



(12) 发明专利申请

(10) 申请公布号 CN 102857369 A

(43) 申请公布日 2013. 01. 02

(21) 申请号 201210279783. 2

(22) 申请日 2012. 08. 07

(71) 申请人 北京鼎震科技有限责任公司

地址 102208 北京市昌平区回龙观龙禧二区  
18 号楼 1 单元 102 室

(72) 发明人 李晓亮

(74) 专利代理机构 北京市盛峰律师事务所  
11337

代理人 赵建刚

(51) Int. Cl.

H04L 12/24 (2006. 01)

H04L 29/08 (2006. 01)

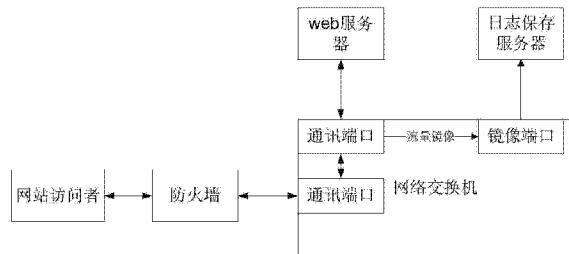
权利要求书 2 页 说明书 6 页 附图 3 页

(54) 发明名称

一种网站日志保存系统及方法和装置

(57) 摘要

本发明提供一种基于旁路镜像的网站日志保存系统及方法和装置, 从而解决现有技术中存在的问题。本发明采用旁路镜像的方式获取访问数据, 对访问网站的数据进行“旁路镜像”, 获得用户访问网站的原始数据包信息, 经由行为分析模块对访问进行行为分类后可以记录成多种格式的网站日志。本发明的技术方案不会对 web 服务器造成任何负担, 且日志格式与 web 服务器的选择完全无关。传统组网模型就是在网络交换机上接入相关的 WEB 服务器, 由 WEB 服务器实体来完成相关的网站日志保存等功能。



1. 一种网站日志保存系统,包括防火墙、网络交换机和 web 服务器,其特征在于,所述网络交换机为具备镜像端口的网络交换机,所述镜像端口上连接有日志保存服务器;所述镜像端口用于通过流量镜像方式获取连接有所述日志保存服务器的通讯端口的通讯数据。

2. 根据权利要求 1 所述的网站日志保存系统,其特征在于,所述日志保存服务器包括流量采集模块、http 协议分析模块、Request 报文分析模块、Response 报文分析模块、行为分析模块、日志条件检查模块和网站日志保存模块;所述流量采集模块、所述 http 协议分析模块、所述 Request 报文分析模块、所述 Response 报文分析模块、所述行为分析模块、所述日志条件检查模块和所述网站日志保存模块顺序连接。

3. 根据权利要求 1 所述的网站日志保存系统,其特征在于,所述网站日志保存系统还包括网站日志筛选模块,所述网站日志筛选模块用于根据请求端指定的条件对网站日志进行筛选并将筛选结果反馈给所述请求端。

4. 一种应用权利要求 1 或 2 或 3 所述的网站日志保存系统进行日志保存的方法,其特征在于,包括以下步骤:

S1,通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包;

S2,分析所述数据包,从所述数据包中获取 http 协议数据包;

S3,分析所述 http 协议数据包中的 Request 报文数据,得到 Request 报文必要信息;

S4,分析所述 http 协议数据包中的 Response 报文数据,得到 Response 报文必要信息;

S5,分析所述 Request 报文必要信息和 / 或 Response 报文必要信息,得到访问行为类型信息;

S6,用所述 Request 报文必要信息和 / 或 Response 报文必要信息和 / 或访问行为类型与预设条件对比,如果符合所述预设条件则缓存所述 Request 报文,并等待获取与该 Request 报文对应的 Response 报文,当获取到与所述 Request 报文对应的 Response 报文后,则将相互对应的 Request 报文和 Response 报文组成完整的访问过程,并将所述完整的访问过程按照预设格式保存到数据库和 / 或日志文件中形成网站日志。

5. 根据权利要求 4 所述的方法,其特征在于,还包括以下步骤:

