



(12) 发明专利

(10) 授权公告号 CN 111027325 B

(45) 授权公告日 2023. 11. 28

(21) 申请号 201911255072.X

G06F 16/583 (2019.01)

(22) 申请日 2019.12.09

G06F 16/58 (2019.01)

(65) 同一申请的已公布的文献号

G06F 16/55 (2019.01)

申请公布号 CN 111027325 A

G06T 3/40 (2006.01)

(43) 申请公布日 2020.04.17

(56) 对比文件

(73) 专利权人 北京知道创宇信息技术股份有限公司

CN 109460434 A, 2019.03.12

CN 107451153 A, 2017.12.08

地址 100000 北京市朝阳区阜通东大街1号院5号楼1单元311501室

CN 106446782 A, 2017.02.22

CN 106569998 A, 2017.04.19

JP H10198763 A, 1998.07.31

(72) 发明人 胡仁伟 陈效友 张会杰

苑全兵; 黄福. 数字字符识别算法研究. 电子测试. 2010, (04), 全文.

(74) 专利代理机构 北京超凡宏宇知识产权代理有限公司 11463

审查员 胡一冰

专利代理师 刘亚飞

(51) Int. Cl.

G06F 40/295 (2020.01)

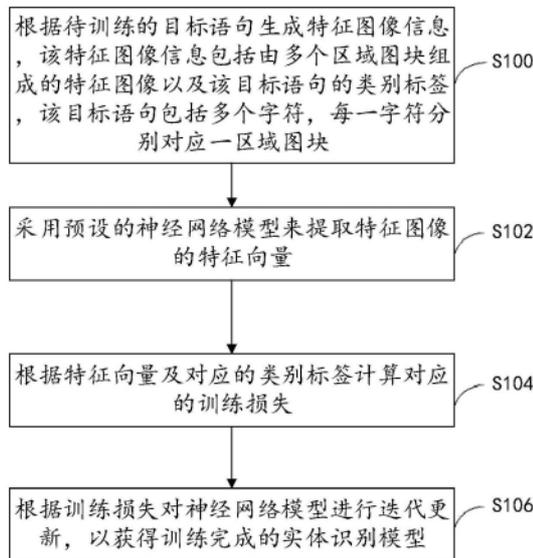
权利要求书2页 说明书10页 附图9页

(54) 发明名称

一种模型生成方法、实体识别方法、装置及电子设备

(57) 摘要

本申请提供一种模型生成方法、实体识别方法、装置及电子设备,该模型生成方法包括:根据待训练的目标语句生成特征图像信息,该特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,该目标语句包括多个字符,每一字符分别对应一区域图块;采用预设的神经网络模型来提取特征图像的特征向量;根据特征向量及对应的类别标签计算对应的训练损失;根据训练损失对神经网络模型进行迭代更新,以获得训练完成的实体识别模型。



1. 一种模型生成方法,其特征在于,所述方法包括:

根据待训练的目标语句生成特征图像信息,所述特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,所述目标语句包括多个字符,每一所述字符分别对应一所述区域图块;

采用预设的神经网络模型来提取所述特征图像的特征向量;

根据所述特征向量及对应的类别标签计算对应的训练损失;

根据所述训练损失对所述神经网络模型进行迭代更新,以获得训练完成的实体识别模型;

所述根据待训练的目标语句生成特征图像信息,包括:

提取所述待训练的目标语句中的每个字符;

根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块预先建立映射关系并存储在数据库中;

根据多个区域图块生成所述特征图像;

所述根据多个区域图块生成所述特征图像,包括:

将查找得到的多个区域图块按照对应的字符在所述待训练的目标语句中的位置依次组合,获得组合图像;

将所述组合图像填充在空白图像的预设区域,并将所述空白图像除所述预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

2. 根据权利要求1所述方法,其特征在于,在所述根据待训练的目标语句生成特征图像信息之前,所述方法还包括:

获取实体数据库中的多个字符以及预设的多个区域图块,其中,所述多个字符中每个字符之间互不重复,预设的多个区域图块中每个区域图块之间互不重复;

建立每个字符与一个预设的区域图块的映射关系并存储在所述数据库中。

3. 根据权利要求1所述方法,其特征在于,所述根据多个区域图块生成所述特征图像,包括:

将查找得到的多个区域图块按照对应的字符在所述待训练的目标语句中的位置依次组合,获得组合图像;

复制多个所述组合图像进行拼接,获得组合拼接图像;

将所述组合拼接图像填充在空白图像的预设区域,并将所述空白图像除所述预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

4. 根据权利要求1所述方法,其特征在于,所述根据多个区域图块生成所述特征图像,包括:

将所述多个区域图块分散填充在空白图像的多个预设区域,所述预设区域的数量与所述区域图块的数量相同;

将所述空白图像除所述多个预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

5. 根据权利要求1所述方法,其特征在于,在所述根据待训练的目标语句生成特征图像信息之后,所述方法还包括:

对所述特征图像进行归一化以及数据增强处理。

6. 一种实体识别方法,其特征在于,所述方法包括:

根据待识别的实体语句生成特征图像,所述特征图像由多个区域图块组成,所述待识别的实体语句包括多个字符,每一所述字符分别对应一所述区域图块;

将所述特征图像输入实体识别模型,所述实体识别模型为权利要求1-5中任一项生成的所述实体识别模型;

获得所述实体识别模型输出的所述待识别的实体语句的预测标签。

7. 一种模型生成装置,其特征在于,所述装置包括:

生成模块,用于根据待训练的目标语句生成特征图像信息,所述特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,所述目标语句包括多个字符,每一所述字符分别对应一所述区域图块;

提取模块,用于采用预设的神经网络模型来提取所述特征图像的特征向量;

计算模块,用于根据所述特征向量及对应的类别标签计算对应的训练损失;

更新模块,用于根据所述训练损失对所述神经网络模型进行迭代更新,以获得训练完成的实体识别模型;

所述生成模块,具体用于提取所述待训练的目标语句中的每个字符;

根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块预先建立映射关系并存储在数据库中;根据多个区域图块生成所述特征图像;所述根据多个区域图块生成所述特征图像,包括:将查找得到的多个区域图块按照对应的字符在所述待训练的目标语句中的位置依次组合,获得组合图像;将所述组合图像填充在空白图像的预设区域,并将所述空白图像除所述预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

8. 一种实体识别装置,其特征在于,所述装置包括:

生成模块,用于根据待识别的实体语句生成特征图像,所述特征图像由多个区域图块组成,所述待识别的实体语句包括多个字符,每一所述字符分别对应一所述区域图块;

输入模块,用于将所述特征图像输入实体识别模型,所述实体识别模型为权利要求1-5中任一项生成的所述实体识别模型;

获得模块,用于获得所述实体识别模型输出的所述待识别的实体语句的预测标签。

9. 一种电子设备,包括存储器和处理器,所述存储器存储有计算机程序,其特征在于,所述处理器执行所述计算机程序时实现权利要求1至6中任一项所述的方法。

一种模型生成方法、实体识别方法、装置及电子设备

技术领域

[0001] 本申请涉及实体识别技术领域,具体而言,涉及一种模型生成方法、实体识别方法、装置及电子设备。

背景技术

[0002] 传统的实体识别方法是通过word2vec将标注语料转化为向量的方式,进而通过神经网络模型对其进行实体识别,但将标注语料转换成向量保存的实体信息较少,进而造成实体识别精度不高的问题。

