



(12) 发明专利申请

(10) 申请公布号 CN 118827151 A

(43) 申请公布日 2024. 10. 22

(21) 申请号 202410783628.7

(22) 申请日 2024.06.17

(71) 申请人 中信国际电讯(信息技术)有限公司

地址 中国香港鲗鱼涌英皇道979号太古坊  
林肯大厦20楼

申请人 中企网络通信技术有限公司

(72) 发明人 李超群

(74) 专利代理机构 北京三聚阳光知识产权代理  
有限公司 11250

专利代理师 胡晓静

(51) Int. Cl.

H04L 9/40 (2022.01)

G06N 20/00 (2019.01)

G06F 18/214 (2023.01)

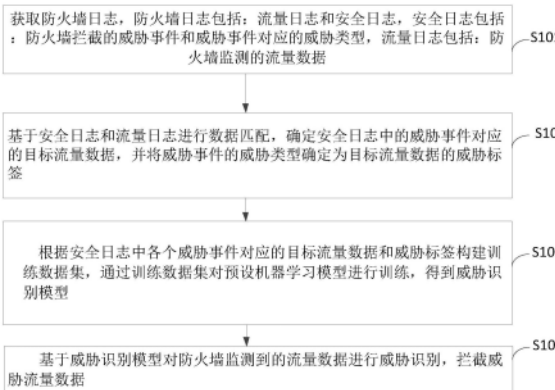
权利要求书2页 说明书11页 附图3页

(54) 发明名称

一种威胁流量数据的识别方法、装置、设备及介质

(57) 摘要

本发明涉及网络安全技术领域,公开了一种威胁流量数据的识别方法、装置、设备及介质,该方法包括:获取包括流量日志和安全日志的防火墙日志,安全日志包括:防火墙拦截的威胁事件和威胁事件对应的威胁类型,流量日志包括:防火墙监测的流量数据;基于安全日志和流量日志进行数据匹配,确定安全日志中的威胁事件对应的目标流量数据,并将威胁事件的威胁类型确定为该目标流量数据的威胁标签;根据安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过训练数据集对预设机器学习模型进行训练,得到威胁识别模型;基于威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据,本发明通过防火墙日志中的数据训练机器学习模型,将训练得到威胁识别模型来识别威胁流量数据,可以更全面的拦截威胁流量以保证网络安全。



1. 一种威胁流量数据的识别方法,其特征在于,所述方法包括:

获取防火墙日志,所述防火墙日志包括:流量日志和安全日志,所述安全日志包括:防火墙拦截的威胁事件和所述威胁事件对应的威胁类型,所述流量日志包括:防火墙监测的流量数据;

基于所述安全日志和流量日志进行数据匹配,确定所述安全日志中的威胁事件对应的目标流量数据,并将所述威胁事件的威胁类型确定为所述目标流量数据的威胁标签;

根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过所述训练数据集对预设机器学习模型进行训练,得到威胁识别模型;

基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

2. 根据权利要求1所述的方法,其特征在于,所述防火墙日志包括:多个不同防火墙对应的日志;

在基于所述安全日志和流量日志进行数据匹配前,所述方法还包括:

根据所述防火墙日志中各个防火墙对应的安全日志,确定各个威胁事件的非结构化文本数据;

对所述各个威胁事件的非结构化文本数据进行数据解析,得到统一结构格式下的各个威胁事件对应的结构化数据,所述数据解析方式为:基于最长公共子序列的流式日志解析方法。

3. 根据权利要求1所述的方法,其特征在于,所述根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,包括:

根据安全日志中各个威胁事件对应的属性信息,确定所述目标流量数据的特征维度;

提取所述目标流量数据中各个特征维度的样本特征,基于所述样本特征和所述目标流量数据对应威胁事件的威胁标签确定各个目标流量数据对应的样本数据;

根据各个目标流量数据对应的样本数据构建训练数据集。

4. 根据权利要求1所述的方法,其特征在于,所述预设机器学习模型为:基于轻量级梯度提升机算法构建的机器学习模型;

所述通过所述训练数据集对预设机器学习模型进行训练,包括:

在通过所述训练数据集对所述基于轻量级梯度提升机算法构建的机器学习模型进行训练的过程中,对所述机器学习模型进行交叉验证;

根据交叉验证结果调整所述机器学习模型的模型参数。

5. 根据权利要求3所述的方法,其特征在于,在基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据后,所述方法还包括:

根据模型解释器确定所述威胁识别模型在对所述威胁流量数据进行识别时,各个特征维度对应的贡献值;

输出所述威胁流量数据的各个特征维度对应的贡献值。

6. 根据权利要求1所述的方法,其特征在于,所述基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,包括:

通过防火墙对监测到的流量数据进行威胁识别,拦截所述流量数据中的第一威胁流量数据,排除所述流量数据中的第一威胁流量数据,得到第二流量数据;

通过所述威胁识别模型对所述第二流量数据进行威胁识别,拦截所述第二流量数据中的第二威胁流量数据。

7. 根据权利要求1所述的方法,其特征在于,所述方法还包括:

输出所述威胁流量数据对应的威胁事件和威胁类型。

8. 一种威胁流量数据的识别装置,其特征在于,所述装置包括:

日志数据获取模块,用于获取防火墙日志,所述防火墙日志包括:流量日志和安全日志,所述安全日志包括:防火墙拦截的威胁事件和所述威胁事件对应的威胁类型,所述流量日志包括:防火墙监测的流量数据;

威胁标签匹配模块,用于基于所述安全日志和流量日志进行数据匹配,确定所述安全日志中的威胁事件对应的目标流量数据,并将所述威胁事件的威胁类型确定为所述目标流量数据的威胁标签;

威胁模型训练模块,用于根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过所述训练数据集对预设机器学习模型进行训练,得到威胁识别模型;

威胁流量拦截模块,用于基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

9. 一种计算机设备,其特征在于,包括:

