



(12) 发明专利申请

(10) 申请公布号 CN 114417841 A

(43) 申请公布日 2022. 04. 29

(21) 申请号 202111622786.7

(22) 申请日 2021.12.28

(71) 申请人 航天科工网络信息发展有限公司
地址 430040 湖北省武汉市临空港经济技术
开发区五环大道666号(21)

(72) 发明人 李慧琦 牛慧博 刘铁生

(74) 专利代理机构 中国航天科工集团公司专利
中心 11024

代理人 张国虹

(51) Int. Cl.

G06F 40/279 (2020.01)

G06F 16/35 (2019.01)

G06N 3/04 (2006.01)

G06N 3/08 (2006.01)

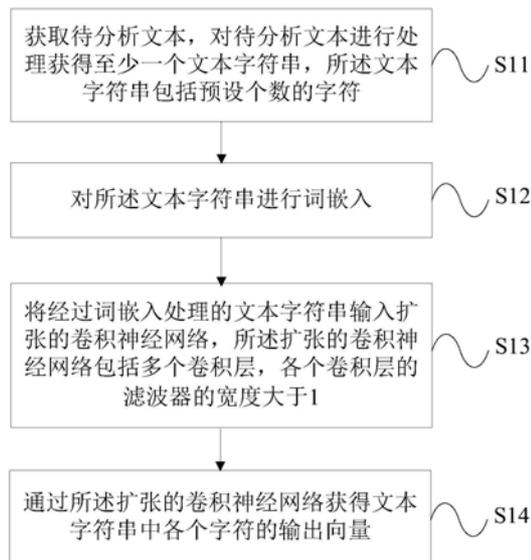
权利要求书1页 说明书7页 附图3页

(54) 发明名称

一种基于扩张卷积神经网络的特征提取方法及装置

(57) 摘要

本发明公开了一种基于扩张卷积神经网络的特征提取方法及装置。该方法包括：获取待分析文本，对待分析文本进行处理获得至少一个文本字符串，所述文本字符串包括预设个数的字符；对所述文本字符串进行词嵌入；将经过词嵌入处理的文本字符串输入扩张的卷积神经网络，所述扩张的卷积神经网络包括多个卷积层，各个卷积层的滤波器的宽度大于1；通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。可见，本发明通过扩张的卷积神经网络对待分析文本进行特征提取，有效输入宽度可以随深度呈指数增长，提高了特征提取的效率。



1. 一种基于扩张卷积神经网络的特征提取方法,其特征在于,包括:
获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;
对所述文本字符串进行词嵌入;
将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;
通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。
2. 根据权利要求1所述的方法,其特征在于,所述通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量,包括:
将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。
3. 根据权利要求1所述的方法,其特征在于,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。
4. 根据权利要求1所述的方法,其特征在于,所述方法还包括:
构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。
5. 一种基于扩张卷积神经网络的特征提取装置,其特征在于,包括:
文本字符串获取单元,用于获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;
词嵌入处理单元,用于对所述文本字符串进行词嵌入;
文本字符串输入单元,用于将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;
输入向量获取单元,用于通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。
6. 根据权利要求5所述的装置,其特征在于,所述输入向量获取单元进一步用于:
将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。
7. 根据权利要求5所述的装置,其特征在于,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。
8. 根据权利要求5所述的装置,其特征在于,还包括:并行处理单元,用于构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。
9. 一种电子设备,其特征在于,该电子设备包括:
处理器;以及,
被安排成存储计算机可执行指令的存储器,所述可执行指令在被执行时使所述处理器执行根据权利要求1-4中任一项所述的方法。
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存储一个或多个程序,所述一个或多个程序当被处理器执行时,实现权利要求1-4中任一项所述的方法。

一种基于扩张卷积神经网络的特征提取方法及装置

技术领域

[0001] 本发明涉及自然语言处理领域,特别涉及一种基于扩张卷积神经网络的特征提取方法、装置、电子设备、计算机可读存储介质。

背景技术

[0002] 信息抽取是自然语言处理领域很重要的任务之一,从文本中抽取特定的结构信息,并形成结构化数据的技术,是形成知识图谱和知识库的基础。实体关系联合抽取是信息抽取任务的重要分支,需要从非结构化文本中抽取出形如(头实体,关系,尾实体)的结构化数据,可以看做同样是信息抽取的子任务命名实体识别(Named Entity Recognition,简称NER)和关系抽取的综合任务。