发明内容

[0003] 本申请实施例的目的在于提供一种模型生成方法、实体识别方法、装置及电子设备,用以解决现有的实体识别方法中将标注语料转换成向量进而通过神经网络模型对其进行实体识别存在的向量保存实体信息较少造成的实体识别精度不高的问题。

[0004] 第一方面,实施例提供一种模型生成方法,所述方法包括:根据待训练的目标语句生成特征图像信息,所述特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,所述目标语句包括多个字符,每一所述字符分别对应一所述区域图块;采用预设的神经网络模型来提取所述特征图像的特征向量;根据所述特征向量及对应的类别标签计算对应的训练损失;根据所述训练损失对所述神经网络模型进行迭代更新,以获得训练完成的实体识别模型。

[0005] 在上述设计的模型生成方法中,通过将待训练的目标语句中的每个字符转换成一区域图块,进而根据多个区域图块生成特征图像,也就是将目标语句转换成了特征图像,通过神经网络模型对该特征图像进行特征提取,进而完成实体识别模型的训练,由于特征提取中图像的方式可以保存更多的实体信息,进而提高了实体识别的精度,解决了现有的实体识别方法中将标注语料转换成向量进而通过神经网络模型对其进行实体识别存在的向量保存实体信息较少造成的实体识别精度不高的问题。

[0006] 在第一方面的可选实施方式中,所述根据待训练的目标语句生成特征图像信息,包括:提取所述待训练的目标语句中的每个字符;根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块预先建立映射关系并存储在数据库中;根据多个区域图块生成所述特征图像。

[0007] 在第一方面的可选实施方式中,在所述根据待训练的目标语句生成特征图像信息之前,所述方法还包括:获取实体数据库中的多个字符以及预设的多个区域图块,其中,所述多个字符中每个字符之间互不重复,预设的多个区域图块中每个区域图块之间互不重复;建立每个字符与一个预设的区域图块的映射关系并存储在所述数据库中。

[0008] 在第一方面的可选实施方式中,所述根据多个区域图块生成所述特征图像,包括:将查找得到的多个区域图块按照对应的字符在所述待训练的目标语句中的位置依次组合,获得组合图像;将所述组合图像填充在空白图像的预设区域,并将所述空白图像除所述预

设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

[0009] 在第一方面的可选实施方式中,所述根据多个区域图块生成所述特征图像,包括:将查找得到的多个区域图块按照对应的字符在所述待训练的目标语句中的位置依次组合,获得组合图像;复制多个所述组合图像进行拼接,获得组合拼接图像;将所述组合拼接图像填充在空白图像的预设区域,并将所述空白图像除所述预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

[0010] 在第一方面的可选实施方式中,所述根据多个区域图块生成所述特征图像,包括:将所述多个区域图块分散填充在空白图像的多个预设区域,所述预设区域的数量与所述区域图块的数量相同;将所述空白图像除所述多个预设区域外的其余区域设置为预设的统一字符,获得所述特征图像。

[0011] 在第一方面的可选实施方式中,在所述根据待训练的目标语句生成特征图像信息之后,所述方法还包括:对所述图像进行归一化以及数据增强处理。

[0012] 在上述设计的实施方式中,通过归一化处理以方便数据处理,加快神经网络学习速度,提高识别的鲁棒性;通过数据增强处理以此来防止深度学习模型过拟合,提高识别的可靠性。

[0013] 第二方面,实施例提供一种实体识别方法,所述方法包括:根据待识别的实体语句生成特征图像,所述特征图像信息包括由多个区域图块组成的特征图像,所述待识别的实体语句包括多个字符,每一所述字符分别对应一所述区域图块;将所述特征图像输入实体识别模型,所述实体识别模型为第一方面中任一可选实施方式生成的所述实体识别模型;获得所述实体识别模型输出的所述待识别的实体语句的预测标签。

[0014] 在上述设计的实体识别方法中,根据该待识别的实体语句生成特征图像,进而通过第一实施例训练完成得到的实体模型对待识别的实体语句生成的特征图像进行预测,进而得到该实体识别模型输出的预测标签,由于该实体识别方法是通过语句转换成特征图像,进而通过第一实施例训练得到的实体模型进行识别,因此,保留的实体识别信息更多,实体识别的精度更高。

[0015] 第三方面,实施例提供一种模型生成装置,所述装置包括:生成模块,用于根据待训练的目标语句生成特征图像信息,所述特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,所述目标语句包括多个字符,每一所述字符分别对应一所述区域图块;提取模块,用于采用预设的神经网络模型来提取所述特征图像的特征向量;计算模块,用于根据所述特征向量及对应的类别标签计算对应的训练损失;更新模块,用于根据所述训练损失对所述神经网络模型进行迭代更新,以获得训练完成的实体识别模型。

[0016] 在上述设计的模型生成装置中,通过将待训练的目标语句中的每个字符转换成一区域图块,进而根据多个区域图块生成特征图像,也就是将目标语句转换成了特征图像,通过神经网络模型对该特征图像进行特征提取,进而完成实体识别模型的训练,由于特征提取中图像的方式可以保存更多的实体信息,进而提高了实体识别的精度,解决了现有的实体识别方法中将标注语料转换成向量进而通过神经网络模型对其进行实体识别存在的向量保存实体信息较少造成的实体识别精度不高的问题。

[0017] 在第三方面的可选实施方式中,所述生成模块具体用于提取所述待训练的目标语句中的每个字符;根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块

预先建立映射关系并存储在数据库中;根据多个区域图块生成所述特征图像。

[0018] 在第三方面的可选实施方式中,所述装置还包括获取模块,用于获取实体数据库中的多个字符以及预设的多个区域图块,其中,所述多个字符中每个字符之间互不重复,预设的多个区域图块中每个区域图块之间互不重复;建立模块,用于建立每个字符与一个预设的区域图块的映射关系并存储在所述数据库中。

[0019] 在第三方面的可选实施方式中,所述装置还包括处理模块,用于对所述特征图像进行归一化以及数据增强处理。

[0020] 第四方面,实施例提供一种实体识别装置,所述装置包括:生成模块,用于根据待识别的实体语句生成特征图像,所述特征图像由多个区域图块组成,所述待识别的实体语句包括多个字符,每一所述字符分别对应一所述区域图块;输入模块,用于将所述特征图像输入实体识别模型,所述实体识别模型为前述实施方式中任一项生成的所述实体识别模型;获得模块,用于获得所述实体识别模型输出的所述待识别的实体语句的预测标签。

[0021] 在上述设计的实体识别装置中,根据该待识别的实体语句生成特征图像,进而通过第一实施例训练完成得到的实体模型对待识别的实体语句生成的特征图像进行预测,进而得到该实体识别模型输出的预测标签,由于该实体识别方法是通过语句转换成特征图像,进而通过第一实施例训练得到的实体模型进行识别,因此,保留的实体识别信息更多,实体识别的精度更高。

[0022] 第五方面,实施例提供一种电子设备,包括存储器和处理器,所述存储器存储有计算机程序,所述处理器执行所述计算机程序时执行第一方面、第一方面的任一可选的实现方式、第二方面、第二方面的任一可选的实现方式中的所述方法。

