间内要重新分布的数据对象的数量来实现,但随后在下一非峰值时间或在适用时通过选择应从哪个缓存输送数据对象来继续分布。

[0094] - 通信系统 59:控制与 CD 中其它网络单元的通信的功能。

[0095] 将领会的是,图 5 所示逻辑单元全部相互交互以允许缓存管理器 26 执行所述功能。具体而言,上述池化布置机制实际上包括拓扑监视器、对象分布器、对象优先化器、对象监视器和资源管理器的功能性。

[0096] 另外(并且图 5 中未出),可存在请求匹配完成器(request match maker),其功能是将请求从发生缺失的缓存重新引导到缓存有数据对象且因此将发生命中的缓存。这样,甚至在存在缓存均衡化延迟的情况下,能够节省带宽,产生相当大的带宽节省。

[0097] 图 6 是示出充当缓存 225 的 RNC 22 和相关联缓存存储单元 224 的功能性的示意图。将领会的是,缓存能够在任何适合的网络单元及 RNC 或不使用 RNC 提供,或者与其相关联,并且描述为由 RNC 22 操作的功能可由单独的实体操作。

[0098] RNC 22 包括用于操作 RNC 的常见功能(例如,与 RBS、路由器进行通信,转发业务等)的控制器 61 和通信系统 62。在 RNC 配置成充当缓存 225 的情况下,如下所述,它也与缓存存储单元 224 相关联,并且包括两个其它功能:DPI 54 和对象代理器 54。

[0099] 缓存 225 中的功能性可描述为:

-DPI 63:检查通过的分组以查找包括诸如 HTTP-GET 等对信息的请求的那些分组的功能。它能够实现为 HTTP 代理或深度分组检查装置。

[0100] - 缓存存储单元 224:诸如硬盘等存储空间。这可与 RNC 22 (如图 6 所示)分开或 RNC 22 的组成部分。

[0101] - 对象代理器 64:朝向缓存管理器 26 的接口(如果缓存管理器 26 不是 RNC 22 的一部分,则经通信系统 62)。对象代理器 64 也能够通过低优先级的 TCP/UDP 会话分布信息,使得其它业务(例如,交互式谈话)在传送期间不受干扰。这能够以两种方式进行:使用

“低于尽力而为型”优先级类或比常见 TCP 更快回退的 TCP 会话。

[0102] 图 7 是示出在图 2 所示网络中能够如何管理缓存数据的分布的一个示例的示意图。由于空间原因,省略了图 2 的 RBS 和一些终端。信息流程如下所述:

S71 拓扑监视器 54 例如使用简单网络管理协议 (SNMP) 监视 CD 中网络单元的拓扑,并且在拓扑数据库 51 中存储更新的知识。

[0103] S72 对象监视器 55 经那些缓存中的对象代理器 64 或者通过截接内容请求,不断从本地缓存 225、235、245 (和图 7 中未示出的缓存 227 和 237) 获取有关缓存 225、235、245、227、237 中最近存储的数据对象的信息(例如,由于来自终端 251 的请求 S72a 的原因)。

[0104] S73 此信息被存储到对象数据库 52 中。

[0105] S74 在特定时间或特定负载条件,对象分布器 56 检查对象数据库 52 是否有新数据对象,并且编译新数据对象列表。

[0106] S75 新数据对象列表被传递到对象优先化器 57。“对象优先化器”指派优先级到数据对象,并且创建数据对象的优先化列表。

[0107] S76 随后将优先化列表传送到资源管理器 58。

[0108] S77 资源管理器 58 从拓扑数据库 51 获取有关拓扑状态的信息,并且基于拓扑信息设置要传送的数据对象的数量。

[0109] S78 随后,将新列表传送回对象分布器 56。

[0110] S79 基于处理的列表,对象分布器 56 创建发送到一个或多个缓存 225、235、245 中对象代理器 64 的对象分布请求集。能够发送可在更长时间内分布的几个请求。请求能够包括有关来源(发现数据对象的本地缓存)、宿(应存储数据对象的本地缓存)和有关分布方法(单播、多播)的信息。