存储器和处理器,所述存储器和所述处理器之间互相通信连接,所述存储器中存储有计算机指令,所述处理器通过执行所述计算机指令,从而执行权利要求1至7中任一项所述的威胁流量数据的识别方法。

10. 一种计算机可读存储介质,其特征在于,所述计算机可读存储介质上存储有计算机指令,所述计算机指令用于使计算机执行权利要求1至7中任一项所述的威胁流量数据的识别方法。

## 一种威胁流量数据的识别方法、装置、设备及介质

### 技术领域

[0001] 本发明涉及网络安全技术领域,具体涉及一种威胁流量数据的识别方法、装置、设备及介质。

### 背景技术

[0002] 随着互联网的飞速发展,人们的生活也越来越依赖于网络,网络威胁的种类也越来越多,威胁网络流量很有可能造成网络用户的设备安全、财产安全受到影响,因此如何识别网络流量数据中的威胁流量变得越来越重要。

[0003] 现有技术中,通常是采用防火墙对威胁流量数据进行拦截,然而防火墙拦截方式通常是经验确定预设规则,以对采集到的流量数据进行威胁判断,从而拦截威胁流量,这样的拦截方式容易导致部分威胁流量无法被防火墙识别,无法全面高效的对流量数据进行威胁识别。

### 发明内容

[0004] 有鉴于此,本发明提供了一种威胁流量数据的识别方法、装置、设备及介质,以解决无法全面高效的对流量数据进行威胁识别的问题。

[0005] 第一方面,本发明提供了一种威胁流量数据的识别方法,所述方法包括:

[0006] 获取防火墙日志,所述防火墙日志包括:流量日志和安全日志,所述安全日志包括:防火墙拦截的威胁事件和所述威胁事件对应的威胁类型,所述流量日志包括:防火墙监测的流量数据;

[0007] 基于所述安全日志和流量日志进行数据匹配,确定所述安全日志中的威胁事件对应的目标流量数据,并将所述威胁事件的威胁类型确定为所述目标流量数据的威胁标签;

[0008] 根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过所述训练数据集对预设机器学习模型进行训练,得到威胁识别模型;

[0009] 基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

[0010] 通过对防火墙的安全日志中拦截的威胁事件和流量日志中的流量数据进行数据匹配,从而确定威胁事件对应的目标流量数据及其对应的威胁标签,以得到训练数据集,通过该训练数据集对预设的机器学习模型进行训练得到威胁识别模型,来对防火墙监测到的流量数据进行识别拦截,提高了威胁流量识别的全面性,进一步保证了网络安全。

[0011] 在一种可选的实施方式中,所述防火墙日志包括:多个不同防火墙对应的日志;

[0012] 在基于所述安全日志和流量日志进行数据匹配前,所述方法还包括:

[0013] 根据所述防火墙日志中各个防火墙对应的安全日志,确定各个威胁事件的非结构化文本数据;

[0014] 对所述各个威胁事件的非结构化文本数据进行数据解析,得到统一结构格式下的各个威胁事件对应的结构化数据,所述数据解析方式为:基于最长公共子序列的流式日志

解析方法。

[0015] 通过对多个防火墙的安全日志进行解析,从而得到统一结构格式的结构化数据,可以获得更多的威胁事件对应的训练数据,保证了训练得到的威胁识别模型的威胁识别效果,同时对安全日志的数据进行解析,可以提高后续安全日志和流量日志的匹配效果,保证了训练集的数据质量。

[0016] 在一种可选的实施方式中,所述根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,包括:

[0017] 根据安全日志中各个威胁事件对应的属性信息,确定所述目标流量数据的特征维度;

[0018] 提取所述目标流量数据中各个特征维度的样本特征,基于所述样本特征和所述目标流量数据对应威胁事件的威胁标签确定各个目标流量数据对应的样本数据;

[0019] 根据各个目标流量数据对应的样本数据构建训练数据集。

[0020] 通过提取目标流量数据各个特征维度对应的样本特征,并将这些样本特征和目标流量数据对应的威胁标签进行匹配从而得到目标流量数据的样本数据,汇总样本数据得到训练数据集,可以进一步保证训练数据的有效性,提高威胁识别模型的训练效果。

[0021] 在一种可选的实施方式中,所述预设机器学习模型为:基于轻量级梯度提升机算法构建的机器学习模型;

[0022] 所述通过所述训练数据集对预设机器学习模型进行训练,包括:

[0023] 在通过所述训练数据集对所述基于轻量级梯度提升机算法构建的机器学习模型进行训练的过程中,对所述机器学习模型进行交叉验证;

[0024] 根据交叉验证结果调整所述机器学习模型的模型参数。

[0025] 通过构建基于轻量级梯度提升机算法构建的机器学习模型,并在对该模型进行训练的过程中通过交叉验证来调整模型参数,可以进一步提高模型稳定性,保证训练得到的威胁识别模型的威胁识别效果。

[0026] 在一种可选的实施方式中,在基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据后,所述方法还包括:

[0027] 根据模型解释器确定所述威胁识别模型在对所述威胁流量数据进行识别时,各个特征维度对应的贡献值;

[0028] 输出所述威胁流量数据的各个特征维度对应的贡献值。

[0029] 通过采用模型解释器来确定威胁识别模型对威胁流量数据进行识别时各个特征维度的贡献度并进行输出,可以使得用户了解具体是哪些维度的原因导致的威胁流量产生,从而及时网络设备进行相应的安全防控。

[0030] 在一种可选的实施方式中,所述基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,包括:

[0031] 通过防火墙对监测到的流量数据进行威胁识别,拦截所述流量数据中的第一威胁流量数据,排除所述流量数据中的第一威胁流量数据,得到第二流量数据;

[0032] 通过所述威胁识别模型对所述第二流量数据进行威胁识别,拦截所述第二流量数据中的第二威胁流量数据。

[0033] 通过防火墙来对流量数据进行初次拦截,接着通过威胁识别模型对拦截后的流量

数据再次进行危险排查,从而实现对流量数据的双重筛查,进一步保证威胁拦截的全面性,保证网络安全。

[0034] 在一种可选的实施方式中,所述方法还包括:

[0035] 输出所述威胁流量数据对应的威胁事件和威胁类型。

[0036] 通过将威胁流量数据的对应的威胁事件和威胁类型进行输出显示,可以使得用户及时了解当前的网络安全状态,从而在后续的网络行为中进行规避。

[0037] 第二方面,本发明提供了一种威胁流量数据的识别装置,所述装置包括:

[0038] 日志数据获取模块,用于获取防火墙日志,所述防火墙日志包括:流量日志和安全日志,所述安全日志包括:防火墙拦截的威胁事件和所述威胁事件对应的威胁类型,所述流量日志包括:防火墙监测的流量数据;

[0039] 威胁标签匹配模块,用于基于所述安全日志和流量日志进行数据匹配,确定所述安全日志中的威胁事件对应的目标流量数据,并将所述威胁事件的威胁类型确定为所述目标流量数据的威胁标签;

[0040] 威胁模型训练模块,用于根据所述安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过所述训练数据集对预设机器学习模型进行训练,得到威胁识别模型;

[0041] 威胁流量拦截模块,用于基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

[0042] 第三方面,本发明提供了一种计算机设备,包括:存储器和处理器,存储器和处理器之间互相通信连接,存储器中存储有计算机指令,处理器通过执行计算机指令,从而执行上述第一方面或其对应的任一实施方式的威胁流量数据的识别方法。

[0043] 第四方面,本发明提供了一种计算机可读存储介质,该计算机可读存储介质上存储有计算机指令,计算机指令用于使计算机执行上述第一方面或其对应的任一实施方式的威胁流量数据的识别方法。

## 附图说明

[0044] 为了更清楚地说明本发明具体实施方式或现有技术中的技术方案,下面将对具体实施方式或现有技术描述中所需要使用的附图作简单地介绍,显而易见地,下面描述中的附图是本发明的一些实施方式,对于本领域普通技术人员来讲,在不付出创造性劳动的前提下,还可以根据这些附图获得其他的附图。

[0045] 图1是根据本发明实施例的一种威胁流量数据的识别方法的流程示意图;

[0046] 图2是根据本发明实施例的另一威胁流量数据的识别方法的流程示意图;

[0047] 图3是根据本发明实施例的一种威胁流量数据的识别装置的结构框图;

[0048] 图4是本发明实施例的计算机设备的硬件结构示意图。

## 具体实施方式

[0049] 为使本发明实施例的目的、技术方案和优点更加清楚,下面将结合本发明实施例中的附图,对本发明实施例中的技术方案进行清楚、完整地描述,显然,所描述的实施例是本发明一部分实施例,而不是全部的实施例。基于本发明中的实施例,本领域技术人员在没

有做出创造性劳动前提下所获得的所有其他实施例,都属于本发明保护的范围。

[0050] 随着互联网的飞速发展,人们的生活也越来越依赖于网络,网络威胁的种类也越来越多,威胁网络流量很有可能造成网络用户的设备安全、财产安全受到影响,因此如何识别网络流量数据中的威胁流量变得越来越重要。

[0051] 现有技术中,通常是采用防火墙对威胁流量数据进行拦截,然而防火墙拦截方式通常是经验确定预设规则,以对采集到的流量数据进行威胁判断,从而拦截威胁流量,这样的拦截方式容易导致部分威胁流量无法被防火墙识别,无法全面高效的对流量数据进行威胁识别。

[0052] 为此,本发明实施例提供了一种威胁流量数据的识别方法,通过对防火墙的安全日志中拦截的威胁事件和流量日志中的流量数据进行匹配,从而得到威胁事件对应的目标流量数据,并确定其对应的威胁标签以得到训练数据集,通过该训练数据集对预设的机器学习模型进行训练得到威胁识别模型,来对防火墙监测到的流量数据进行识别拦截,提高了威胁流量识别的全面性,进一步保证了网络安全。

[0053] 根据本发明实施例,提供了一种威胁流量数据的识别方法实施例,需要说明的是,在附图的流程图示出的步骤可以在诸如一组计算机可执行指令的计算机系统中执行,并且,虽然在流程图中示出了逻辑顺序,但是在某些情况下,可以以不同于此处的顺序执行所示出或描述的步骤。

[0054] 在本实施例中提供了一种威胁流量数据的识别方法,可用于上述的网络威胁识别,图1是根据本发明实施例的一种威胁流量数据的识别方法的流程图,如图1所示,该流程包括如下步骤:

[0055] 步骤S101,获取防火墙日志,防火墙日志包括:流量日志和安全日志,安全日志包括:防火墙拦截的威胁事件和威胁事件对应的威胁类型,流量日志包括:防火墙监测的流量数据。

[0056] 在需要连接网络的设备上,用户通常会安装现有的防火墙软件来对设备进行网络防护,这些防火墙软件可以对用户设备的流量数据进行监测,从而生成流量日志,在流量日志中会记录用户设备在一段事件内各种流量行为所对应的流量数据。对于这些监测到的流量数据,防火墙可以按照设置好的拦截规则进行威胁拦截,例如,根据预设规则对流量数据的五元组进行处理,来确定其是否为威胁事件从而进行拦截。因此防火墙会将一些流量数据所对应的网络事件判定为威胁事件,并根据设置好的规则来确定该威胁事件所对应的威胁类型,例如,病毒、危险网站、异常软件安装等威胁类型。在防火墙识别并拦截对应的威胁事件后,会生成对应的安全日志,在安全日志中会记录拦截的威胁事件和这些威胁事件对应的威胁类型。

[0057] 步骤S102,基于安全日志和流量日志进行数据匹配,确定安全日志中的威胁事件对应的目标流量数据,并将威胁事件的威胁类型确定为目标流量数据的威胁标签。

[0058] 由于并不是所有流量数据对应的事件都是具有威胁的,在安全日志中记录的威胁事件只是和流量日志中的部分流量数据对应,因此需要对安全日志和流量日志中记录的数据进行匹配,从而确定各个威胁事件对应的目标流量数据,同时根据安全日志中记录的各个威胁事件的威胁类型来为目标流量数据打上威胁标签,即对目标流量数据进行数据标注,从而得到训练数据。