S7、根据请求端设置的筛选条件从所述数据库和 / 或日志文件中筛选出符合条件的日志记录,并将该符合条件的日志记录保存为新文件再反馈给请求端。

6. 根据权利要求 4 或 5 所述的方法,其特征在于,所述预设条件、所述预设格式、所述筛选条件均通过 web 页面设置。

7. 根据权利要求 4 或 5 所述的方法,其特征在于,所述 Request 报文必要信息包括访问者的 IP 地址、访问的具体域名、访问的具体 URL、Refrence 信息、UserAgent 和携带的 Cookies;所述 Response 报文必要信息包括应答状态码、携带的内容类型和报文长度。

8. 根据权利要求 4 或 5 所述的方法,其特征在于,

S1 具体为,通过所述镜像端口获取,得到所有发送到所述 web 服务器以及从所述 web 服务器发出的报文,并将所述报文分离成上行和下行流量;和 / 或

S2 具体为,通过对所述上行和下行流量中 TCP 载荷的内容分析区分,获取得到 http 协议报文;和 / 或

S3 具体为,对所述 http 协议报文中的 Request 报文进行解码处理,分离出 Request 必要信息,并将所述 Request 必要信息缓冲;和 / 或

S4 具体为,对所述 http 协议报文中的 Response 报文进行解码处理,分离出 Response 必要信息,并将所述 Response 必要信息缓冲;和/或

S5 具体为,根据所述 Request 报文和所述 Response 报文所携带的信息对访问者的访问行为进行分析,确定所述访问行为的行为类型;和/或

S6 具体为,用所述 Request 必要信息和/或所述 Response 必要信息和/或所述访问行为类型与预设日志条件比对,如果符合所述预设日志条件,则把包含有所述 Request 必要信息的 Request 报文缓存,并等待与该 Request 报文相对应的 Response 报文,当获取到与该 Request 报文对应的 Response 报文后,则将相互对应的 Request 报文中的 Request 必要信息和 Response 报文中的 Response 必要信息合并组成一个完整的访问过程,再根据预设的日志格式和日志条目组合成最终的一条网站日志并写入数据库和/或日志文件中并建立该条网站日志的查询索引。

9. 一种应用权利要求 1 或 2 或 3 所述的网站日志保存系统进行日志保存的装置,其特征在于,包括:

流量采集模块,用于通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包;

http 协议分析模块,用于分析所述数据包,从所述数据包中获取 http 协议数据包;

Request 报文分析模块,用于分析所述 http 协议数据包中的 Request 报文数据,得到 Request 报文必要信息;

Response 报文分析模块,用于分析所述 http 协议数据包中的 Response 报文数据,得到 Response 报文必要信息;

行为分析模块,用于分析所述 Request 报文必要信息和/或 Response 报文必要信息,得到访问行为类型信息;

日志条件检查模块,用于用所述 Request 报文必要信息和/或 Response 报文必要信息和/或访问行为类型与预设条件对比,如果符合所述预设条件则送入下一处理步骤;

网站日志保存模块,用于将完整的访问过程按照预设格式保存到数据库和/或日志文件中形成网站日志。

10. 根据权利要求 9 所述的装置,其特征在于所述装置还包括网站日志筛选模块,所述网站日志筛选模块用于根据指定条件对网站日志进行筛选。

## 一种网站日志保存系统及方法和装置

### 技术领域

[0001] 本发明涉及通信技术领域,尤其涉及一种网站日志保存系统及方法和装置。

### 背景技术

[0002] 网站日志是记录 web 服务器接收处理请求以及运行时错误等各种原始信息的以 .log 结尾的文件。通过网站日志可以了解谁在什么时间使用什么工具访问了网站的哪些内容,它是网站分析和网站数据仓储的最基础来源。因为能够完整无误的保存网站日志成为了保证 web 服务器正常运行的一个必要基础。

[0003] 在现有技术中,网站日志由 WEB 服务器自身记录,当访问产生时 WEB 服务器按照预先设置的日志格式以文本的形式把该次访问的某些信息记录在本地或者某台网络服务器上。

