



(12)发明专利申请

(10)申请公布号 CN 107180628 A

(43)申请公布日 2017.09.19

(21)申请号 201710361210.7

(22)申请日 2017.05.19

(71)申请人 百度在线网络技术(北京)有限公司

地址 100085 北京市海淀区上地十街10号
百度大厦

(72)发明人 李超 马啸空 蒋兵 李先刚

(74)专利代理机构 北京鸿德海业知识产权代理
事务所(普通合伙) 11412

代理人 袁媛

(51) Int. Cl.

G10L 15/02(2006.01)

G10L 15/06(2013.01)

G10L 15/16(2006.01)

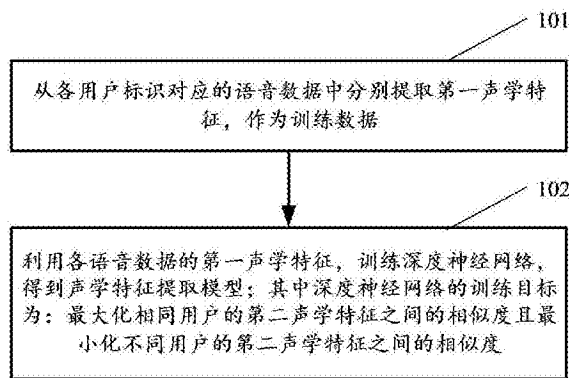
权利要求书2页 说明书11页 附图5页

(54)发明名称

建立声学特征提取模型的方法、提取声学特征的方法、装置

(57)摘要

本发明提供了一种建立声学特征提取模型的方法、提取声学特征的方法、装置。其中建立声学特征提取模型的方法包括:将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;利用所述训练数据训练深度神经网络,得到声学特征提取模型;其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。本发明的声学特征提取模型能够自学习到达到训练目标的最优声学特征。相比较现有预设特征类型和变换方式的声学特征提取方式,实现更加灵活,准确性更高。



1. 一种建立声学特征提取模型的方法,其特征在于,该方法包括:
将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;
利用所述训练数据训练深度神经网络,得到声学特征提取模型;
其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。
2. 根据权利要求1所述的方法,其特征在于,所述第一声学特征包括:FBank64声学特征。
3. 根据权利要求1所述的方法,其特征在于,所述深度神经网络包括:卷积神经网络CNN、残差卷积神经网络ResCNN或者门控递归单元GRU。
4. 根据权利要求1所述的方法,其特征在于,利用所述训练数据训练深度神经网络,得到声学特征提取模型包括:
利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征;
利用所述各语音数据的第二声学特征计算三元组损失,利用所述三元组损失对所述深度神经网络进行调参,以最小化所述三元组损失;
其中,所述三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。
5. 根据权利要求4所述的方法,其特征在于,所述利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征包括:
利用深度神经网络对各语音数据的第一声学特征进行学习,输出帧级别的第二声学特征;
对帧级别的第二声学特征进行池化和语句标准化处理,输出句子级别的第二声学特征;
在计算三元组损失时利用的所述各语音数据的第二声学特征为各语音数据的句子级别的第二声学特征。
6. 一种提取声学特征的方法,其特征在于,该方法包括:
提取待处理语音数据的第一声学特征;
将所述第一声学特征输入声学特征提取模型,得到待处理语音数据的第二声学特征;
其中所述声学特征提取模型是采用如权利要求1至5任一权利要求所述方法预先建立的。
7. 根据权利要求6所述的方法,其特征在于,该方法还包括:
利用所述待处理语音数据的第二声学特征,注册所述待处理语音数据所对应用户标识的声纹模型;或者,
将所述待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配,确定所述待处理语音数据对应的用户标识。
8. 一种建立声学特征提取模型的装置,其特征在于,该装置包括:
数据获取单元,用于将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;
模型训练单元,用于利用所述训练数据训练深度神经网络,得到声学特征提取模型;其