[0023] 第六方面,实施例提供一种非暂态可读存储介质,其上存储有计算机程序,所述计算机程序被处理器执行时执行第一方面、第一方面的任一可选的实现方式、第二方面、第二方面的任一可选的实现方式中的所述方法。

[0024] 第七方面,实施例提供了一种计算机程序产品,所述计算机程序产品在计算机上运行时,使得计算机执行第一方面、第一方面的任一可选的实现方式、第二方面、第二方面的任一可选的实现方式中的所述方法。

附图说明

[0025] 为了更清楚地说明本申请实施例的技术方案,下面将对本申请实施例中所需要使用的附图作简单地介绍,应当理解,以下附图仅示出了本申请的某些实施例,因此不应被看作是对范围的限定,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相关的附图。

[0026] 图1为本申请第一实施例提供的模型生成方法第一流程图;

[0027] 图2为本申请第一实施例提供的模型生成方法第二流程图;

[0028] 图3为本申请第一实施例提供的模型生成方法第三流程图;

[0029] 图4为本申请第一实施例提供的模型生成方法第四流程图;

[0030] 图5为本申请第一实施例提供的特征图像第一示例图;

[0031] 图6为本申请第一实施例提供的模型生成方法第五流程图;

[0032] 图7为本申请第一实施例提供的特征图像第二示例图;

- [0033] 图8为本申请第一实施例提供的模型生成方法第六流程图；
- [0034] 图9为本申请第一实施例提供的特征图像第三示例图；
- [0035] 图10为本申请第一实施例提供的模型生成方法第五流程图；
- [0036] 图11为本申请第二实施例提供的实体识别方法流程图；
- [0037] 图12为本申请第三实施例提供的模型生成装置结构图；
- [0038] 图13为本申请第四实施例提供的实体识别装置结构图；
- [0039] 图14为本申请第五实施例提供的电子设备结构图。
- [0040] 图标：300-生成模块；302-提取模块；304-计算模块；306-更新模块；308-获取模块；310-建立模块；312-处理模块；400-生成模块；402-输入模块；404-获得模块；5-电子设备；501-处理器；502-存储器；503-通信总线。

具体实施方式

[0041] 下面将结合本申请实施例中的附图，对本申请实施例中的技术方案进行描述。

[0042] 第一实施例

[0043] 如图1所示，本申请实施例提供一种模型生成方法，该方法具体包括如下步骤：

[0044] 步骤S100：根据待训练的目标语句生成特征图像信息，该特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签，该目标语句包括多个字符，每一字符分别对应一区域图块。

[0045] 步骤S102：采用预设的神经网络模型来提取特征图像的特征向量。

[0046] 步骤S104：根据特征向量及对应的类别标签计算对应的训练损失。

[0047] 步骤S106：根据训练损失对神经网络模型进行迭代更新，以获得训练完成的实体识别模型。

[0048] 在步骤S100中，待训练的目标语句即为待训练的实体识别语句，该实体识别语句中包含了多个字符，例如，该待训练的实体识别语句可为一句具有多个中文词的语句，如“ABXY门”、“CB EF庙”等；也可以为一串字符串，例如“139xxxxxxx”、“6125xxxxxxxxxxxxx”等。其中，每一个中文字表示前述所说的字符。可提前对该待训练的目标语句进行类别标签标注，例如，该“ABXY门”语句可标注为地名，“139xxxxxxx”可标注为手机号，“6125xxxxxxxxxxxxx”可标注为身份证号码，该地名、手机号以及身份证号码表示该目标语句的类别标签。在此基础上，步骤S100中的根据待训练的目标语句生成特征图像信息可理解为将该待训练的目标语句中的每一字符都对应有一区域图块，该区域图像可提前配置。由于该目标语句中有多个字符，进而可以得到多个区域图块，根据这多个区域图块生成一个特征图像。在前述的基础上，可以基于多类型的待训练的目标语句来获得多类型的特征图像，对这些图像进行类型标注，并可以将标注好的多个特征图像分为训练集、测试集以及验证集，其中，该训练集、测试集以及验证集的数量比例可自行设定，例如训练集为总数量的60%、测试集为总数量的20%以及验证集为总数量的20%，在进行上述操作之后进而执行步骤S102。

[0049] 步骤S102中的采用预设的神经网络模型来提取特征图像的特征向量可理解为：将步骤S100中生成的特征图像输入预设的神经网络模型，通过该预设的神经网络模型来提取该特征图像对应的特征向量。具体的，该预设的神经网络模型可为ResNet模型，可将该训练

集中的特征图像打乱并分批(如每批为64张特征图像)输入到该ResNet 模型中,经过该神经网络模型的多次卷积、多次池化、多次激活后得到每个特征图像对应的特征向量,进而执行步骤S104。

[0050] 步骤S104中根据特征向量及对应的类别标签计算对应的训练损失可理解为:在步骤S102中得到特征图像对应的特征向量之后,可通过预设的分类函数将提取的该特征图像的特征进行分类,得到初步识别结果,进而将初步识别结果和标注的类别标签对比计算,得到一次训练的损失值,进而执行步骤S106。

[0051] 步骤S106中根据训练损失对神经网络模型进行迭代更新,在步骤S104得到训练的损失值的基础上,会根据反向传播算法进行反向传播,进而对神经网络模型的参数进行更新优化。在更新优化后又进入下一次的训练过程,进而得到第二个训练损失值,循环前述步骤对该神经网络模型的参数进行不断的迭代更新,当得到的训练损失值满足预设值要求或训练达到设定的次数上限之后,根据满足要求时的神经网络模型参数或达到设定次数上限时的参数,得到训练完成的实体识别模型。在得到上述模型之后,可通过前述所说的验证集去验证得到的实体模型是否已经达到要求,并且可以通过测试集去测试训练好的实体识别模型的准确率。

[0052] 在上述设计的模型生成方法中,通过将待训练的目标语句中的每个字符转换成一区域图块,进而根据多个区域图块生成特征图像,也就是将目标语句转换成了特征图像,通过神经网络模型对该特征图像进行特征提取,进而完成实体识别模型的训练,由于特征提取中图像的方式可以保存更多的实体信息,进而提高了实体识别的精度,解决了现有的实体识别方法中将标注语料转换成向量进而通过神经网络模型对其进行实体识别存在的向量保存实体信息较少造成的实体识别精度不高的问题。

[0053] 在本实施例的可选实施方式中,在步骤S100根据待训练的目标语句生成特征图像信息之前,如图2所示,该方法还包括:

[0054] 步骤S90:获取实体数据库中的多个字符以及预设的多个区域图块。

[0055] 步骤S92:建立每个字符与一个预设的区域图块的映射关系并存储在数据库中。

[0056] 在步骤S90中,可预先从实体数据库中获取多个字符,例如,可从字典中获取多个字;该预设的区域图块可为阿拉伯数字或预设的图案等。以阿拉伯数字为例,上述过程具体如下:将获取的每个字符与一阿拉伯数字进行编码,使得每个字符对应一个阿拉伯数字,其中,相同的字符对应的阿拉伯数字相同,不同的字符对应的阿拉伯数字不同。例如,前述举例中的“ABXY门”中的字符“A”、“B”、“X”、“Y”、“门”可分别对应阿拉伯数字“1”、“3”、“5”、“7”、“9”;“CBEF庙”中的字符“C”、“B”、“E”、“F”、“庙”可分别对应“11”、“3”、“4”、“12”、“13”。另外,该预设的图案可为图形符号或希腊数字等。以前述的方式建立字符与对应的区域图块的映射关系,并建立映射关系之后存储在数据库中。