[0004] 但是,不同的 WEB 服务器一般仅支持自己特定的日志格式,如 apache 支持的 NCSA 日志格式和 IIS 支持的 W3C 日志格式,大多数的日志分析工具都提供对 NCSA 和 W3C 至少一种格式的支持。另有一些 WEB 服务器如 nginx 有自己默认的日志格式,一般需要手工配置成 NCSA 格式以方便使用日志分析软件。总体上现有技术存在以下问题:

[0005] 1. 访问日志由 web 服务器负责记录,web 服务器不仅需要响应访客的请求还需要记录访问日志,增加了 web 服务器的负担。获得每一次访问的信息都是由 web 服务器在处理请求时同步进行,影响 web 服务器的性能。

[0006] 2. 日志的格式与使用的 web 服务器有关,这极大限制了网站的日志分析工具的选择范围。传统的网站日志格式受使用的 web 服务器制约,选定了某种服务器也就选定了某种日志格式,或者说为了可以使用某种日志格式不得不选用某种服务器。

[0007] 3. 日志配置过程繁琐复杂,某些 web 服务器甚至仅能透过配置文件才能完成日志配置,这需要有较高的计算机知识才能顺利完成。另外 web 服务器一般不提供对已生成的日志的筛选功能,无法对已生成的日志进行筛选处理。

[0008] 4. 日志记录不具备智能性,现有的网站日志只是单纯的记录 web 报文所携带的固有信息,不具备任何的行为分析能力,不管是攻击还是正常访问对于现有网站日志而言没有什么区别,一般都需要专业技术人员进行分析来推测访问的行为,假如网站遭到攻击,在大量的访问日志中寻找攻击线索犹如大海捞针。

### 发明内容

[0009] 针对传统网站日志模式的上述缺点,本发明的目的在于提供一种基于旁路镜像的网站日志保存系统及方法和装置,从而解决现有技术中存在的前述问题。本发明采用旁路镜像的方式获取访问数据,对访问网站的数据进行“旁路镜像”,获得用户访问网站的原始数据包信息,经由行为分析模块对访问进行行为分类后可以记录成多种格式的网站日志。本发明的技术方案不会对 web 服务器造成任何负担,且日志格式与 web 服务器的选择完全无关。传统组网模型就是在网络交换机上接入相关的 WEB 服务器,由 WEB 服务器实体来完成

成相关的网站日志保存等功能；而本发明的技术组网方案是在交换机上旁路部署了一个设备实体，由该设备实体来完成保存网站日志和查询网站日志的功能，WEB 服务器实体仅需要完成网站的信息应答功能。

[0010] 本发明公开的技术方案具体如下：

[0011] 一种网站日志保存系统，包括防火墙、网络交换机和 web 服务器，所述网络交换机为具备镜像端口的网络交换机，所述镜像端口上连接有日志保存服务器；所述镜像端口用于通过流量镜像方式获取连接有所述日志保存服务器的通讯端口的通讯数据。

[0012] 优选的，所述日志保存服务器包括流量采集模块、http 协议分析模块、Request 报文分析模块、Response 报文分析模块、行为分析模块、日志条件检查模块和网站日志保存模块；所述流量采集模块、所述 http 协议分析模块、所述 Request 报文分析模块、所述 Response 报文分析模块、所述行为分析模块、所述日志条件检查模块和所述网站日志保存模块顺序连接。

[0013] 优选的，所述网站日志保存系统还包括网站日志筛选模块，所述网站日志筛选模块用于根据请求端指定的条件对网站日志进行筛选并将筛选结果反馈给所述请求端。

[0014] 一种应用网站日志保存系统进行日志保存的方法，包括以下步骤：

[0015] S1，通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包；

[0016] S2，分析所述数据包，从所述数据包中获取 http 协议数据包；

[0017] S3，分析所述 http 协议数据包中的 Request 报文数据，得到 Request 报文必要信息；

[0018] S4，分析所述 http 协议数据包中的 Response 报文数据，得到 Response 报文必要信息；

[0019] S5，分析所述 Request 报文必要信息和 / 或 Response 报文必要信息，得到访问行为类型信息；

[0020] S6，用所述 Request 报文必要信息和 / 或 Response 报文必要信息和 / 或访问行为类型与预设条件对比，如果符合所述预设条件则缓存所述 Request 报文，并等待获取与所述 Request 报文对应的 Response 报文，当获取到与所述 Request 报文对应的 Response 报文后，则将相互对应的 Request 报文和 Response 报文组成完整的访问过程，并将所述完整的访问过程按照预设格式保存到数据库和 / 或日志文件中形成网站日志。