[0111] S80 对象代理器 64 在接收此类请求后启动传送。此类传送可以许多方式实现。作为第一直接示例,在管理器的对象分布器 56 可指示在目的地缓存的对象代理器发出普通请求消息(即,HTTP-GET 消息),并确保这些请求发送到来源本地缓存的对象代理器。作为第二更高级示例,在管理器的对象分布器 46 可指示在来源缓存的对象代理器设置广播,并且指示在目的地缓存的指定集的对象代理器收听这些广播。

[0112] 上述顺序是在分布式缓存体系结构中更新的示例。作为补充示例,对象监视器 55、对象分布器 56、对象优先化器 57 和资源管理器 59 提供的一些功能性可实现为一个或多个对象传送队列 (OTQ)。

[0113] OTQ 包含在一个或多个缓存要更新的数据对象(或数据对象的引用)。根据基于优先级的原则服务于请求,其中,最先服务最紧急的更新。

[0114] OTQ 中的数据对象可以是包括到来源(例如,缓存 225)的指针和到接收器(例如,缓存 235、245)的一个或多个指针的更新。我们能够设想单个全局 OTQ:用于从来源获取数据对象的一个 OTQ 和用于分发数据对象到接收器的另一 OTQ;或可选的是,用于每个本地缓存的“获取 OTQ”和“分发 OTQ”。将注意到的是,分开的队列是逻辑记法而不是物理实现。

[0115] 对象监视器 55 和对象优先化器 57 可负责在观察到新或更新的数据对象时添加元素到 OTQ。在观察到更近的更新或更受欢迎的数据对象时,可用其它条目替换排队的条目或者将其它条目置于排队的条目之前。

[0116] 资源管理器 58 可监查来自每个缓存的业务,并且在业务低于某个阈值时激活对

应缓存的获取 OTQ 中的第一条目。类似地,它可监查到每个缓存的业务,并且在业务低于某个阈值时激活对应缓存的分发 OTQ 中的第一条目。同时,将由管理器本身或者由带有空闲空间的任何中间缓存暂时存储已获取但尚未分发的数据对象。此“全局待机调度”应使用所有可用带宽以便最小化传送时间,并且因此最大化在本地缓存的命中率。

[0117] 如上所述,术语“待机”与在检测到空闲容量(可例如根据链路负载来识别)时服务于请求的实际情况有关。然而,该含意可扩展,并且不但指请求而且指单一分组。这能够通过两个端节点的链路调度中 MPLS 中的业务区分来完成,或者借助于比其它业务更快和更强产生拥塞的“软”端对端协议(即,比 TCP 更具响应性的协议)来完成。

[0118] 将领会的是,在某种意义上,所有本地缓存 225、235、245 具有代理器,即,收听并且服务“主控”(即,缓存管理器 26)的“从属代理器”。缓存管理器的分布式实现可扩展这些本地从属代理器,或者将新的本地管理器实现为“主控代理器”。

[0119] 此外,将领会的是,在 RNC 22、23、24 和 RBS 221、231 与缓存管理器 26 之间稍微不同模式的通信是可能的。例如,考虑在 RNC 22 的本地缓存中的两种情况:主要缓存命中和次要缓存缺失(如在图 3 的步骤 S4 和 S5 中一样)。

[0120] 如果有命中(即,在缓存存储单元 224 中存在请求的内容),则由于大小或安全性原因,RNC 22 可将请求通知或不通知缓存管理器 26。一些数据对象可包含病毒或者门户不信任。这能够由缓存中的本地策略规则描述。通知管理器可帮助管理器改进其在其它缓存放置内容的计划(更受欢迎的内容可能应越快共享到其它缓存),并且具体而言,改进用于识别在缓存已满时能够删除的数据对象的其建议(缓存已满时,应先删除不那么受欢迎的内容,并且缓存管理器的“全局”视野可有助于以统计上准确的方式识别这些数据对象)。