[0057] 在本实施例的可选实施方式中,步骤S100中的根据待训练的目标语句生成特征图像信息,如图3所示,具体可为:

[0058] 步骤S1000:提取该待训练的目标语句中的每个字符。

[0059] 步骤S1002:根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块预先建立映射关系并存储在数据库中。

[0060] 步骤S1004:根据多个区域图块生成特征图像。

[0061] 在前述步骤S100中已经提到该待训练的目标语句中包含了多个字符,步骤S1000中可理解为提取该待训练的目标语句中的每个字符,例如,当待训练的目标语句为“ABXY门”时,步骤S1000可为提取该目标语句中的字符“A”、“B”、“X”、“Y”、“门”,进而执行步骤S1002。

[0062] 步骤S1002可理解为:根据前述步骤S90~S92建立的映射关系,根据提取的字符“A”、“B”、“X”、“Y”、“门”可分别在数据库中查找到对应阿拉伯数字“1”、“3”、“5”、“7”、“9”,进而执行步骤S1004,根据查找到的阿拉伯数字“1”、“3”、“5”、“7”、“9”生成特征图像。

[0063] 在本实施例的可选实施方式中,步骤S1004根据多个区域图块生成特征图像,可通过word2Image模型生成对应的特征图像,如图4所示,可具体为:

[0064] 步骤S10040:将查找得到的多个区域图块按照对应的字符在待训练的目标语句中的位置依次组合,获得组合图像。

[0065] 步骤S10042:将组合图像填充在空白图像的预设区域,并将空白图像除预设区域外的其余区域设置为预设的统一字符,获得特征图像。

[0066] 在前述步骤S1002的基础上,执行步骤S10040,在前述举例的描述上,“A”、“B”、“X”、“Y”、“门”在数据库中查找到对应阿拉伯数字“1”、“3”、“5”、“7”、“9”之后,将该“1”、“3”、“5”、“7”、“9”按照顺序依次组合,形成“1、3、5、7、9”的组合图像,进而执行步骤S10042。

[0067] 在步骤S10042中将组合图像填充在空白图像的预设区域,可理解为将上述图像填充在空白图像的预设区域内,其中,该空白图像的大小可提前设置,例如,该空白图像可为64*64的空白图像或其他大小的空白图像;该预设区域可提前设置,例如,可将该组合图像填充在该空白图像的中央,进而将该空白图像的其余区域设置为预设的统一字符。例如,如图5所示,将前述所说的“1、3、5、7、9”填充在10*10大小的空白图像的3*5、4*5、5*5、6*5以及7*5位置之后,该“1、3、5、7、9”为5个位置区域,除了该中央预设区域内呈现“1、3、5、7、9”的五个区域图块的组合图像之外,该空白图像的其余区域均设置为字符0,进而获得特征图像。

[0068] 这里需要说明的是,前述所说的预设区域可为该空白图像的任意区域,不仅限于图像中央。

[0069] 另外,考虑到通常情况下目标语句的长度较短,因此,在按照顺序依次组合得到组合图像之后,可以将多个相同的该组合图像拼接设置在该空白图像的预设区域,以此来增加特征图像的特征信息,具体如图6所示,包括以下步骤:

[0070] 步骤S10044:将查找得到的多个区域图块按照对应的字符在待训练的目标语句中的位置依次组合,获得组合图像。

[0071] 步骤S10045:复制多个组合图像进行拼接,获得组合拼接图像。

[0072] 步骤S10046:将组合拼接图像填充在空白图像的预设区域,并将空白图像除预设区域外的其余区域设置为预设的统一字符,获得特征图像。

[0073] 上述步骤具体可以如图7的示例所示,在形成“1、3、5、7、9”的组合图像之后,可将该多个(图中为6个)形成“1、3、5、7、9”的组合图像拼接设置在该10*10空白图像的中央,以此来增加特征图像的特征信息。

[0074] 在本实施例的可选实施方式中,除了按照步骤S10040的方式将字符对应的区域图块按照顺序组合获得组合图像以外,还可以直接随机将字符对应的区域图块分散分布在该

空白图像内。具体地,如图8所示,包括如下步骤:

[0075] 步骤S10047:将多个区域图块分散填充在空白图像的多个预设区域,预设区域的数量与区域图块的数量相同。

[0076] 步骤S10048:将空白图像除多个预设区域外的其余区域设置为预设的统一字符,获得特征图像。

[0077] 上述步骤具体如图9的示例所示,以10*10的空白图像为例,将该“A”对应的“1”设置在1*1位置,将该“B”对应的“3”设置在10*1位置;将该“X”对应的“5”设置在5*5位置;将该“Y”对应的“7”设置在1*10位置;将该“门”对应的“9”设置在10*10位置。

[0078] 在本实施例的可选实施方式中,在步骤S100根据待训练的目标语句生成特征图像信息之后,如图10所示,该方法还包括:

[0079] 步骤S101:对该特征图像进行归一化以及数据增强处理。

[0080] 在上述步骤中,对该特征图像进行归一化处理可理解为:因为图像的像素值区间在[0,255],而实体数据库中的实体数量远超过255这个数值,按照前述的方式,生成的特征图像的像素值很可能会超过255,所以对转换后的特征图像进行归一化处理,例如,将超过255的范围转换成一个预设的范围或将该特征图像所有像素的灰度值都映射到一个预设范围(如[0,1]这个区间),确保图像中的每个像素点对训练结果的贡献是相同的,以方便数据处理,加快神经网络学习速度,提高识别的鲁棒性。

[0081] 对该特征图像进行数据增强处理表示为对该生成的特征图像进行加噪声/翻转/增强等操作处理,以此来防止深度学习模型过拟合,提高识别的可靠性。

[0082] 第二实施例

[0083] 本申请提供一种实体识别方法,如图11所示,该方法具体包括如下步骤:

[0084] 步骤S200:根据待识别的实体语句生成特征图像,该特征图像信息包括由多个区域图块组成的特征图像,该待识别的实体语句包括多个字符,每一字符分别对应一区域图块。

[0085] 步骤S202:将该特征图像输入实体识别模型,该实体识别模型为第一实施例中任一可选实施方式生成的实体识别模型。

[0086] 步骤S204:获得实体识别模型输出的待识别的实体语句的预测标签。

[0087] 上述步骤中的步骤S200根据待识别模型的实体语句生成特征图像的方式与第一实施例中的步骤S100的方式一致,在这里不再赘述。

[0088] 在步骤S200生成特征图像之后,执行步骤S202将该特征图像输入该实体识别模型,该实体识别模型为第一实施例得到的训练完成的实体识别模型。在将该待识别的实体对应的特征图像输入该训练完成的实体识别模型之后,该实体识别模型会输出该特征图像对应的预测标签,也就是该待识别的实体语句的预测标签。例如,该待识别的实体语句为“UVW滩”,根据该“UVW滩”转换成对应的特征图像,进而将该特征图像输入第一实施例训练得到的实体识别模型中,该实体识别模型会输出预测的标签可能为地名。