[0059] 步骤S103,根据安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过训练数据集对预设机器学习模型进行训练,得到威胁识别模型。

[0060] 在确定安全日志中所有威胁事件对应的目标流量数据并对这些目标流量数据进行数据标注后,将这些标注后的目标流量数据进行汇总,构建训练数据集。通过该训练数据集对预先建立的机器学习模型进行训练,在训练过程中调整模型的训练参数,从而得到威胁识别模型。具体的,可以将训练数据集划分为训练集和测试集,通过测试集来对训练后的模型进行测试,确保训练得到的威胁识别模型满足要求,具有预设要求的威胁识别能力。

[0061] 步骤S104,基于威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

[0062] 在训练得到威胁识别模型后,通过该威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截流量数据中于威胁事件对应的威胁流量数据。

[0063] 本实施例提供的威胁流量数据的识别方法,通过对防火墙的安全日志中拦截的威胁事件和流量日志中的流量数据进行匹配,从而得到威胁事件对应的目标流量数据,并确定其对应的威胁标签以得到训练数据集,通过该训练数据集对预设的机器学习模型进行训练得到威胁识别模型,来对防火墙监测到的流量数据进行识别拦截,提高了威胁流量识别的全面性,进一步保证了网络安全。

[0064] 根据本发明实施例,提供了另一种威胁流量数据的识别方法实施例,可用于上述的威胁识别,图2是根据本发明实施例的另一威胁流量数据的识别方法的流程图,如图2所示,该流程包括如下步骤:

[0065] 步骤S201,获取防火墙日志,防火墙日志包括:流量日志和安全日志,安全日志包括:防火墙拦截的威胁事件和威胁事件对应的威胁类型,流量日志包括:防火墙监测的流量数据。

[0066] 具体的,防火墙日志包括:多个不同防火墙对应的日志。

[0067] 可以理解为,用户为了进一步保证设备的网络安全,可能会在设备上安装不同种类的防火墙软件。这些防火墙软件对应的拦截规则上可能会存在差异,因此可能会拦截到更多的威胁事件,不同的防火墙具有对应的安全日志,用以记录其拦截的威胁事件和对应的威胁类型。同样的,对于不同的防火墙来说,虽然其监测的流量数据基本上不会存在差异,但是各个防火墙也可以生成其对应的流量日志。在后续匹配时,可以将不同防火墙对应的安全日志中的威胁事件和流量日志中的流量数据分别进行汇总,再汇总之后在进行数据匹配。通过对多个防火墙日志对应的威胁事件和流量数据进行匹配,可以得到更加丰富的训练数据集,进一步提升后续模型的训练效果。

[0068] 步骤S202,根据防火墙日志中各个防火墙对应的安全日志,确定各个威胁事件的非结构化文本数据;

[0069] 对各个威胁事件的非结构化文本数据进行数据解析,得到统一结构格式下的各个威胁事件对应的结构化数据,数据解析方式为:基于最长公共子序列的流式日志解析方法。

[0070] 在防火墙的安全日志中的数据记录格式通常是非结构化文本数据,即安全日志中记录的威胁事件以非结构化文本数据的形式存在。在防火墙日志中确定各个防火墙对应的安全日志,并确定各个安全日志中的各个威胁事件对应的非结构化文本数据。由于防火墙的类型差异,这些非结构化文本数据在格式上也会存在差异,因此需要对这些数据进行格



式统一,可以通过基于最长公共子序列的流式日志解析方法来对这些威胁事件对应的数据进行解析,得到统一结构格式下的各个威胁事件对应的结构化数据,从而保证后续安全日志和流量日志的匹配准确性。

[0071] 通过对多个防火墙的安全日志进行解析,从而得到统一结构格式的结构化数据,可以获得更多的威胁事件对应的训练数据,保证了训练得到的威胁识别模型的威胁识别效果,并且对安全日志的数据进行解析,可以提升后续安全日志和流量日志的匹配效果,保证了训练集的数据质量。

[0072] 在一个例子中,通过基于最长公共子序列的流式日志解析方法对安全日志进行解析的流程可以为:步骤1,对安全日志的数据进行初始化处理,确定日志对象LCS Object、日志模板LCS seq和行数列表line Ids,以及确定存放日志对象的列表LCS Map。步骤2,采用流式的读取方式对日志进行读取。步骤3,当读取到新的日志条目后遍历LCS Map,寻找该日志与所有LCS Object的最大公共子序列,如果子序列的长度大于日志序列长度的一半,则认为该日志该与日志键匹配。如果找到匹配的日志对象,执行步骤5,如果没有,或者LCS Map为空,执行步骤4。步骤4,将该行日志初始化为一个新的LCS Object,放入列表LCS Map中。步骤5,将该行日志更新到匹配的LCS Object的行数列表line Ids中,并且更新LCS seq。步骤6,跳转到步骤2,直到日志读取完毕。该日志解析方法只是一种示例性的日志解析方式,具体实施时,可以根据实际情况确定合适的日志解析方式。

[0073] 步骤S203,基于安全日志和流量日志进行数据匹配,确定安全日志中的威胁事件对应的目标流量数据,并将威胁事件的威胁类型确定为目标流量数据的威胁标签。具体实施方式参考图1所示实施例的步骤S202,此处不在进行赘述。

[0074] 步骤S204,根据安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过训练数据集对预设机器学习模型进行训练,得到威胁识别模型。

[0075] 具体的,在步骤S204中包括:

[0076] 步骤S2041,根据安全日志中各个威胁事件对应的属性信息,确定目标流量数据的特征维度。

[0077] 在安全日志中会记录各个威胁事件对应的属性信息,例如,威胁事件对应的访问IP地址、访问方向、危险等级等属性,结合这些属性对应的具体信息,确定其对应目标流量数据的特征维度,例如,特征维度可以是源ip、源端口、目的IP、目的端口、流长度的标准差等特征维度。

[0078] 步骤S2042,提取目标流量数据中各个特征维度的样本特征,基于样本特征和目标流量数据对应威胁事件的威胁标签确定各个目标流量数据对应的样本数据。

[0079] 在确定特征维度后,对各个目标流量数据进行特征提取,得到目标流量数据在各个特征维度对应的样本特征数据,将这些不同特征维度的样本特征进行汇聚,同时结合该目标流量数据对应的威胁事件的威胁标签,将这些样本特征数据进行标签标注,得到各个目标流量数据对应的样本数据。

[0080] 步骤S2043,根据各个目标流量数据对应的样本数据构建训练数据集,通过训练数据集对预设机器学习模型进行训练,得到威胁识别模型。

[0081] 将各个目标流量数据对应的样本数据进行汇总,得到训练数据集,在该训练数据集中的每条样本数据都是对目标流量数据进行特征提取后的数据,可以进一步保证机器学

习的训练效率以及识别效果。在对预设的机器学习模型进行训练、测试等常规训练流程后，得到训练好的威胁识别模型。

[0082] 通过提取目标流量数据各个特征维度对应的样本特征，并将这些样本特征和目标流量数据对应的威胁标签进行匹配从而得到目标流量数据的样本数据，汇总样本数据得到训练数据集，可以进一步保证训练数据的有效性，提高威胁识别模型的训练效果。

[0083] 在一个例子中，可以对上述得到的训练数据集进行预处理，具体可以包括：填充数据缺失值、对缺失值进行处理、对特征进行相关性分析从而进行特征筛选以及对分类变量和IP地址进行编码并进行数据类型转换等预处理操作，保证构造的训练数据集质量。

[0084] 具体的，在步骤S204中，预设机器学习模型为：基于轻量级梯度提升机算法构建的机器学习模型；

[0085] 通过训练数据集对预设机器学习模型进行训练，包括：

[0086] 在通过训练数据集对基于轻量级梯度提升机算法构建的机器学习模型进行训练的过程中，对机器学习模型进行交叉验证；

[0087] 根据交叉验证结果调整机器学习模型的模型参数。

[0088] 轻量级梯度提升机算法即为Light GBM算法，是一种用于解决分类和回归问题的梯度提升机算法，基于该算法构建的机器学习模型对于威胁识别具有较高的准确性，同时在训练的过程中采用交叉验证的方式来调整参数，可以进一步保证训练得到的模型稳定性，保证对在流量数据进行威胁识别时的效果。

[0089] 步骤S205，基于威胁识别模型对防火墙监测到的流量数据进行威胁识别，拦截威胁流量数据。

[0090] 具体的，在步骤S205中，基于威胁识别模型对防火墙监测到的流量数据进行威胁识别，包括：

[0091] 通过防火墙对监测到的流量数据进行威胁识别，拦截流量数据中的第一威胁流量数据，排除流量数据中的第一威胁流量数据，得到第二流量数据；

[0092] 通过威胁识别模型对第二流量数据进行威胁识别，拦截第二流量数据中的第二威胁流量数据。

[0093] 通过防火墙来对流量数据进行初次拦截，将初次拦截到的威胁事件对应的威胁流量数据进行筛选，得到第二流量数据，接着再通过训练好的威胁识别模型来进行威胁识别，拦截第二流量数据中的第二威胁流量数据，可以实现对流量数据的双重筛查，进一步保证威胁拦截的全面性，保证网络安全。

[0094] 步骤S206，根据模型解释器确定所述威胁识别模型在对所述威胁流量数据进行识别时，各个特征维度对应的贡献值；

[0095] 输出所述威胁流量数据的各个特征维度对应的贡献值。

[0096] 可以理解为，威胁识别模型在识别到威胁流量数据后，可以根据模型解释器来确定该威胁识别模型在对威胁流量数据的进行识别的过程中，各个特征维度的贡献度，即各个维度的因素对于流量数据被识别为威胁流量数据的关键性。并将各个特征维度的贡献度进行输出。可以使得用户可以更直观的了解威胁流量数据具体是由于哪些维度的因素导致的，从而及时做出相应的防控。

[0097] 具体的，模型解释器可以为SHAP解释器，在一个例子中，可以通过SHAP解释器来对

各个威胁流量数据被识别过程中,不同特征维度对应的贡献度。SHAP解释器在进行贡献度计算时的计算原理如下:假设模型的输入有3个特征A、B、C,输出预测值为0或1,可以确定模型输入只有A、B或C时,模型对应的输出预测值;以及模型输入为AB、AC或BC时,模型对应的输出预测值;以及模型输入为ABC时的模型输出预测值。接着通过对比不同子集中模型输出的差异,来计算每个特征的边际效应,最后根据shapley值公式,得到每个特征A,B,C对预测结果的独立贡献,也就是它们的SHAP值。

[0098] 通过采用模型解释器来确定威胁识别模型对威胁流量数据进行识别时各个特征维度的贡献度并进行输出,可以使得用户了解具体是哪些维度的原因导致的威胁流量产生,从而及时网络设备进行相应的安全防控。

[0099] 步骤S207,输出威胁流量数据对应的威胁事件和威胁类型。

[0100] 通过将威胁流量数据的对应的威胁事件和威胁类型进行输出显示,可以使得用户及时了解当前的网络安全状态,从而在后续的网络行为中进行规避。具体的,可以根据预先设置的威胁等级判定规则,来确定威胁流量数据所对应的威胁等级并进行输出。例如,可以预先设置威胁等级表,在该威胁等级表中包括不同威胁等级所对应的威胁事件,根据威胁流量数据所对应的威胁事件结合该威胁等级表,来确定具体的威胁等级,便于用户进一步了解威胁的严重情况,从而及时做出相应防控。