[0121] 如果有缺失(即,缓存存储单元 224 中不存在请求的内容),则请求的本地缓存能够从别处(一般为更高层缓存或内容服务器 25)获取内容,并且随后通知缓存管理器 26 其新的获得。备选,在请求之前,它能够询问缓存管理器 26 应从何处获取内容,并且缓存管理器 26 能够引导 RNC 22 到另一本地缓存。在又一备选,所有请求能够通过缓存管理器 26 处理,以便请求的本地缓存不知道内容的来源。

[0122] 上述布置允许分布式缓存体系结构执行与大缓存类似的功能,甚至在资源有限的情况下也允许。与缓存位于拓扑中的位置无关,能够提供高缓存命中率。

[0123] 缓存的更新能够得以执行而不干扰其它业务,例如在非忙时间执行,以确保分布不要求传输网络中的更多容量。

[0124] 应明白的是,上述方案允许分布式缓存体系结构以与大的中心缓存类似的方式执行。与缓存位于拓扑中的位置无关,能够提供高缓存命中率。缓存的更新得以执行而不干扰其它业务,例如在非忙时间执行,以确保分布不要求传输网络中的更多容量。该机制对于低频率内容请求是最佳的,这是因为内容能够从所有连接的 RNC 的缓存池获取。

[0125] 此外,通过将对象的受欢迎度细分,能够优化分布。特定内容的频率受欢迎,足以分布到 RNC 中的所有缓存时,机制启动将识别的受欢迎对象分布到 RNC 中所有缓存的均衡化机制。