[0089] 在上述设计的实体识别方法中,根据该待识别的实体语句生成特征图像,进而通过第一实施例训练完成得到的实体模型对待识别的实体语句生成的特征图像进行预测,进而得到该实体识别模型输出的预测标签,由于该实体识别方法是通过语句转换成特征图像,进而通过第一实施例训练得到的实体模型进行识别,因此,保留的实体识别信息更多,

实体识别的精度更高。

[0090] 第三实施例

[0091] 图12出示了本申请提供的模型生成装置的示意性结构框图,应理解,该装置与上述图1至图10中的方法实施例对应,能够执行第一实施例中服务器执行的方法涉及的步骤,该装置具体的功能可以参见上文中的描述,为避免重复,此处适当省略详细描述。该装置包括至少一个能以软件或固件(firmware)的形式存储于存储器中或固化在装置的操作系统(operating system,OS)中的软件功能模块。具体地,该装置包括:生成模块300,用于根据待训练的目标语句生成特征图像信息,该特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,该目标语句包括多个字符,每一字符分别对应一区域图块;提取模块302,用于采用预设的神经网络模型来提取特征图像的特征向量;计算模块304,用于根据特征向量及对应的类别标签计算对应的训练损失;更新模块306,用于根据训练损失对神经网络模型进行迭代更新,以获得训练完成的实体识别模型。

[0092] 在上述设计的模型生成装置中,通过将待训练的目标语句中的每个字符转换成一区域图块,进而根据多个区域图块生成特征图像,也就是将目标语句转换成了特征图像,通过神经网络模型对该特征图像进行特征提取,进而完成实体识别模型的训练,由于特征提取中图像的方式可以保存更多的实体信息,进而提高了实体识别的精度,解决了现有的实体识别方法中将标注语料转换成向量进而通过神经网络模型对其进行实体识别存在的向量保存实体信息较少造成的实体识别精度不高的问题。

[0093] 在本实施例的可选实施方式中,生成模块300具体用于提取所述待训练的目标语句中的每个字符;根据提取的每个字符查找对应的区域图块,每个字符与对应的区域图块预先建立映射关系并存储在数据库中;根据多个区域图块生成所述特征图像。

[0094] 在本实施例的可选实施方式中,该装置还包括获取模块308,用于获取实体数据库中的多个字符以及预设的多个区域图块,其中,多个字符中每个字符之间互不重复,预设的多个区域图块中每个区域图块之间互不重复;建立模块310,用于建立每个字符与一个预设的区域图块的映射关系并存储在数据库中。

[0095] 在本实施例的可选实施方式中,该装置还包括处理模块312,用于对特征图像进行归一化以及数据增强处理。

[0096] 第四实施例

[0097] 图13出示了本申请提供的实体识别装置的示意性结构框图,应理解,该装置与上述图11中的方法实施例对应,能够执行第一实施例中服务器执行的方法涉及的步骤,该装置具体的功能可以参见上文中的描述,为避免重复,此处适当省略详细描述。该装置包括至少一个能以软件或固件(firmware)的形式存储于存储器中或固化在装置的操作系统(operating system,OS)中的软件功能模块。具体地,该装置包括:生成模块400,用于根据待识别的实体语句生成特征图像,该特征图像由多个区域图块组成,待识别的实体语句包括多个字符,每一字符分别对应一区域图块;输入模块402,用于将特征图像输入实体识别模型,该实体识别模型为第一实施例中任一实施方式生成的实体识别模型;获得模块404,用于获得实体识别模型输出的待识别的实体语句的预测标签。

[0098] 在上述设计的实体识别装置中,根据该待识别的实体语句生成特征图像,进而通过第一实施例训练完成得到的实体模型对待识别的实体语句生成的特征图像进行预测,进

而得到该实体识别模型输出的预测标签,由于该实体识别方法是通过语句转换成特征图像,进而通过第一实施例训练得到的实体模型进行识别,因此,保留的实体识别信息更多,实体识别的精度更高。

[0099] 第五实施例

[0100] 如图14示,本申请提供一种电子设备5,包括:处理器501和存储器502,处理器501和存储器502通过通信总线503和/或其他形式的连接机构(未标出)互连并相互通讯,存储器502存储有处理器501可执行的计算机程序,当计算设备运行时,处理器501执行该计算机程序,以执行时执行第一实施例、第一实施例的任一可选的实现方式、第二实施例、第二实施例的任一可选的实现方式中的方法,例如步骤S100~步骤S106:根据待训练的目标语句生成特征图像信息,该特征图像信息包括由多个区域图块组成的特征图像以及该目标语句的类别标签,该目标语句包括多个字符,每一字符分别对应一区域图块;采用预设的神经网络模型来提取所述特征图像的特征向量;根据特征向量及对应的类别标签计算对应的训练损失;根据训练损失对神经网络模型进行迭代更新,以获得训练完成的实体识别模型。

[0101] 本申请提供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程序,该计算机程序被处理器运行时执行第一实施例、第一实施例的任一可选的实现方式、第二实施例、第二实施例的任一可选的实现方式中的方法。

[0102] 其中,存储介质可以由任何类型的易失性或非易失性存储设备或者它们的组合实现,如静态随机存取存储器(Static Random Access Memory,简称SRAM),电可擦除可编程只读存储器(Electrically Erasable Programmable Read-Only Memory,简称EEPROM),可擦除可编程只读存储器(Erasable Programmable Read Only Memory,简称 EPROM),可编程只读存储器(Programmable Red-Only Memory,简称 PROM),只读存储器(Read-Only Memory,简称ROM),磁存储器,快闪存储器,磁盘或光盘。

[0103] 本申请提供一种计算机程序产品,该计算机程序产品在计算机上运行时,使得计算机执行第一实施例、第一实施例的任一可选的实现方式、第二实施例、第二实施例的任一可选的实现方式中的所述方法。

[0104] 在本申请所提供的实施例中,应该理解到,所揭露装置和方法,可以通过其它的方式实现。以上所描述的装置实施例仅仅是示意性的,例如,所述单元的划分,仅仅为一种逻辑功能划分,实际实现时可以有另外的划分方式,又例如,多个单元或组件可以结合或者可以集成到另一个系统,或一些特征可以忽略,或不执行。另一点,所显示或讨论的相互之间的耦合或直接耦合或通信连接可以是通过一些通信接口,装置或单元的间接耦合或通信连接,可以是电性,机械或其它的形式。

[0105] 另外,作为分离部件说明的单元可以是或者也可以不是物理上分开的,作为单元显示的部件可以是或者也可以不是物理单元,即可以位于一个地方,或者也可以分布到多个网络单元上。可以根据实际的需要选择其中的部分或者全部单元来实现本实施例方案的目的。

[0106] 需要说明的是,功能如果以软件功能模块的形式实现并作为独立的产品销售或使用,可以存储在一个计算机可读取存储介质中。基于这样的理解,本申请的技术方案本质上或者说对现有技术做出贡献的部分或者该技术方案的部分可以以软件产品的形式体现出来,该计算机软件产品存储在一个存储介质中,包括若干指令用以使得一台计算机设备