[0101] 本发明实施例提供的威胁流量数据的识别方法,通过对防火墙的安全日志中拦截的威胁事件和流量日志中的流量数据进行匹配,从而得到威胁事件对应的目标流量数据,并确定其对应的威胁标签,同时对目标流量数据进行特征提取,从而得到训练数据集,通过该训练数据集对预设的机器学习模型进行训练得到威胁识别模型,来对防火墙监测到的流量数据进行识别拦截,提高了威胁流量识别的全面性,进一步保证了网络安全,并且在威胁识别模型识别到威胁流量数据后,根据模型解释器确定识别过程中各个特征维度对应的识别贡献度,从而便于用户可以更有针对性的进行网络防控,提升了用户体验。

[0102] 在本实施例中还提供了一种威胁流量数据的识别装置,该装置用于实现上述实施例及优选实施方式,已经进行过说明的不再赘述。如以下所使用的,术语“模块”可以实现预定功能的软件和/或硬件的组合。尽管以下实施例所描述的装置较佳地以软件来实现,但是硬件,或者软件和硬件的组合的实现也是可能并被构想的。

[0103] 本实施例提供一种威胁流量数据的识别装置,如图3所示,包括:

[0104] 日志数据获取模块401,用于获取防火墙日志,防火墙日志包括:流量日志和安全日志,安全日志包括:防火墙拦截的威胁事件和威胁事件对应的威胁类型,流量日志包括:防火墙监测的流量数据。

[0105] 威胁标签匹配模块402,用于基于安全日志和流量日志进行数据匹配,确定安全日志中的威胁事件对应的目标流量数据,并将威胁事件的威胁类型确定为目标流量数据的威胁标签。

[0106] 威胁模型训练模块403,用于根据安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集,通过训练数据集对预设机器学习模型进行训练,得到威胁识别模型。

[0107] 威胁流量拦截模块404,用于基于威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据。

[0108] 在一些可选的实施方式中,日志数据获取模块401,获取的防火墙日志包括:多个不同防火墙对应的日志,日志数据获取模块401,还用于根据防火墙日志中各个防火墙对应的安全日志,确定各个威胁事件的非结构化文本数据;

[0109] 对各个威胁事件的非结构化文本数据进行数据解析,得到统一结构格式下的各个威胁事件对应的结构化数据,数据解析方式为:基于最长公共子序列的流式日志解析方法。

[0110] 在一些可选的实施方式中,威胁模型训练模块403,在根据安全日志中各个威胁事件对应的目标流量数据和威胁标签构建训练数据集时,包括:根据安全日志中各个威胁事件对应的属性信息,确定目标流量数据的特征维度;

[0111] 提取目标流量数据中各个特征维度的样本特征,基于样本特征和目标流量数据对应威胁事件的威胁标签确定各个目标流量数据对应的样本数据;

[0112] 根据各个目标流量数据对应的样本数据构建训练数据集。

[0113] 在一些可选的实施方式中,威胁模型训练模块403,所训练的预设机器学习模型为:基于轻量级梯度提升机算法构建的机器学习模型。

[0114] 威胁模型训练模块403,在通过训练数据集对预设机器学习模型进行训练时,包括:

[0115] 在通过训练数据集对基于轻量级梯度提升机算法构建的机器学习模型进行训练的过程中,对机器学习模型进行交叉验证;

[0116] 根据交叉验证结果调整机器学习模型的模型参数。

[0117] 在一些可选的实施方式中,威胁流量拦截模块404,在基于所述威胁识别模型对防火墙监测到的流量数据进行威胁识别,拦截威胁流量数据后,还用于:

[0118] 根据模型解释器确定所述威胁识别模型在对所述威胁流量数据进行识别时,各个特征维度对应的贡献值;

[0119] 输出所述威胁流量数据的各个特征维度对应的贡献值。

[0120] 在一些可选的实施方式中,威胁流量拦截模块404,在基于威胁识别模型对防火墙监测到的流量数据进行威胁识别时,包括:

[0121] 通过防火墙对监测到的流量数据进行威胁识别,拦截流量数据中的第一威胁流量数据,排除流量数据中的第一威胁流量数据,得到第二流量数据;

[0122] 通过威胁识别模型对第二流量数据进行威胁识别,拦截第二流量数据中的第二威胁流量数据。

[0123] 在一些可选的实施方式中,威胁流量拦截模块404,还用于输出威胁流量数据对应的威胁事件和威胁类型。

[0124] 上述各个模块和单元的功能描述与上述对应实施例相同,在此不再赘述。

[0125] 本实施例中的威胁流量数据的识别装置是以功能单元的形式来呈现,这里的单元是指ASIC(Application Specific Integrated Circuit,专用集成电路)电路,执行一个或多个软件或固定程序的处理器和存储器,和/或其他可以提供上述功能的器件。

[0126] 本发明实施例还提供一种计算机设备,具有上述图3所示的威胁流量数据的识别装置。

[0127] 请参阅图4,图4是本发明可选实施例提供的一种计算机设备的结构示意图,如图4

所示,该计算机设备包括:一个或多个处理器10、存储器20,以及用于连接各部件的接口,包括高速接口和低速接口。各个部件利用不同的总线互相通信连接,并且可以被安装在公共主板上或者根据需要以其它方式安装。处理器可以对在计算机设备内执行的指令进行处理,包括存储在存储器中或者存储器上以在外部输入/输出装置(诸如,耦合至接口的显示设备)上显示GUI的图形信息的指令。在一些可选的实施方式中,若需要,可以将多个处理器和/或多条总线与多个存储器和多个存储器一起使用。同样,可以连接多个计算机设备,各个设备提供部分必要的操作(例如,作为服务器阵列、一组刀片式服务器、或者多处理器系统)。图4中以一个处理器10为例。

[0128] 处理器10可以是中央处理器,网络处理器或其组合。其中,处理器10还可以进一步包括硬件芯片。上述硬件芯片可以是专用集成电路,可编程逻辑器件或其组合。上述可编程逻辑器件可以是复杂可编程逻辑器件,现场可编程逻辑门阵列,通用阵列逻辑或其任意组合。