[0021] 优选的，还包括以下步骤：

[0022] S7、根据请求端设置的筛选条件从所述数据库和 / 或日志文件中筛选出符合条件的日志记录，并将该符合条件的日志记录保存为新文件再反馈给请求端。

[0023] 优选的，所述预设条件、所述预设格式、所述筛选条件均通过 web 页面设置。

[0024] 优选的，所述 Request 报文必要信息包括访问者的 IP 地址、访问的具体域名、访问的具体 URL、Refrence 信息、UserAgent 和携带的 Cookies；所述 Response 报文必要信息包括应答状态码、携带的内容类型和报文长度。

[0025] 优选的，

[0026] S1 具体为，通过所述镜像端口获取，得到所有发送到所述 web 服务器以及从所述 web 服务器发出的报文，并将所述报文分离成上行和下行流量；和 / 或

[0027] S2 具体为，通过对所述上行和下行流量中 TCP 载荷的内容分析区分，获取得到

http 协议报文 ;和 / 或

[0028] S3 具体为,对所述 http 协议报文中的 Request 报文进行解码处理,分离出 Request 必要信息,并将所述 Request 必要信息缓冲 ;和 / 或

[0029] S4 具体为,对所述 http 协议报文中的 Response 报文进行解码处理,分离出 Response 必要信息,并将所述 Response 必要信息缓冲 ;和 / 或

[0030] S5 具体为,根据所述 Request 报文和所述 Response 报文所携带的信息对访问者的访问行为进行分析,确定所述访问行为的行为类型 ;和 / 或

[0031] S6 具体为,用所述 Request 必要信息和 / 或所述 Response 必要信息和 / 或所述访问行为类型与预设日志条件比对,如果符合所述预设日志条件,则把包含有所述 Request 必要信息的 Request 报文缓存,并等待与该 Request 报文相对应的 Response 报文,当获取到与该 Request 报文对应的 Response 报文后,则将相互对应的 Request 报文中的 Request 必要信息和 Response 报文中的 Response 必要信息合并组成一个完整的访问过程,再根据预设的日志格式和日志条目组合成最终的一条网站日志并写入数据库和 / 或日志文件中并建立该条网站日志的查询索引。

[0032] 一种应用网站日志保存系统进行日志保存的装置,包括 :

[0033] 流量采集模块,用于通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包 ;

[0034] http 协议分析模块,用于分析所述数据包,从所述数据包中获取 http 协议数据包 ;

[0035] Request 报文分析模块,用于分析所述 http 协议数据包中的 Request 报文数据,得到 Request 报文必要信息 ;

[0036] Response 报文分析模块,用于分析所述 http 协议数据包中的 Response 报文数据,得到 Response 报文必要信息 ;

[0037] 行为分析模块,用于分析所述 Request 报文必要信息和 / 或 Response 报文必要信息,得到访问行为类型信息 ;

[0038] 日志条件检查模块,用于用所述 Request 报文必要信息和 / 或 Response 报文必要信息和 / 或访问行为类型与预设条件对比,如果符合所述预设条件则送入下一处理步骤 ;

[0039] 网站日志保存模块,用于将完整的访问过程按照预设格式保存到数据库和 / 或日志文件中形成网站日志。

[0040] 优选的,所述装置还包括网站日志筛选模块,所述网站日志筛选模块用于根据指定条件对网站日志进行筛选。

[0041] 本发明的有益效果是 :

[0042] 1. 在记录并保存网站日志的同时,对网站没有任何的影响,无需修改网站任何的配置,无需改写网站的网页,可以做到即插即用 ;