(可以是个人计算机,服务器,或者网络设备等)执行本申请各个实施例所述方法的全部或部分步骤。而前述的存储介质包括:U盘、移动硬盘、只读存储器(Read-Only Memory,ROM)随机存取存储器(Random Access Memory,RAM)、磁碟或者光盘等各种可以存储程序代码的介质。

[0107] 在本文中,诸如第一和第二等之类的关系术语仅仅用来将一个实体或者操作与另一个实体或操作区分开来,而不一定要求或者暗示这些实体或操作之间存在任何这种实际的关系或者顺序。

[0108] 以上所述仅为本申请的实施例而已,并不用于限制本申请的保护范围,对于本领域的技术人员来说,本申请可以有各种更改和变化。凡在本申请的精神和原则之内,所作的任何修改、等同替换、改进等,均应包含在本申请的保护范围之内。

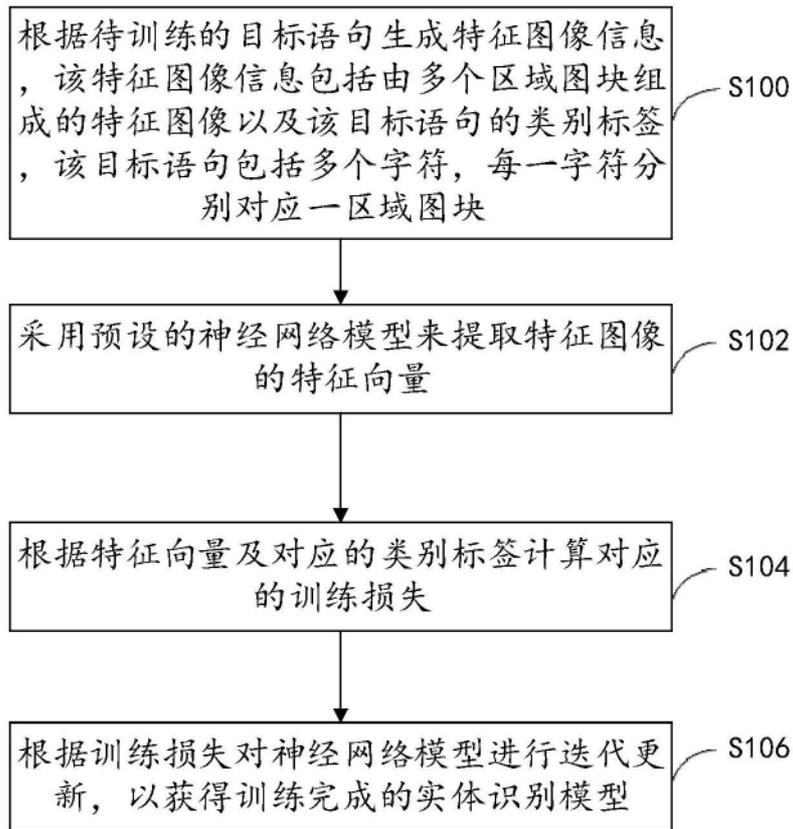


图1

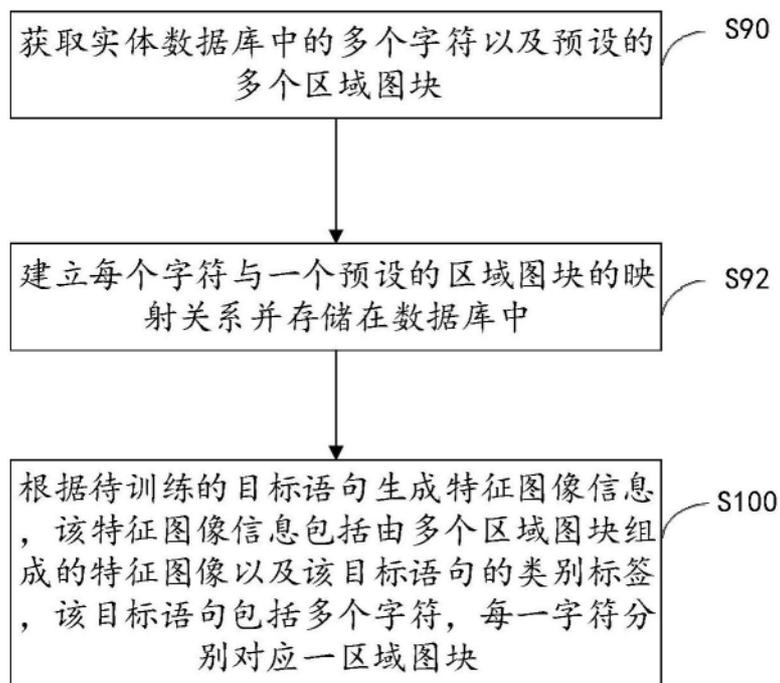


图2

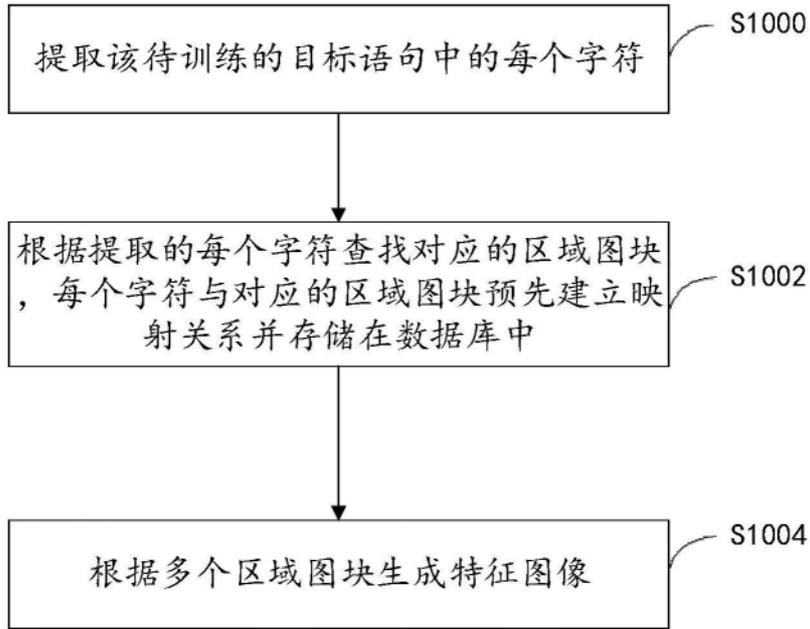