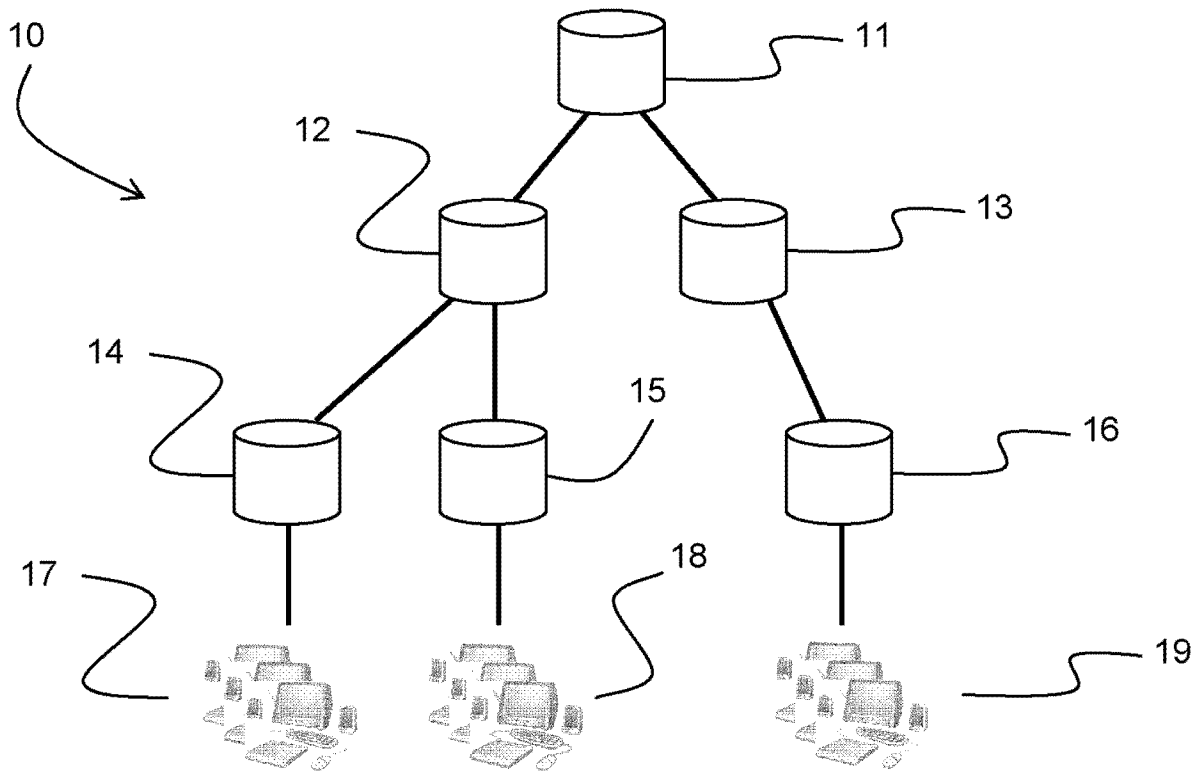


图 1