[0043] 2. 本方案由置于旁路设备上的流量采集模块完成数据采集,不会损伤 web 的性能,使 web 服务器可以节省出资源提高并发请求量和计算速度。

[0044] 3. 本方案由行为分析模块对访问行为进行了智能分类,攻击、爬虫、正常访问等一目了然。

[0045] 4. 本方案的日志记录格式与使用什么 web 服务器没有任何关系,使用 apache 服务

器也可以得到 W3C 格式的日志。

[0046] 5. 本发明的日志筛选模块可以直接向用户输出符合用户需求的日志内容。

## 附图说明

[0047] 图 1 是本发明公开的网站日志保存系统结构示意图；

[0048] 图 2 是本发明公开的应用网站日志保存系统进行日志保存的方法的步骤流程图；

[0049] 图 3 是本发明公开的应用网站日志保存系统进行日志保存的装置的示意框图。

## 具体实施方式

[0050] 为了使本发明所解决的技术问题、技术方案及有益效果更加清楚明白，以下结合附图，对本发明进行进一步详细说明。应当理解，此处所描述的具体实施方式仅仅用以解释本发明，并不用于限定本发明。

[0051] 如图 1 所示，本发明公开了一种网站日志保存系统，包括防火墙、网络交换机和 web 服务器，所述网络交换机为具备镜像端口的网络交换机，所述镜像端口上连接有日志保存服务器；所述镜像端口用于通过流量镜像方式获取连接有所述日志保存服务器的通讯端口的通讯数据。所述日志保存服务器包括流量采集模块、http 协议分析模块、Request 报文分析模块、Response 报文分析模块、行为分析模块、日志条件检查模块和网站日志保存模块；所述流量采集模块、所述 http 协议分析模块、所述 Request 报文分析模块、所述 Response 报文分析模块、所述行为分析模块、所述日志条件检查模块和所述网站日志保存模块顺序连接。所述网站日志保存系统还包括网站日志筛选模块，所述网站日志筛选模块用于根据请求端指定的条件对网站日志进行筛选并将筛选结果反馈给所述请求端。

[0052] 如图 2 所示，本发明公开了一种应用网站日志保存系统进行日志保存的方法，包括以下步骤：

[0053] S1，通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包；具体为，通过所述镜像端口获取，得到所有发送到所述 web 服务器以及从所述 web 服务器发出的报文，并将所述报文分离成上行和下行流量；

[0054] S2，分析所述数据包，从所述数据包中获取 http 协议数据包；具体为，通过对所述上行和下行流量中 TCP 载荷的内容分析精确区分属于 http 协议的报文，获取得到 http 协议报文；因为 http 协议是由 Request 报文发起的，因此 http 协议分析系统首先分离出 Request 报文，然后再找到针对这个 Request 报文的应答，分别将 Request 报文和 Response 报文传递到 Request 分析系统和 Response 分析系统，并形成 Request 报文和 Response 报文的对应关系。

[0055] S3，分析所述 http 协议数据包中的 Request 报文数据，得到 Request 报文必要信息；具体为，对所述 http 协议报文中的 Request 报文进行解码处理，分离出 Request 必要信息，并将所述 Request 必要信息缓冲；所述 Request 报文必要信息包括访问者的 IP 地址、访问的具体域名、访问的具体 URL、Refrence 信息、UserAgent 和携带的 Cookies 等信息；

[0056] S4，分析所述 http 协议数据包中的 Response 报文数据，得到 Response 报文必要信息；具体为，对所述 http 协议报文中的 Response 报文进行解码处理，分离出 Response 必要信息，并将所述 Response 必要信息缓冲；所述 Response 报文必要信息包括应答状态码、

携带的内容类型和报文长度等信息。

[0057] S5,分析所述 Request 报文必要信息和 / 或 Response 报文必要信息,得到访问行为类型信息 ;具体为,根据所述 Request 报文和所述 Response 报文所携带的信息对访问者的访问行为进行分析,确定所述访问行为的行为类型 ;所述访问行为类型包括 :正常访问、爬虫和攻击等多种行为类型。