图3

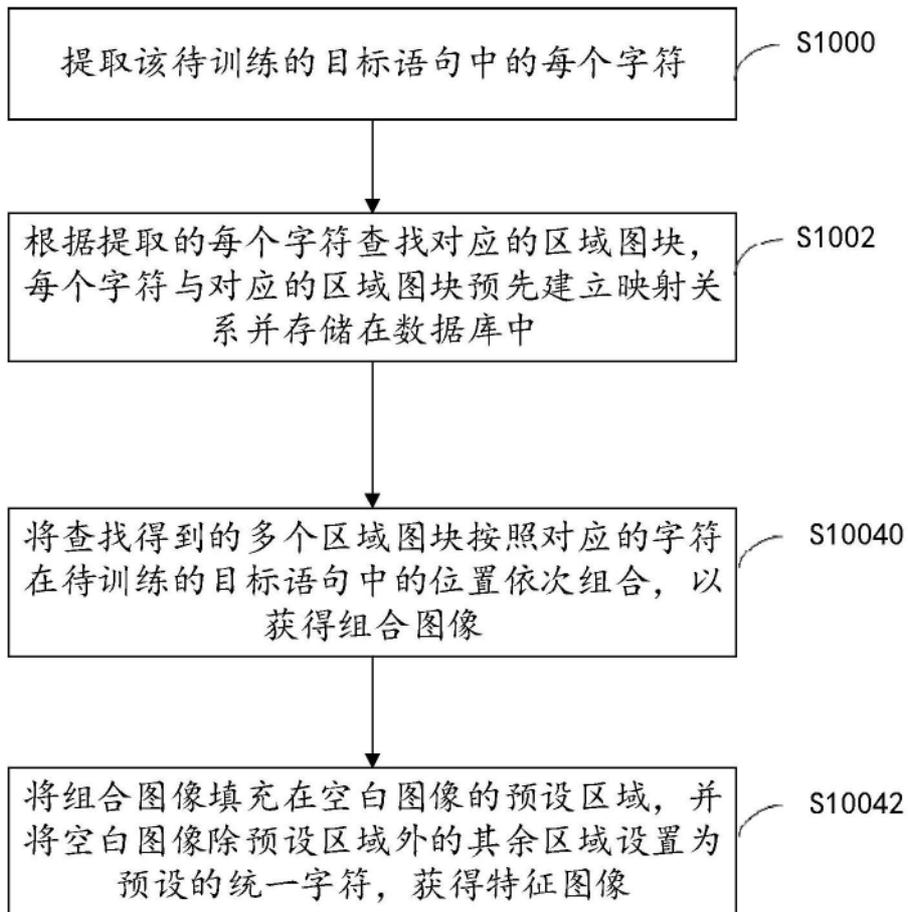


图4

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

图5

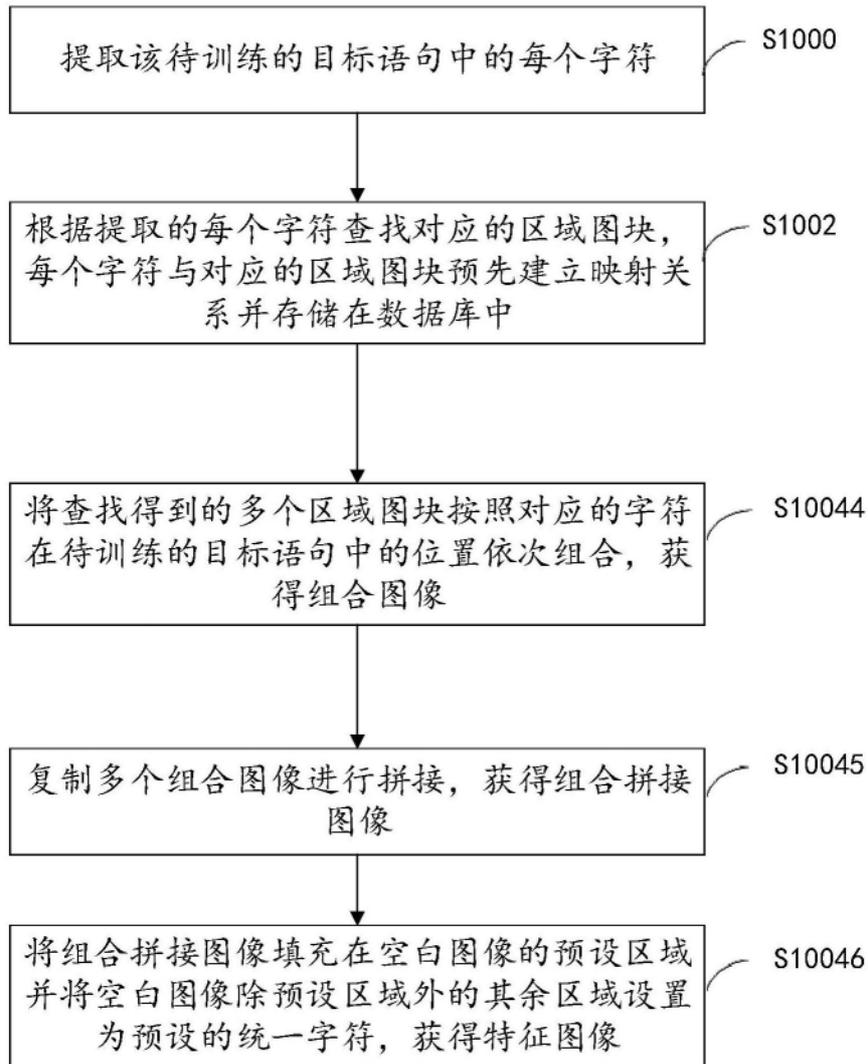


图6

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	1	3	5	7	9	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

图7

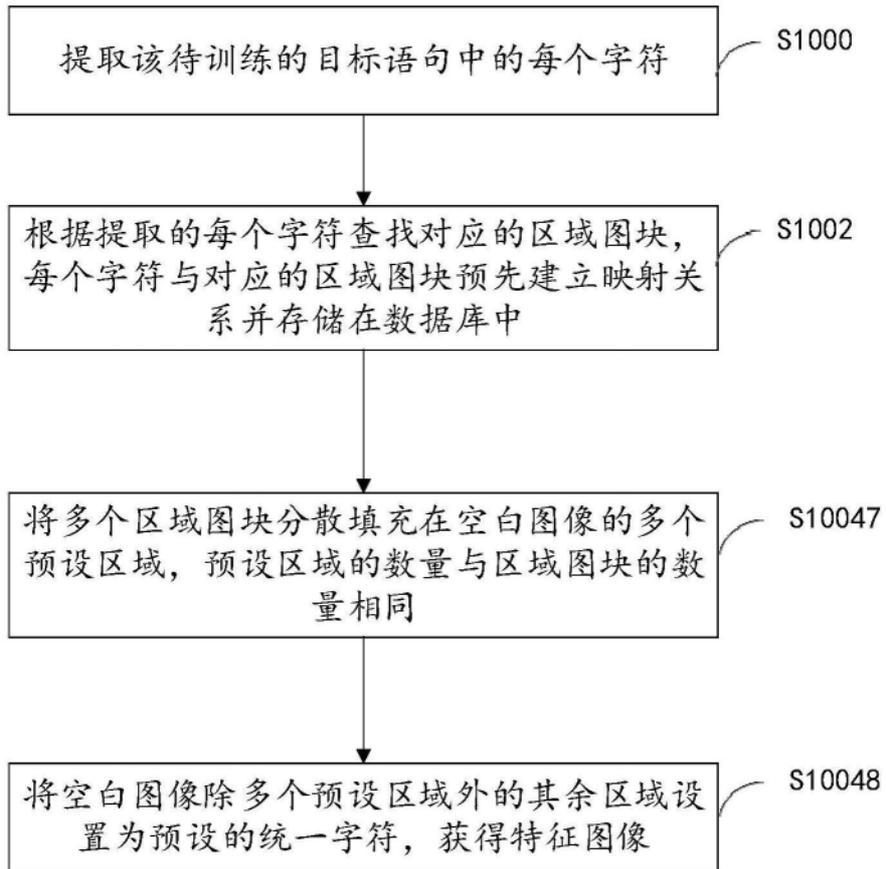


图8

1	0	0	0	0	0	0	0	0	0	3
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	5	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	9