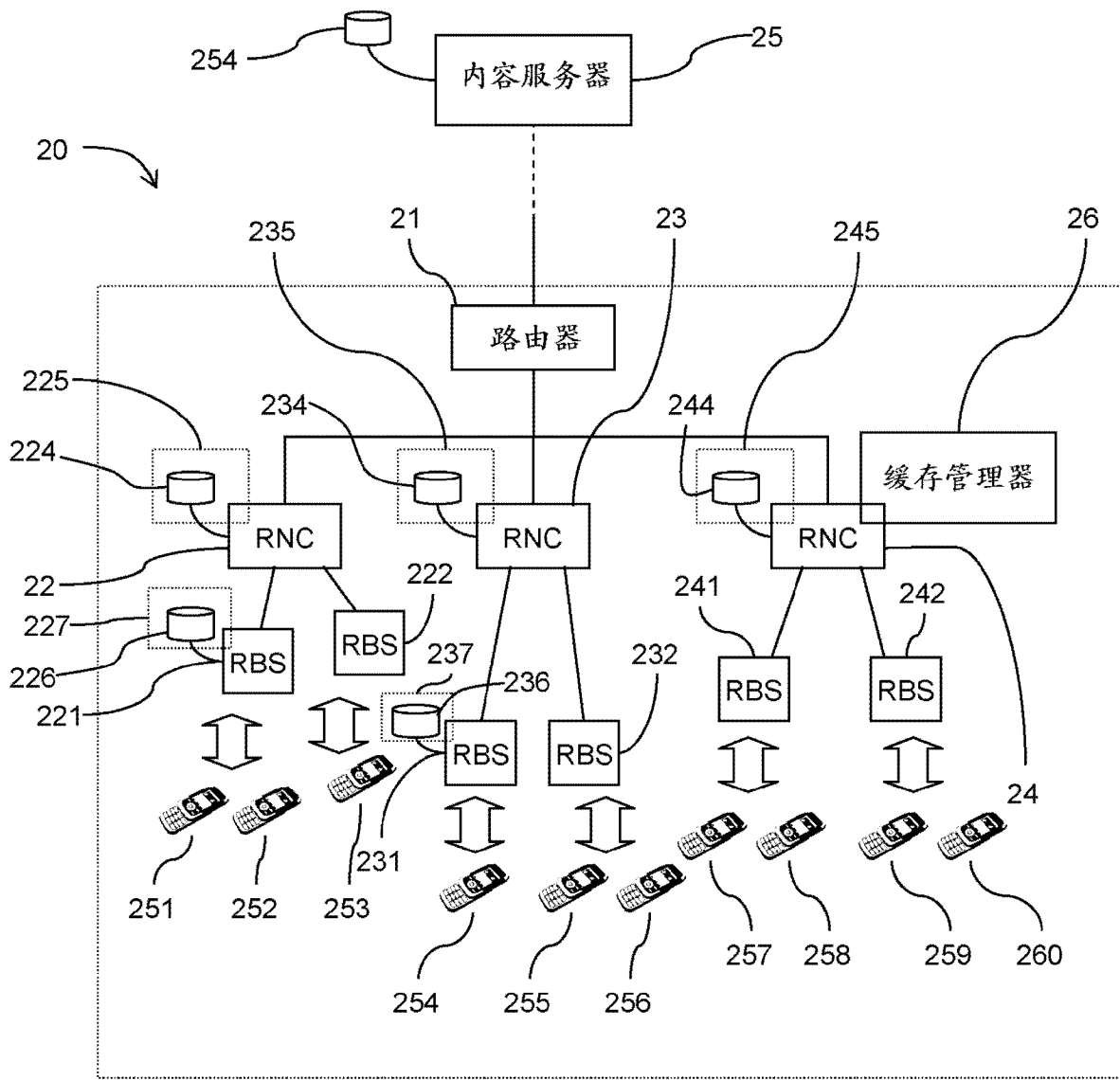


图 2

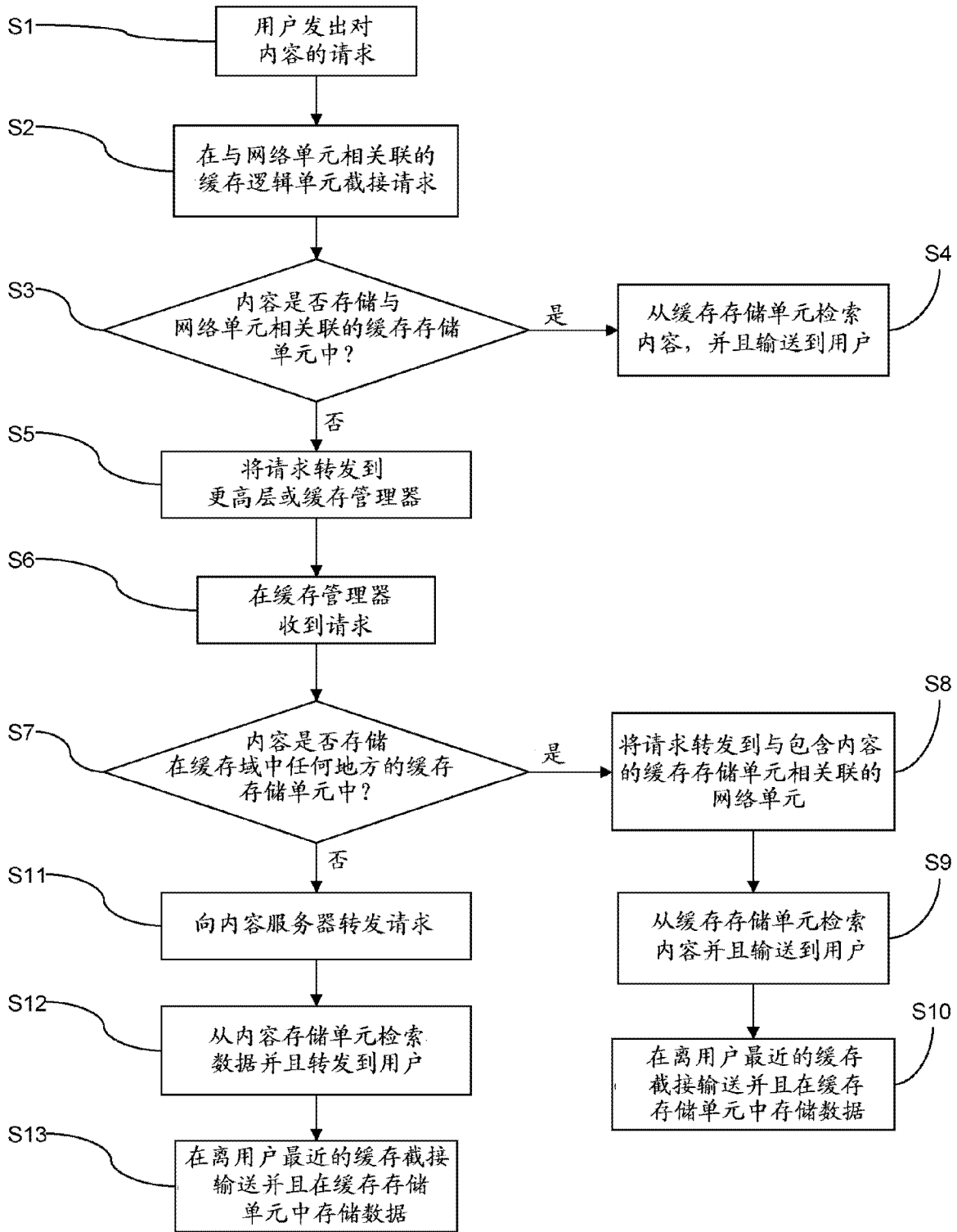