[0058] S6,用所述 Request 报文必要信息和 / 或 Response 报文必要信息和 / 或访问行为类型与预设条件对比,如果符合所述预设条件则缓存所述 Request 报文,并等待获取与该 Request 报文对应的 Response 报文,当获取到与所述 Request 报文对应的 Response 报文后,则将相互对应的 Request 报文和 Response 报文组成完整的访问过程,并将所述完整的访问过程按照预设格式保存到数据库和 / 或文件中形成网站日志 ;具体为,用所述 Request 必要信息和 / 或所述 Response 必要信息和 / 或所述访问行为类型与预设日志条件比对,如果符合所述预设日志条件,则把包含有所述 Request 必要信息的 Request 报文缓存,并等待与该 Request 报文相对应的 Response 报文,当获取到与该 Request 报文对应的 Response 报文后,则将相互对应的 Request 报文中的 Request 必要信息和 Response 报文中的 Response 必要信息合并组成一个完整的访问过程,再根据预设的日志格式和日志条目组合成最终的一条网站日志并写入数据库和 / 或日志文件中并建立该条网站日志的查询索引。

[0059] 为了让获取保存的网站日志具有更大的可用性,在通过上述步骤保存网站日志后,还可以通过以下步骤对日志进行筛选。

[0060] S7、根据请求端设置的筛选条件从所述数据库和 / 或文件中筛选出符合条件的日志记录,并将该符合条件的日志记录保存为新文件再反馈给请求端。

[0061] 所述日志格式 :一条日志中需要记录的条目、条目的出现顺序及其格式。目前常见的网站日志格式主要有 NCSA 日志格式和 W3C 日志格式,分别被 apache 和 IIS 采用,这两种格式下又有更细的分类不做介绍。

[0062] 另外由于本方案中使用了一台专用的日志保存服务器做为日志保存设备,所以通过该服务器上的 web 管理页面就可以对所述预设条件、所述预设格式、所述筛选条件等进行设置。所述预设格式可以是 NCSA common, NCSA combined, W3C 模版, Apache 自定义和 W3C 自定义格式等,所述筛选条件可以是响应状态 (如 200, 304)、请求方法 (如 Get)、源 IP、目的 IP、排除 IP、URL 规则、内容类型 (如图片) 和行为分类 (如正常访问、爬虫、攻击等) 等 ;这些条件也可以组合使用。通过方便的设置筛选条件,进而可以快速获取所需要的日志内容,从而不必像大海捞针一样的查找日志,提高了工作效率。

[0063] 如图 3 所示,本发明公开了一种应用网站日志保存系统进行日志保存的装置,包括 :

[0064] 流量采集模块,用于通过所述镜像端口获取所述 web 服务器收到和发出的全部数据包 ;

[0065] http 协议分析模块,用于分析所述数据包,从所述数据包中获取 http 协议数据包 ;

[0066] Request 报文分析模块,用于分析所述 http 协议数据包中的 Request 报文数据,得到 Request 报文必要信息 ;

[0067] Response 报文分析模块,用于分析所述 http 协议数据包中的 Response 报文数据,



得到 Response 报文必要信息；

[0068] 行为分析模块,用于分析所述 Request 报文必要信息和 / 或 Response 报文必要信息,得到访问行为类型信息；

[0069] 日志条件检查模块,用于用所述 Request 报文必要信息和 / 或 Response 报文必要信息和 / 或访问行为类型与预设条件对比,如果符合所述预设条件则送入下一处理步骤；