[0003] 随着信息化时代的发展,越来越多的文本充斥着各种平台,有效的分析这些数据可以获得之中隐藏的信息,促进各行各业的发展,也就是自然语言处理任务的初衷和目的。然而大多数的文本数据是非结构化的,这样的数据处理困难,提取信息量稀少,是不理想的数据源,而结构化的数据处理起来方便简单有条理,且能最大程度地挖掘文本中的信息。因此将非结构化数据转换为结构化数据是一个重要课题,一般将从非结构化文本中抽取出特定的三元组等结构化数据的过程叫做信息提取。

[0004] 命名实体识别(Named Entity Recognition,NER)旨在从文本中抽取出预先定义好类型的命名实体,比如人名、地名、机构名等,是自然语言处理领域信息抽取任务中很重要的子任务,可以帮助后续的文本语义理解。区别于简单的分词,拥有着丰富类型的命名实体是最基本的语义单元,也是绝大多数自然语言处理任务的基础,因此NER任务在研究和应用领域都有很大的作用和价值。

[0005] 命名实体识别的实体分类有很多种,较为通用的有人名(Person)、地名(Location)、机构名(Organization)、日期(Date)、时间(Time)等,根据任务需求和应用场景不同也会有更加细粒度的分类,例如人可以细分成工人、律师、政府人员等,地名包括国家、州或城市,还有特定领域的实体名称,或者各个公司机构的特定专有名词,如品牌名称、产品名称等。命名实体识别任务需要从给定的非结构化文本中,抽取出所有可能是预定义的实体类型的文本所在。

[0006] NER问题也可看作是序列标注问题的一种。序列标注问题可以认为是分类问题的一个推广,或者是更复杂的结构预测(structure prediction)问题的简单形式。序列标注问题的输入是一个观测序列,输出是一个标记序列或状态序列。问题的目标在于学习一个模型,使它能够对观测序列给出标记序列作为预测。

[0007] 解决序列标注问题的方式大体可分为两种,一种是概率图模型,另一种是深度学习模型。

[0008] 概率图模型方法主要有HMM(隐马尔可夫模型,Hidden Markov Model)、MEMM(最大熵马尔可夫模型,Maximum Entropy Markov Model)、CRF(条件随机场,Conditional Random Field)三种方法。将HMM应用至序列标注时,状态对应着标记,此时假定待标注数据

是由HMM生成的,这样可以利用HMM的学习与预测算法进行标注:学习阶段采用极大似然估计,预测阶段采用Viterbi算法。MEMM是对HMM和ME的结合,取消了观测独立性假设,且是判别式模型。HMM和MEMM存在标注偏置(label bias)问题。CRF也是判别式模型,它和MEMM的区别在于将局部归一化变成了全局归一化(将标记序列作为整体看待),因此进一步解决了标注偏置问题。相比于HMM、MEMM,CRF需要训练的参数更多,存在训练代价大、复杂度高的缺点。

[0009] 深度学习模型更多是改变模型结构,通常使用的结构有BiLSTM、transformer等网络结构,也会使用CRF对结果进行整理。但是BiLSTM、transformer这类的结构很难并行化处理,需要的运行时间通常会很长,存在序列标注效率低的问题。

发明内容

[0010] 鉴于上述问题,提出了本发明以便提供一种克服上述问题或者至少部分地解决上述问题的基于扩张卷积神经网络的特征提取方法、装置、电子设备、计算机可读存储介质。

[0011] 本发明的一个实施例提供一种基于扩张卷积神经网络的特征提取方法,该方法包括:

[0012] 获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;

[0013] 对所述文本字符串进行词嵌入;

[0014] 将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;

[0015] 通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。

[0016] 可选地,所述通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量,包括:

[0017] 将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。

[0018] 可选地,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。

[0019] 可选地,所述方法还包括:

[0020] 构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。

[0021] 本发明的另一个实施例提供一种基于扩张卷积神经网络的特征提取装置,包括:

[0022] 文本字符串获取单元,用于获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;

[0023] 词嵌入处理单元,用于对所述文本字符串进行词嵌入;

[0024] 文本字符串输入单元,用于将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;

[0025] 输入向量获取单元,用于通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。

[0026] 可选地,所述输入向量获取单元进一步用于:

[0027] 将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。

[0028] 可选地,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。

[0029] 可选地,还包括:并行处理单元,用于构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。

[0030] 本发明的另一个实施例提供一种电子设备,其中,该电子设备包括:

[0031] 处理器;以及,

[0032] 被安排成存储计算机可执行指令的存储器,所述可执行指令在被执行时使所述处理器执行上述的基于扩张卷积神经网络的特征提取方法。

[0033] 本发明的另一个实施例提供一种计算机可读存储介质,其中,所述计算机可读存储介质存储一个或多个程序,所述一个或多个程序当被处理器执行时,实现上述的基于扩张卷积神经网络的特征提取方法。

[0034] 本发明的有益效果是,本发明通过扩张的卷积神经网络对待分析文本进行特征提取,有效输入宽度可以随深度呈指数增长,提高了特征提取的效率。

[0035] 本发明迭代扩张的卷积神经网络,相比传统的卷积神经网络可更好地用于大上下文和结构化预测,在提高效率的同时保持了类似的F1性能。

[0036] 本发明还可以构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络,进一步提高了特征提取的效率。

附图说明

[0037] 图1为本发明一个实施例的基于扩张卷积神经网络的特征提取方法的流程示意图;

[0038] 图2为本发明一个实施例的扩张卷积神经网络的结构示意图;

[0039] 图3为本发明一个实施例的基于扩张卷积神经网络的特征提取装置的结构示意图;

[0040] 图4示出了根据本发明一个实施例的电子设备的结构示意图;

[0041] 图5示出了根据本发明一个实施例的计算机可读存储介质的结构示意图。

具体实施方式

[0042] 为使本发明的目的、技术方案和优点更加清楚,下面将结合附图对本发明实施方式作进一步地详细描述。

[0043] 图1为本发明一个实施例的基于扩张卷积神经网络的特征提取方法的流程示意图。如图1所示,该方法包括:

[0044] S11:获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;

[0045] 可理解的是,对于待分析文本,首先通过NLP工具获取文本字符串(token串,中文的token通常表示一个字符),例如token串 $x=[x_1, x_2, x_3, \dots]$ 。为了满足神经网络的计算格式,对文本字符串进行裁剪或补充,设定超参数文本长度,小于该长度的需要补齐,大于

的需要截断,最终获取保护预设个数的字符的token串。

[0046] S12:对所述文本字符串进行词嵌入;

[0047] 在实际应用中,通过一层嵌入(embedding)层对文本字符串进行词嵌入,将文本字符串中的各个字符转化为计算机能够处理的数值向量。

[0048] S13:将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;

[0049] S14:通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。

[0050] 在得到各个字符的输出向量后,可以直接进行分类,或者通过CRF进行进一步的加工。CRF为条件随机场,可以对输出的概率分布进行进一步学习输出,学习到顺序化的数据。

[0051] 可理解的是,独立于序列长度的并行化运行可以节省时间和能源成本,最大化GPU资源使用并最大限度地减少培训和评估模型所需的时间。卷积神经网络(CNN)正好提供了这种特性。它们不是在序列中的每个标记上递增地组成表示,而是一次在整个序列中并行应用滤波器,此处的滤波器表示CNN里的卷积核。

[0052] 对于扩张的卷积神经网络,在同等参数量和计算量下,增加卷积核处理时的间隔距离,有效输入宽度可以随深度呈指数增长,每层的分辨率没有损失,并且估计的参数数量适中。与典型的CNN层一样,扩张的卷积在序列上下文的滑动窗口上操作,但与传统的卷积不同,上下文不需要是连续的;扩张的窗口会遍历每个扩张宽度为d的输入。有效输入宽度会随着扩张卷积层的数量进行指数型的增长可以扩展有效输入宽度的大小以仅使用几个卷积层覆盖大多数序列的整个长度。