图 3

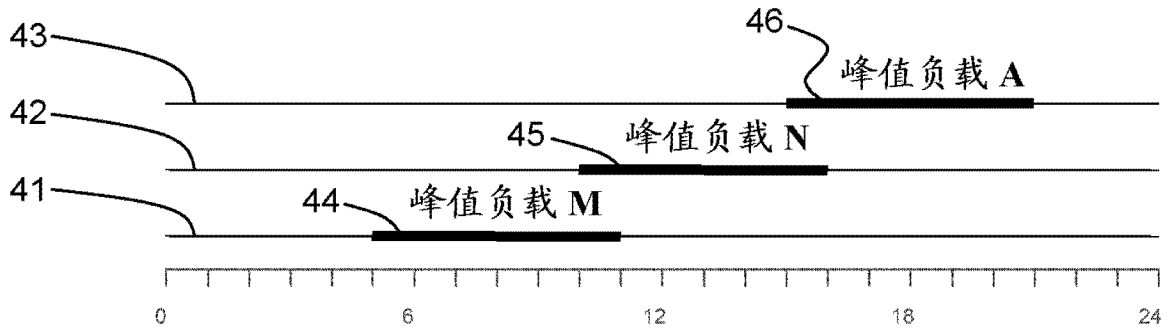


图 4