[0070] 网站日志保存模块,用于将完整的访问过程按照预设格式保存到数据库和 / 或日志文件中形成网站日志。

[0071] 还包括网站日志筛选模块,所述网站日志筛选模块用于根据指定条件对网站日志进行筛选。

[0072] 通过采用本发明公开的上述技术方案,得到了如下有益的效果：

[0073] 1. 在记录并保存网站日志的同时,对网站没有任何的影响,无需修改网站任何的配置,无需改写网站的网页,可以做到即插即用；

[0074] 2. 本方案由置于旁路设备上的流量采集模块完成数据采集,不会损伤 web 的性能,使 web 服务器可以节省出资源提高并发请求量和计算速度。

[0075] 3. 本方案由行为分析模块对访问行为进行了智能分类,攻击、爬虫、正常访问等一目了然。

[0076] 4. 本方案的日志记录格式与使用什么 web 服务器没有任何关系,使用 apache 服务器也可以得到 W3C 格式的日志。

[0077] 本发明的日志筛选模块可以直接向用户输出符合用户需求的日志内容。

[0078] 以上所述仅是本发明的优选实施方式,应当指出,对于本技术领域的普通技术人员来说,在不脱离本发明原理的前提下,还可以做出若干改进和润饰,这些改进和润饰也应视本发明的保护范围。