[0053] 本发明实施例的基于扩张卷积神经网络的特征提取方法,通过扩张的卷积神经网络对待分析文本进行特征提取,有效输入宽度可以随深度呈指数增长,提高了特征提取的效率。

[0054] 在本发明实施例的一种可选的实施方式中,所述通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量,包括:

[0055] 将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。

[0056] 迭代扩张的卷积神经网络,相比传统的卷积神经网络可更好地用于大上下文和结构化预测,在提高效率的同时保持了类似的F1性能。F1性能作为一个评估指标,等于精度乘以召回率除以精度和召回率的和,精度为实际是正类并且被预测为正类的样本占有所有预测为正类的比例,召回率为实际是正类并且被预测为正类的样本占有所有实际为正类样本的比例。当使用独立分类执行预测时,始终优于双向LSTM(Bi-LSTM),并且与具有来自Bi-LSTM(Bi-LSTM-CRF)的对数的CRF中的推断相同。

[0057] 整体迭代扩张CNN架构重复地将相同的扩张卷积块应用于以令牌方式表示,令牌即token,按token级别进行卷积操作。与使用RNN生成的logits的模型类似,该迭代扩张的卷积神经网络提供了两种执行预测的方法:可以独立地预测每个标记的标签,或者通过在链式结构化图形模型中运行Viterbi推理。

[0058] 优选地,如图2所示,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。

[0059] 可理解的是,本发明实施例的卷积层的滤波器的宽度对应卷积层半径,滤波器的

宽度为3相当于卷积层半径为1,滤波器的宽度为5相当于卷积层的半径为2。

[0060] 进一步地,所述方法还包括:

[0061] 构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。

[0062] 在实际应用中,为进一步提高特征提取的效率,可以构建4个结构相同的扩张的卷积神经网络,将4个文本字符串分别并行输入到各个扩张的卷积神经网络,由4个扩张的卷积神经网络并行对相应的文本字符串进行特征提取。

[0063] 图3为本发明一个实施例的基于扩张卷积神经网络的特征提取装置的结构示意图。如图3所示,该装置包括:

[0064] 文本字符串获取单元31,用于获取待分析文本,对待分析文本进行处理获得至少一个文本字符串,所述文本字符串包括预设个数的字符;

[0065] 词嵌入处理单元32,用于对所述文本字符串进行词嵌入;

[0066] 文本字符串输入单元33,用于将经过词嵌入处理的文本字符串输入扩张的卷积神经网络,所述扩张的卷积神经网络包括多个卷积层,各个卷积层的滤波器的宽度大于1;

[0067] 输入向量获取单元34,用于通过所述扩张的卷积神经网络获得文本字符串中各个字符的输出向量。

[0068] 本发明实施例的基于扩张卷积神经网络的特征提取装置,通过扩张的卷积神经网络对待分析文本进行特征提取,有效输入宽度可以随深度呈指数增长,提高了特征提取的效率。

[0069] 进一步地,输入向量获取单元34进一步用于:

[0070] 将前一个卷积层的输出作为后一个卷积层的输入,层层迭代获得文本字符串中各个字符的输出向量。

[0071] 优选地,所述扩张的卷积神经网络包括3个卷积层,其中第1个卷积层的滤波器的宽度为3,第2个卷积层的滤波器的宽度为3,第3个卷积层的滤波器的宽度为5。

[0072] 进一步地,该装置还包括:并行处理单元,用于构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络。

[0073] 需要说明的是,上述实施例中的基于扩张卷积神经网络的特征提取装置可分别用于执行前述实施例中的方法,因此不再一一进行具体的说明。

[0074] 综上所述,本发明通过扩张的卷积神经网络对待分析文本进行特征提取,有效输入宽度可以随深度呈指数增长,提高了特征提取的效率。

[0075] 本发明迭代扩张的卷积神经网络,相比传统的卷积神经网络可更好地用于大上下文和结构化预测,在提高效率的同时保持了类似的F1性能。

[0076] 本发明还可以构建多个扩张的卷积神经网络,将获得的各个文本字符串分别并行输入到各个扩张的卷积神经网络,进一步提高了特征提取的效率。

[0077] 本领域内的技术人员应明白,本发明的实施例可提供为方法、系统、或计算机程序产品。因此,本发明可采用完全硬件实施例、完全软件实施例、或结合软件和硬件方面的实施例的形式。而且,本发明可采用在一个或多个其中包含有计算机可用程序代码的计算机可用存储介质(包括但不限于磁盘存储器、CD-ROM、光学存储器等)上实施的计算机程序产品的形式。