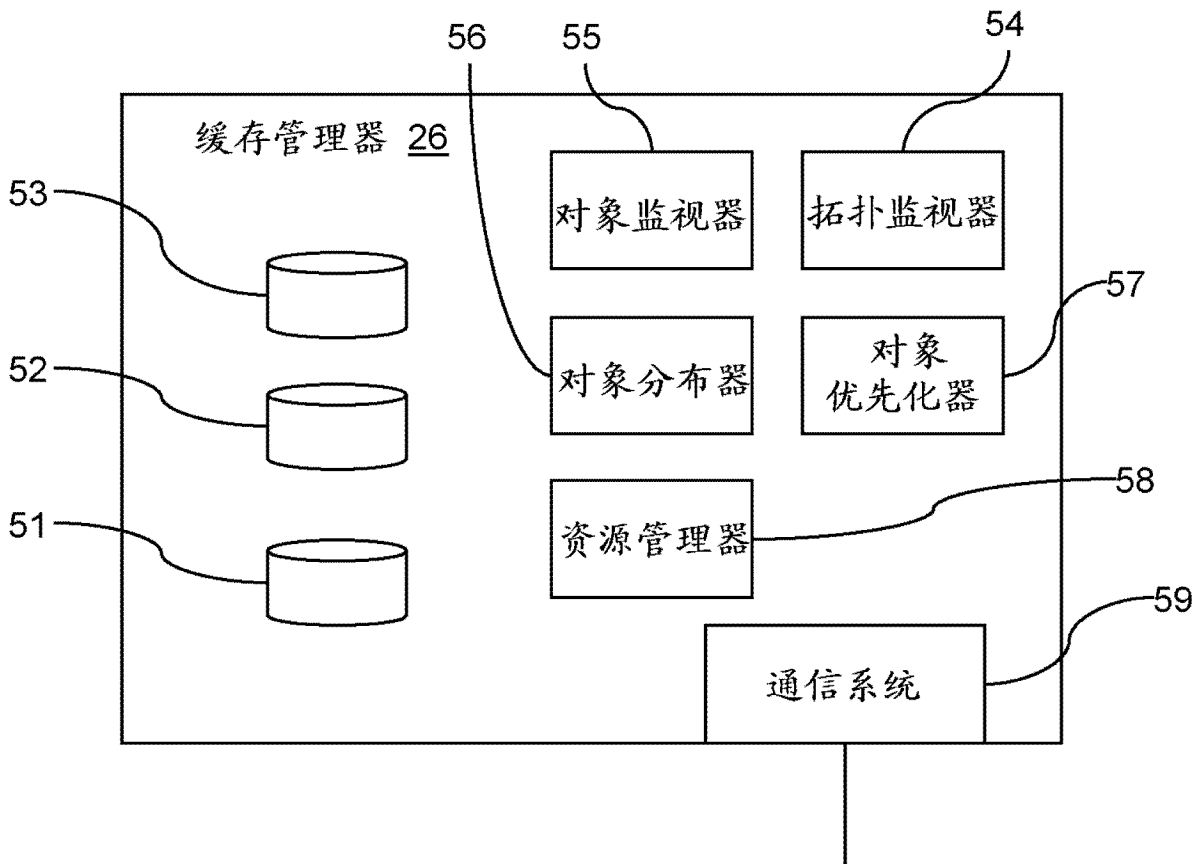


图 5

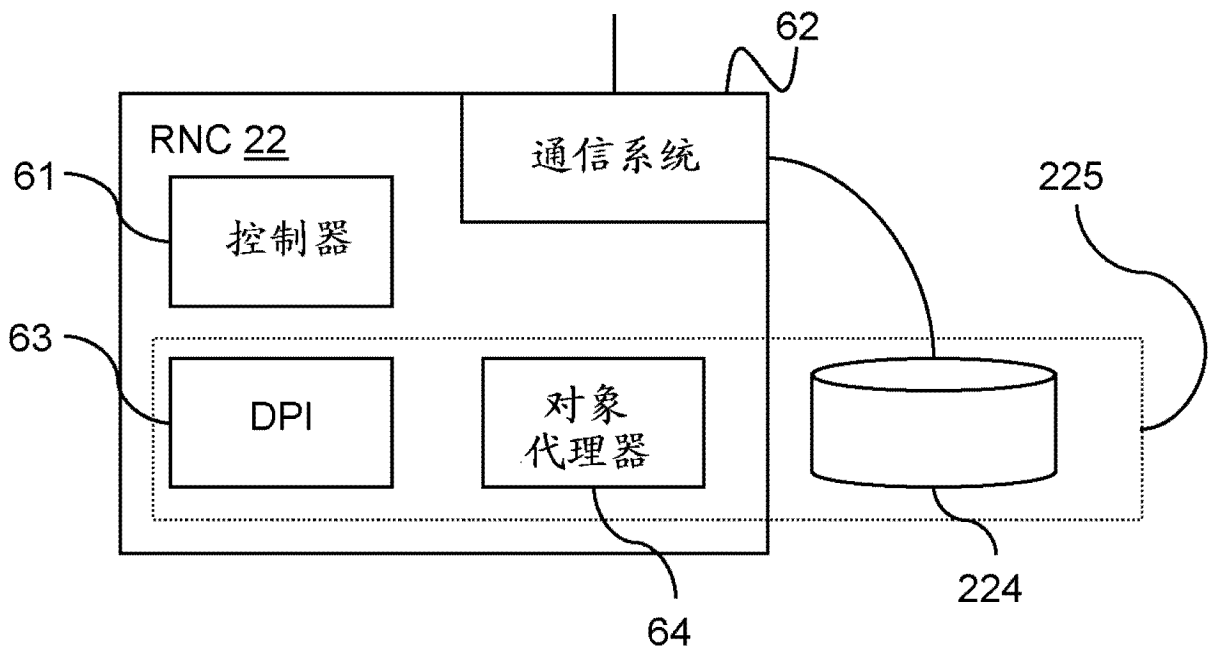