[0129] 其中,存储器20存储有可由至少一个处理器10执行的指令,以使至少一个处理器10执行实现上述实施例示出的方法。

[0130] 存储器20可以包括存储程序区和存储数据区,其中,存储程序区可存储操作系统、至少一个功能所需要的应用程序;存储数据区可存储根据计算机设备的使用所创建的数据等。此外,存储器20可以包括高速随机存取存储器,还可以包括非瞬时存储器,例如至少一个磁盘存储器件、闪存器件、或其他非瞬时固态存储器件。在一些可选的实施方式中,存储器20可选包括相对于处理器10远程设置的存储器,这些远程存储器可以通过网络连接至该计算机设备。上述网络的实例包括但不限于互联网、企业内部网、局域网、移动通信网及其组合。

[0131] 存储器20可以包括易失性存储器,例如,随机存取存储器;存储器也可以包括非易失性存储器,例如,快闪存储器,硬盘或固态硬盘;存储器20还可以包括上述种类的存储器的组合。

[0132] 该计算机设备还包括输入装置30和输出装置40。处理器10、存储器20、输入装置30和输出装置40可以通过总线或者其他方式连接,图4中以通过总线连接为例。

[0133] 输入装置30可接收输入的数字或字符信息,以及产生与该计算机设备的用户设置以及功能控制有关的键信号输入,例如触摸屏、小键盘、鼠标、轨迹板、触摸板、指示杆、一个或者多个鼠标按钮、轨迹球、操纵杆等。输出装置40可以包括显示设备、辅助照明装置(例如,LED)和触觉反馈装置(例如,振动电机)等。上述显示设备包括但不限于液晶显示器,发光二极管,显示器和等离子体显示器。在一些可选的实施方式中,显示设备可以是触摸屏。

[0134] 本发明实施例还提供了一种计算机可读存储介质,上述根据本发明实施例的方法可在硬件、固件中实现,或者被实现为可记录在存储介质,或者被实现通过网络下载的原始存储在远程存储介质或非暂时机器可读存储介质中并将被存储在本地存储介质中的计算机代码,从而在此描述的方法可被存储在使用通用计算机、专用处理器或者可编程或专用硬件的存储介质上的这样的软件处理。其中,存储介质可为磁碟、光盘、只读存储记忆体、随机存储记忆体、快闪存储器、硬盘或固态硬盘等;进一步地,存储介质还可以包括上述种类的存储器的组合。可以理解,计算机、处理器、微处理器控制器或可编程硬件包括可存储或接收软件或计算机代码的存储组件,当软件或计算机代码被计算机、处理器或硬件访问且

执行时,实现上述实施例示出的方法。

[0135] 虽然结合附图描述了本发明的实施例,但是本领域技术人员可以在不脱离本发明的精神和范围的情况下做出各种修改和变型,这样的修改和变型均落入由所附权利要求所限定的范围之内。