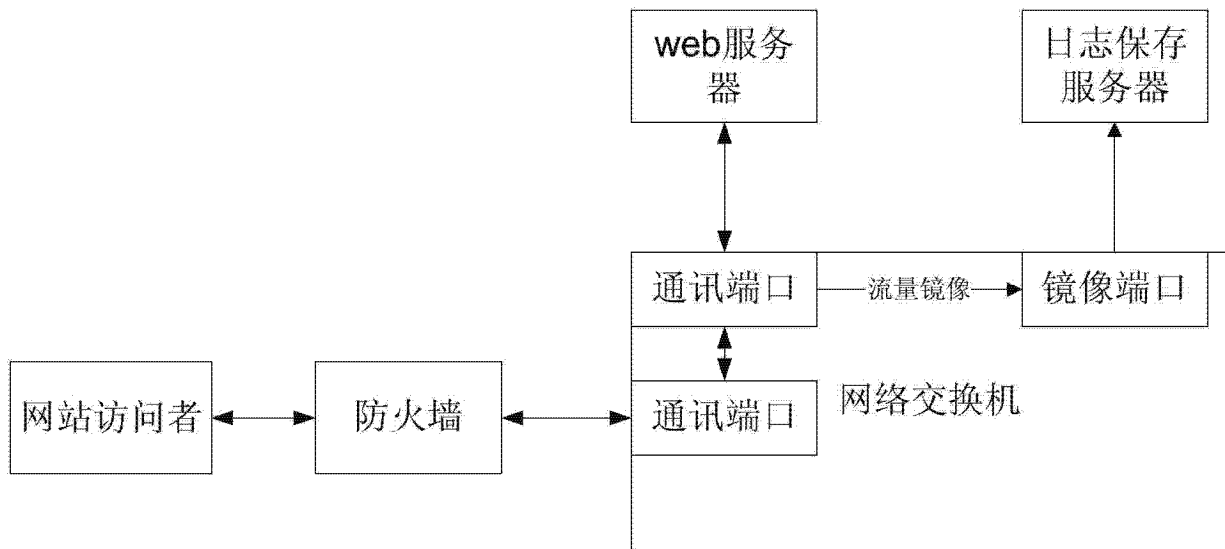


图 1

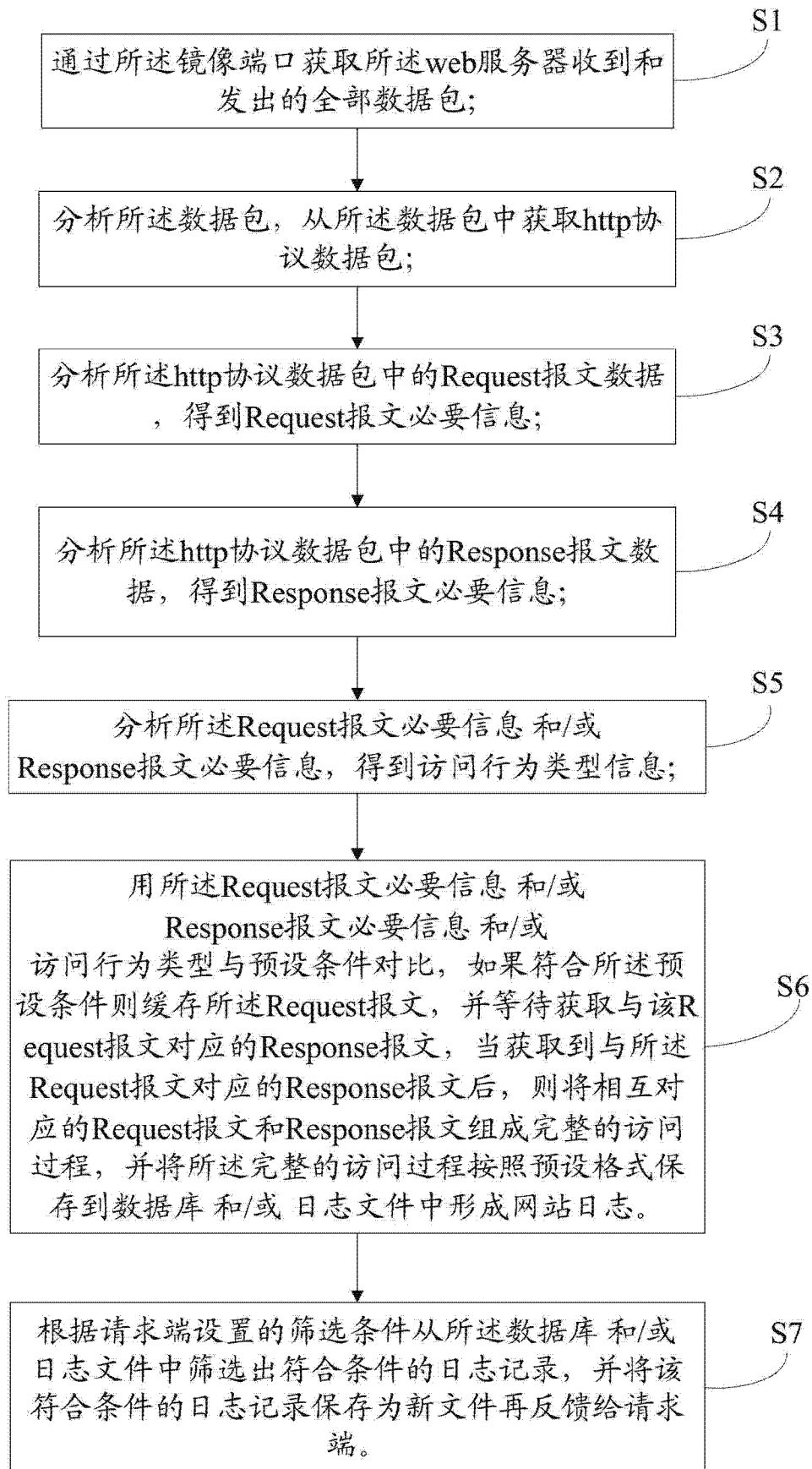


图 2

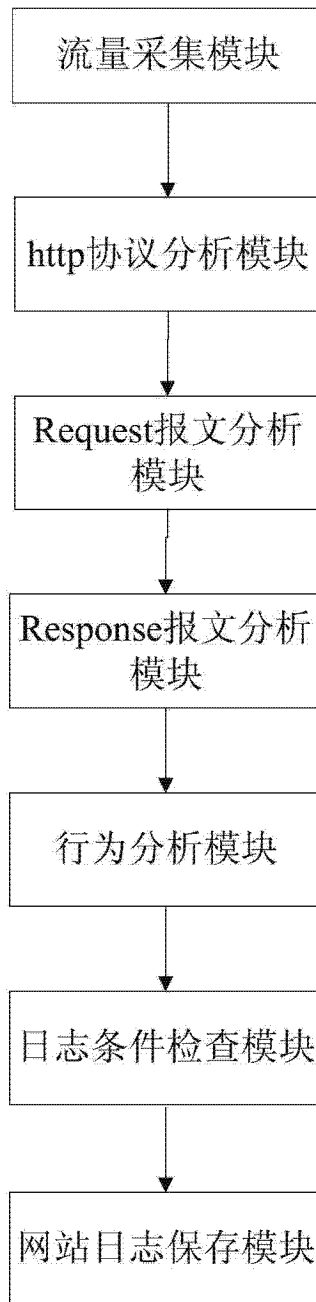


图 3