图 6

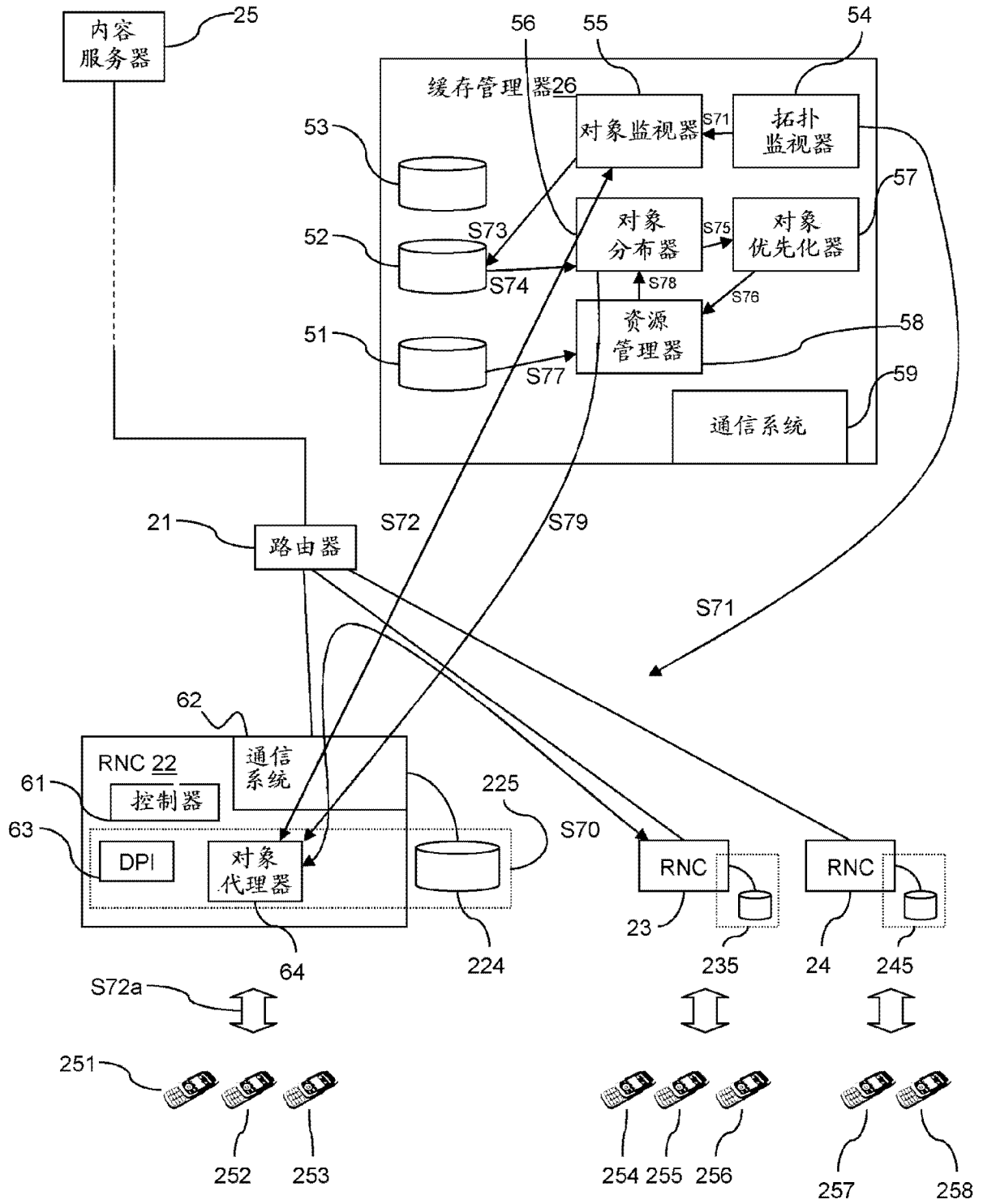


图 7