[0078] 需要说明的是：

[0079] 在此提供的算法和显示不与任何特定计算机、虚拟装置或者其它设备固有相关。各种通用装置也可以与基于在此的示教一起使用。根据上面的描述，构造这类装置所要求的结构是显而易见的。此外，本发明也不针对任何特定编程语言。应当明白，可以利用各种编程语言实现在此描述的本发明的内容，并且上面对特定语言所做的描述是为了披露本发明的最佳实施方式。

[0080] 在此处所提供的说明书中，说明了大量具体细节。然而，能够理解，本发明的实施例可以在没有这些具体细节的情况下实践。在一些实例中，并未详细示出公知的方法、结构和技术，以便不模糊对本说明书的理解。

[0081] 类似地，应当理解，为了精简本发明并帮助理解各个发明方面中的一个或多个，在上面对本发明的示例性实施例的描述中，本发明的各个特征有时被一起分组到单个实施例、图、或者对其的描述中。然而，并不应将该公开的方法解释成反映如下意图：即所要求保护的本发明要求比在每个权利要求中所明确记载的特征更多的特征。更确切地说，如下面的权利要求书所反映的那样，发明方面在于少于前面公开的单个实施例的所有特征。因此，遵循具体实施方式的权利要求书由此明确地并入该具体实施方式，其中每个权利要求本身都作为本发明的单独实施例。

[0082] 本领域那些技术人员可以理解，可以对实施例中的设备中的模块进行自适应性地改变并且把它们设置在与该实施例不同的一个或多个设备中。可以把实施例中的模块或单元或组件组合成一个模块或单元或组件，以及此外可以把它分成多个子模块或子单元或子组件。除了这样的特征和/或过程或者单元中的至少一些是相互排斥之外，可以采用任何组合对本说明书(包括伴随的权利要求、摘要和附图)中公开的所有特征以及如此公开的任何方法或者设备的所有过程或单元进行组合。除非另外明确陈述，本说明书(包括伴随的权利要求、摘要和附图)中公开的每个特征可以由提供相同、等同或相似目的的替代特征来代替。

[0083] 此外，本领域的技术人员能够理解，尽管在此所述的一些实施例包括其它实施例中包括的某些特征而不是其它特征，但是不同实施例的特征的组合意味着处于本发明的范围之内并且形成不同的实施例。例如，在下面的权利要求书中，所要求保护的实施例的任意之一都可以以任意的组合方式来使用。

[0084] 本发明的各个部件实施例可以以硬件实现，或者以在一个或者多个处理器上运行的软件模块实现，或者以它们的组合实现。本领域的技术人员应当理解，可以在实践中使用微处理器或者数字信号处理器(DSP)来实现根据本发明实施例的检测电子设备的佩戴状态的装置中的一些或者全部部件的一些或者全部功能。本发明还可以实现为用于执行这里所描述的方法的一部分或者全部的设备或者装置程序(例如，计算机程序和计算机程序产品)。这样的实现本发明的程序可以存储在计算机可读介质上，或者可以具有一个或者多个信号的形式。这样的信号可以从因特网网站上下载得到，或者在载体信号上提供，或者以任何其他形式提供。

[0085] 例如，图4示出了根据本发明一个实施例的电子设备的结构示意图。该电子设备传统上包括处理器41和被安排成存储计算机可执行指令(程序代码)的存储器42。存储器42可以是诸如闪存、EEPROM(电可擦除可编程只读存储器)、EPROM、硬盘或者ROM之类的电子存储