中所述深度神经网络的训练目标为：最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

9. 根据权利要求8所述的装置，其特征在于，所述第一声学特征包括：FBank64声学特征。

10. 根据权利要求8所述的装置，其特征在于，所述深度神经网络包括：卷积神经网络CNN、残差卷积神经网络ResCNN或者门控递归单元GRU。

11. 根据权利要求8所述的装置，其特征在于，所述模型训练单元，具体用于：

利用深度神经网络对各语音数据的第一声学特征进行学习，输出各语音数据的第二声学特征；

利用所述各语音数据的第二声学特征计算三元组损失，利用所述三元组损失对所述深度神经网络进行调参，以最小化所述三元组损失；

其中，所述三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。

12. 根据权利要求11所述的装置，其特征在于，所述模型训练单元在利用深度神经网络对各语音数据的第一声学特征进行学习，输出各语音数据的第二声学特征时，具体执行：利用深度神经网络对各语音数据的第一声学特征进行学习，输出帧级别的第二声学特征；对帧级别的第二声学特征进行池化和语句标准化处理，输出句子级别的第二声学特征；

在模型训练单元计算三元组损失时利用的所述各语音数据的第二声学特征为各语音数据的句子级别的第二声学特征。

13. 一种提取声学特征的装置，其特征在于，该装置包括：

预处理单元，用于提取待处理语音数据的第一声学特征；

特征提取单元，用于将所述第一声学特征输入声学特征提取模型，得到待处理语音数据的第二声学特征；

其中所述声学特征提取模型是由权利要求8至12任一权利要求所述装置预先建立的。

14. 根据权利要求13所述的装置，其特征在于，该装置还包括：

声纹注册单元，用于利用所述待处理语音数据的第二声学特征，注册所述待处理语音数据所对应用户标识的声纹模型；或者，

声纹匹配单元，用于将所述待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配，确定所述待处理语音数据对应的用户标识。

15. 一种设备，包括

存储器，包括一个或者多个程序；

一个或者多个处理器，耦合到所述存储器，执行所述一个或者多个程序，以实现如权利要求1至7任一权项所述方法中执行的操作。

16. 一种计算机存储介质，所述计算机存储介质被编码有计算机程序，所述程序在被一个或多个计算机执行时，使得所述一个或多个计算机执行如权利要求1至7任一权项所述方法中执行的操作。

建立声学特征提取模型的方法、提取声学特征的方法、装置

【技术领域】

[0001] 本发明涉及计算机应用技术领域,特别涉及一种建立声学特征提取模型的方法、提取声学特征的方法及对应装置。

【背景技术】

[0002] 随着人工智能的不断发展,语音交互已经成为最自然的交互方式之一得到日益推广,语音识别技术也越来越得到人们的重视。在语音识别技术中,声学特征的提取是核心技术,其可以用于用户识别、验证或分类等。

[0003] 现有声学特征提取方式,主要是依据预设的特征类型,对语音数据进行预设方式的变换后,从中提取对应类型的特征。这种声学特征提取方式很大程度上依靠特征类型的设置和变换方式的设置,准确性和灵活性较低。

【发明内容】

[0004] 本发明提供了一种建立声学特征提取模型的方法、提取声学特征的方法、装置、设备和计算机存储介质,以便于提高所提取声学特征的准确性和灵活性。

[0005] 具体技术方案如下:

[0006] 本发明提供了一种建立声学特征提取模型的方法,该方法包括:

[0007] 将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;

[0008] 利用所述训练数据训练深度神经网络,得到声学特征提取模型;

[0009] 其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0010] 根据本发明一优选实施方式,所述第一声学特征包括:FBank64声学特征。

[0011] 根据本发明一优选实施方式,所述深度神经网络包括:卷积神经网络CNN、残差卷积神经网络ResCNN或者门控递归单元GRU。

[0012] 根据本发明一优选实施方式,利用所述训练数据训练深度神经网络,得到声学特征提取模型包括:

[0013] 利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征;

[0014] 利用所述各语音数据的第二声学特征计算三元组损失,利用所述三元组损失对所述深度神经网络进行调参,以最小化所述三元组损失;

[0015] 其中,所述三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。

[0016] 根据本发明一优选实施方式,所述利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征包括:

[0017] 利用深度神经网络对各语音数据的第一声学特征进行学习,输出帧级别的第二声学特征;