图9

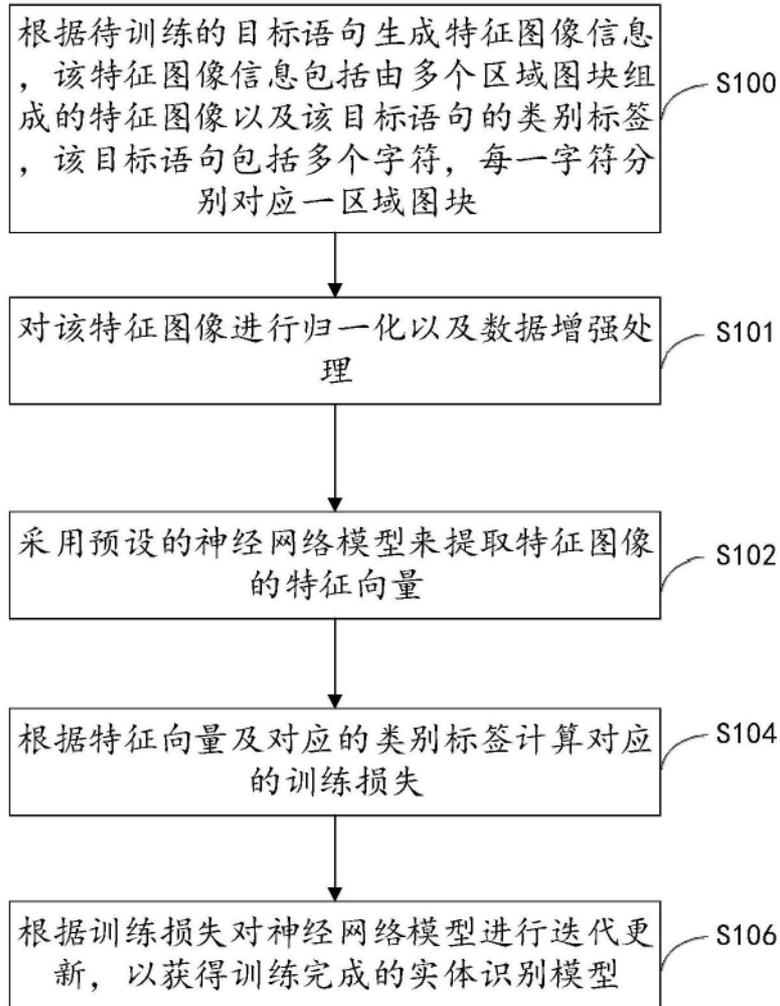


图10

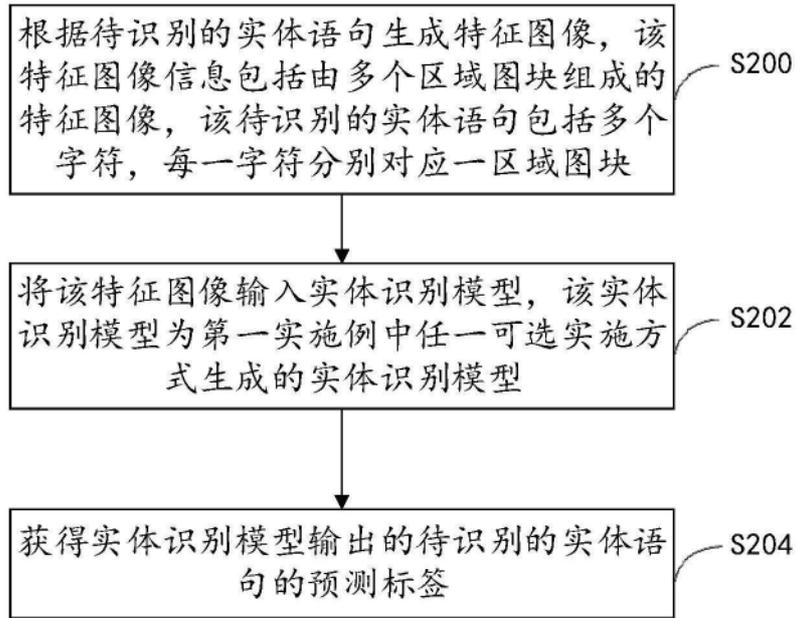


图11

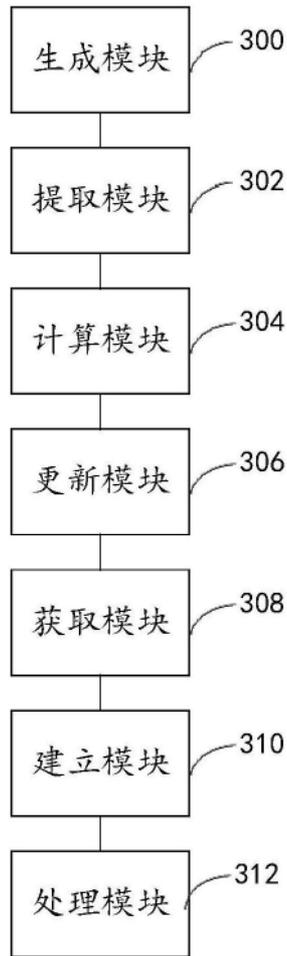


图12

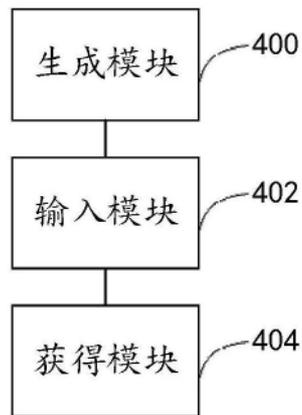


图13

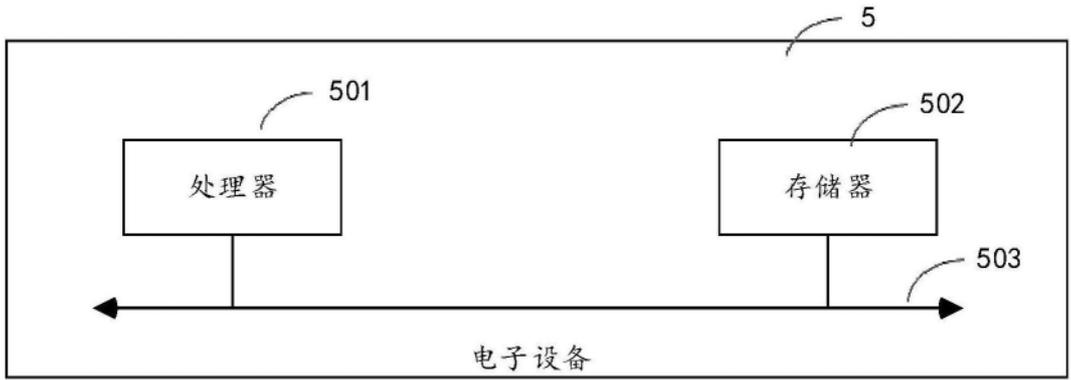


图14