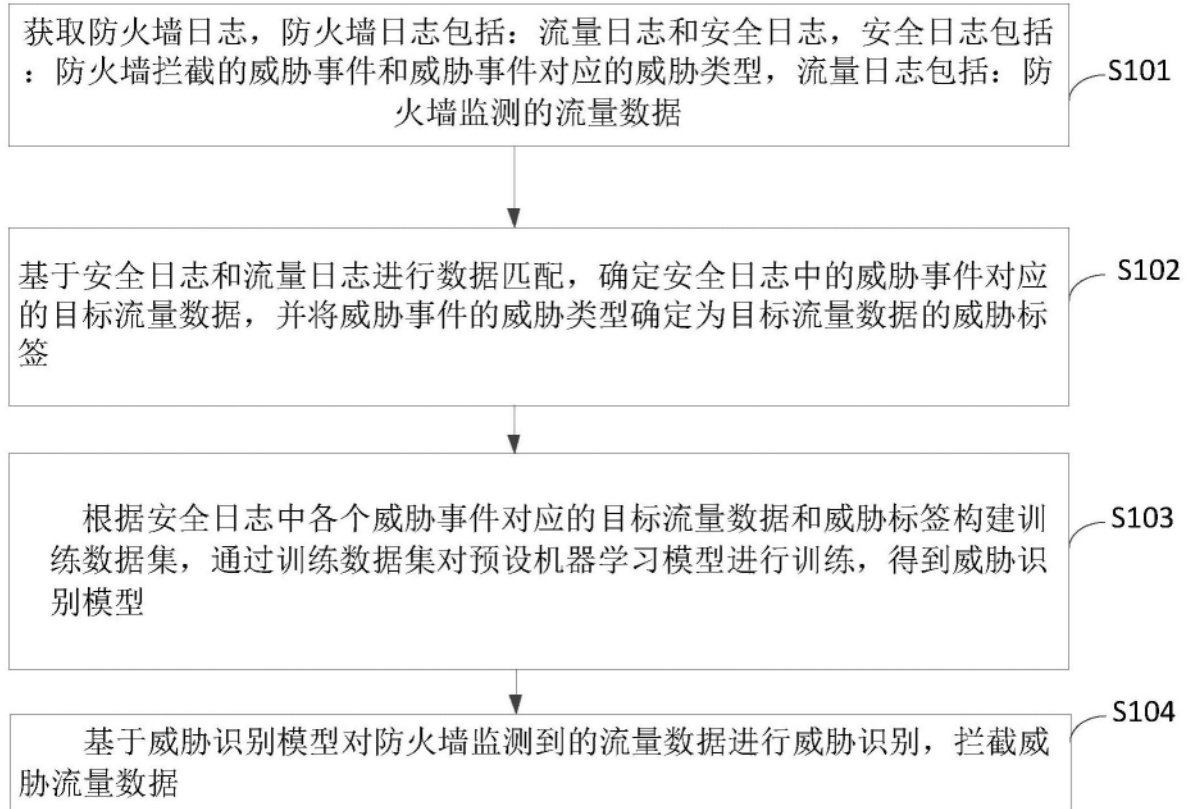


图1

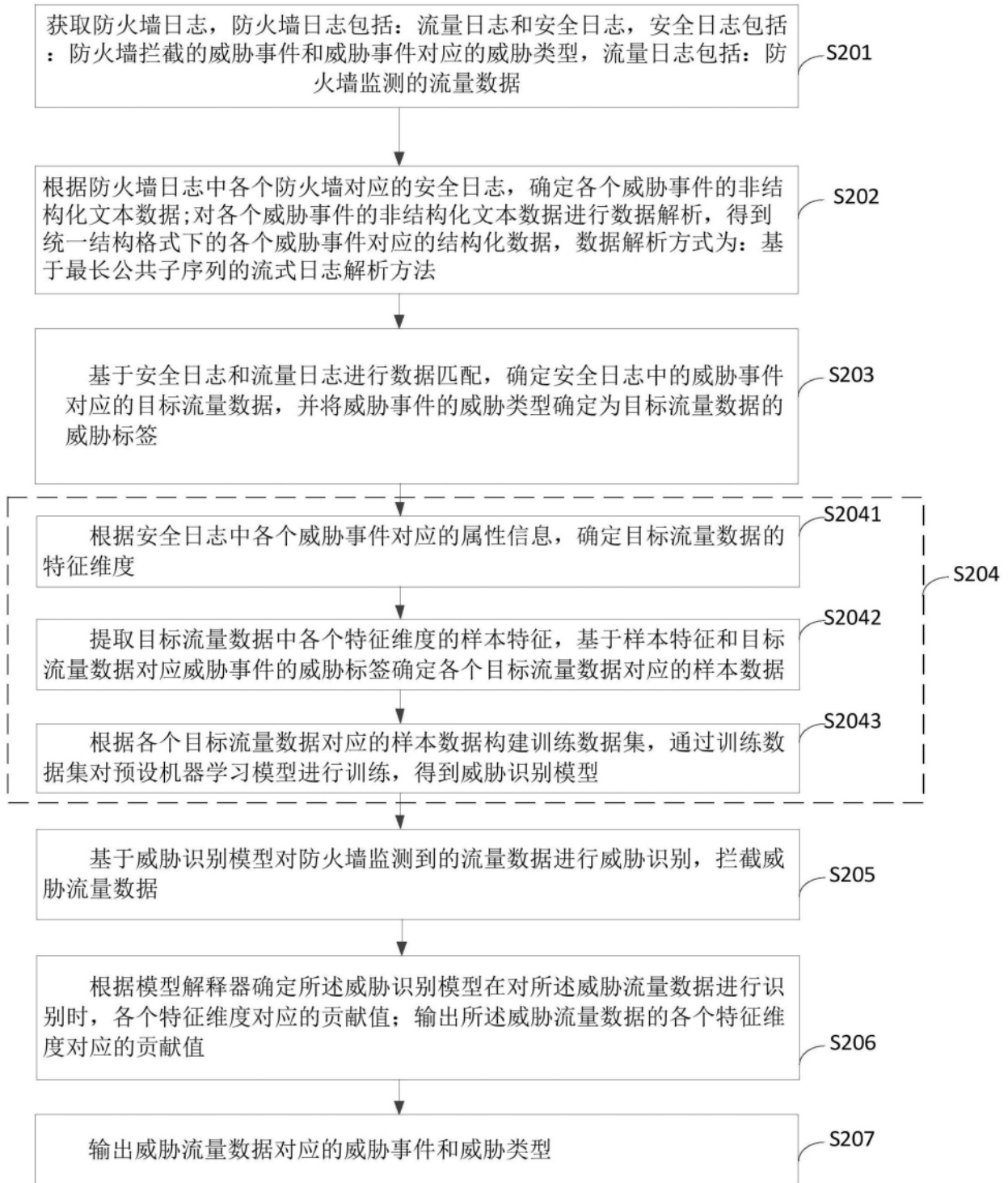


图2



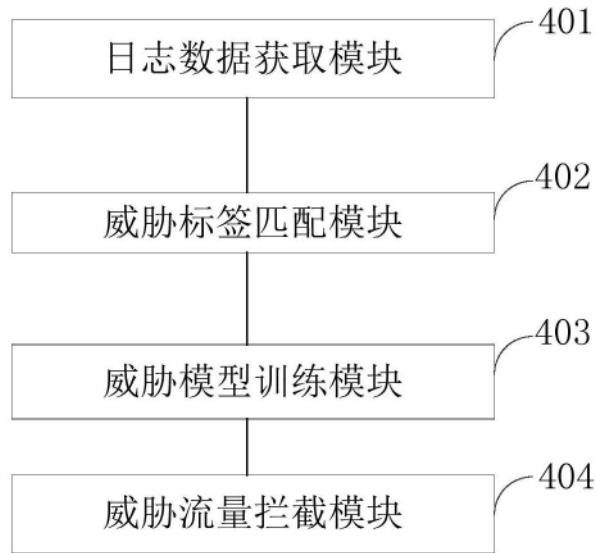


图3

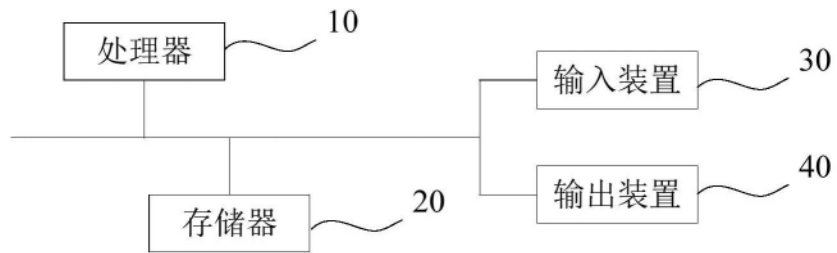


图4