器。存储器42具有存储用于执行图1所示的以及各实施例中的任何方法步骤的程序代码44的存储空间43。例如,用于存储程序代码的存储空间43可以包括分别用于实现上面的方法中的各种步骤的各个程序代码44。这些程序代码可以从一个或者多个计算机程序产品中读出或者写入到这一个或者多个计算机程序产品中。这些计算机程序产品包括诸如硬盘,紧致盘(CD)、存储卡或者软盘之类的程序代码载体。这样的计算机程序产品通常为例如图5所述的计算机可读存储介质。该计算机可读存储介质可以具有与图4的电子设备中的存储器42类似布置的存储段、存储空间等。程序代码可以例如以适当形式进行压缩。通常,存储空间存储有用于执行根据本发明的方法步骤的程序代码51,即可以有诸如处理器41读取的程序代码,当这些程序代码由电子设备运行时,导致该电子设备执行上面所描述的方法中的各个步骤。

[0086] 以上所述,仅为本发明的具体实施方式,在本发明的上述教导下,本领域技术人员可以在上述实施例的基础上进行其他的改进或变形。本领域技术人员应该明白,上述的具体描述只是更好的解释本发明的目的,本发明的保护范围应以权利要求的保护范围为准。

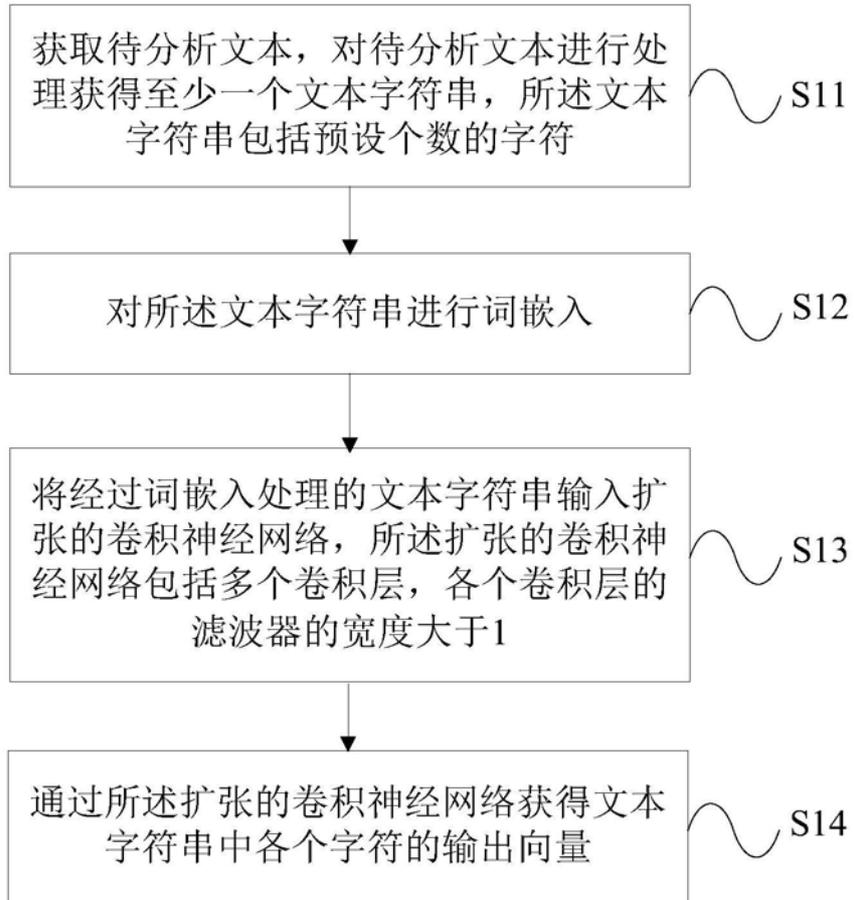


图1

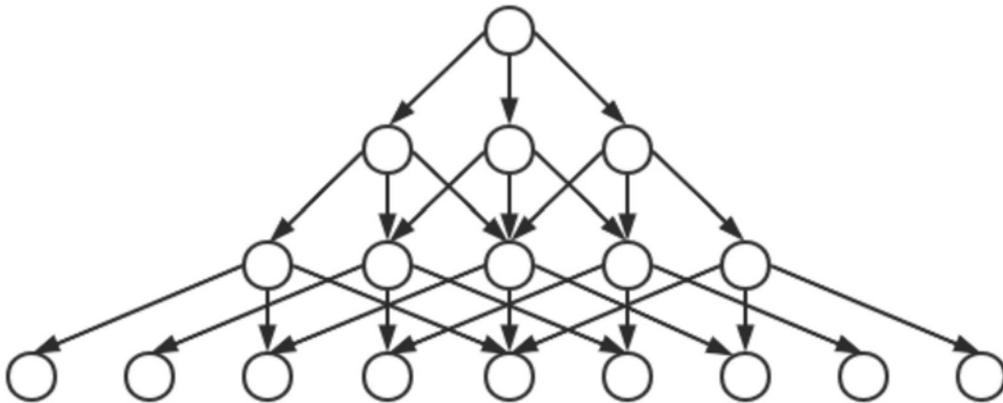


图2

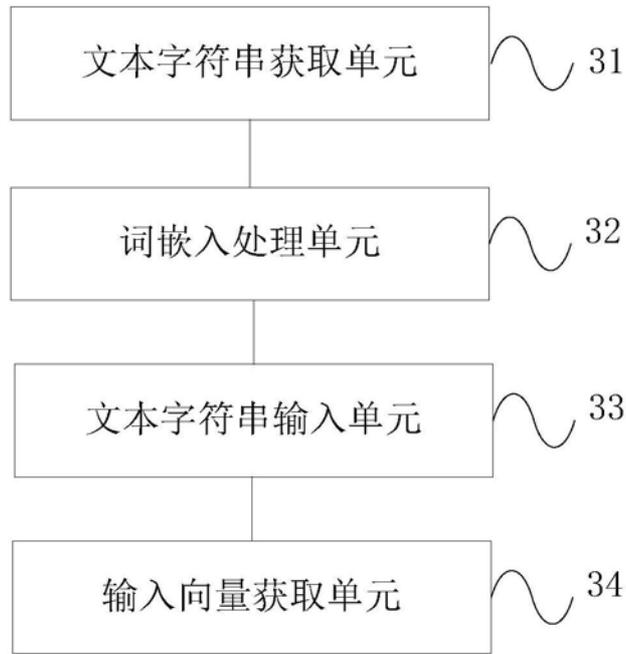


图3

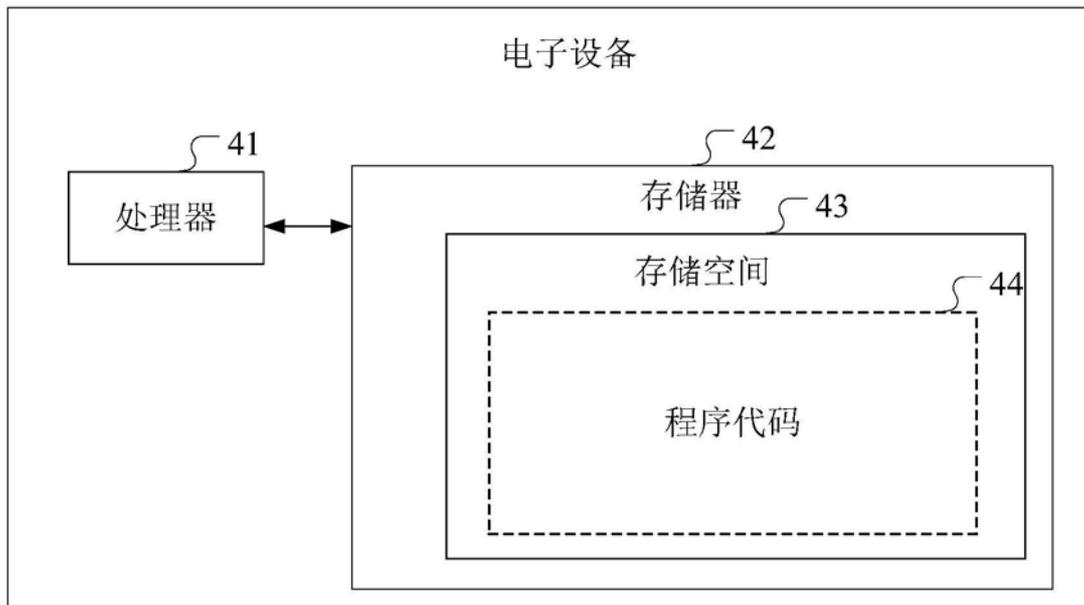


图4



图5