- [0018] 对帧级别的第二声学特征进行池化和语句标准化处理,输出句子级别的第二声学特征;
- [0019] 在计算三元组损失时利用的所述各语音数据的第二声学特征为各语音数据的句子级别的第二声学特征。
- [0020] 本发明还提供了一种提取声学特征的方法,该方法包括:
- [0021] 提取待处理语音数据的第一声学特征;
- [0022] 将所述第一声学特征输入声学特征提取模型,得到待处理语音数据的第二声学特征;
- [0023] 其中所述声学特征提取模型是采用上述建立声学特征提取模型的方法预先建立的。
- [0024] 根据本发明一优选实施方式,该方法还包括:
- [0025] 利用所述待处理语音数据的第二声学特征,注册所述待处理语音数据所对应用户标识的声纹模型;或者,
- [0026] 将所述待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配,确定所述待处理语音数据对应的用户标识。
- [0027] 本发明还提供了一种建立声学特征提取模型的装置,该装置包括:
- [0028] 数据获取单元,用于将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;
- [0029] 模型训练单元,用于利用所述训练数据训练深度神经网络,得到声学特征提取模型;其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。
- [0030] 根据本发明一优选实施方式,所述第一声学特征包括:FBank64声学特征。
- [0031] 根据本发明一优选实施方式,所述深度神经网络包括:卷积神经网络CNN、残差卷积神经网络ResCNN或者门控递归单元GRU。
- [0032] 根据本发明一优选实施方式,所述模型训练单元,具体用于:
- [0033] 利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征;
- [0034] 利用所述各语音数据的第二声学特征计算三元组损失,利用所述三元组损失对所述深度神经网络进行调参,以最小化所述三元组损失;
- [0035] 其中,所述三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。
- [0036] 根据本发明一优选实施方式,所述模型训练单元在利用深度神经网络对各语音数据的第一声学特征进行学习,输出各语音数据的第二声学特征时,具体执行:利用深度神经网络对各语音数据的第一声学特征进行学习,输出帧级别的第二声学特征;对帧级别的第二声学特征进行池化和语句标准化处理,输出句子级别的第二声学特征;
- [0037] 在模型训练单元计算三元组损失时利用的所述各语音数据的第二声学特征为各语音数据的句子级别的第二声学特征。
- [0038] 本发明还提供了一种提取声学特征的装置,该装置包括:
- [0039] 预处理单元,用于提取待处理语音数据的第一声学特征;

[0040] 特征提取单元,用于将所述第一声学特征输入声学特征提取模型,得到待处理语音数据的第二声学特征;

[0041] 其中所述声学特征提取模型是由上述建立声学特征提取模型的装置预先建立的。

[0042] 根据本发明一优选实施方式,该装置还包括:

[0043] 声纹注册单元,用于利用所述待处理语音数据的第二声学特征,注册所述待处理语音数据所对应用户标识的声纹模型;或者,

[0044] 声纹匹配单元,用于将所述待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配,确定所述待处理语音数据对应的用户标识。

[0045] 本发明提供了一种设备,包括

[0046] 存储器,包括一个或者多个程序;

[0047] 一个或者多个处理器,耦合到所述存储器,执行所述一个或者多个程序,以实现上述方法中执行的操作。

[0048] 本发明还提供了一种计算机存储介质,所述计算机存储介质被编码有计算机程序,所述程序在被一个或多个计算机执行时,使得所述一个或多个计算机执行上述方法中执行的操作。

[0049] 由以上技术方案可以看出,本发明基于神经网络,以最小化相同用户的第二声学特征之间的相似度且最大化不同用户的第二声学特征之间的相似度为目标,训练得到声学特征提取模型。也就是说,本发明的声学特征提取模型能够自学习以达到训练目标的最优声学特征。相比较现有预设特征类型和变换方式的声学特征提取方式,实现更加灵活,准确性更高。

【附图说明】

[0050] 图1为本发明实施例提供的建立声学特征提取模型的方法流程图;

[0051] 图2为本发明实施例提供的声学特征提取模型的结构示意图;

[0052] 图3为本发明实施例提供的堆叠残差块的结构示意图;

[0053] 图4为本发明实施例提供的提取声学特征的方法流程图;

[0054] 图5为本发明实施例提供的建立声学特征提取模型的装置结构图;

[0055] 图6为本发明实施例提供的提取声学特征的装置结构图;

[0056] 图7为实现本发明实施方式的示例性计算机系统/服务器的框图。

【具体实施方式】

[0057] 为了使本发明的目的、技术方案和优点更加清楚,下面结合附图和具体实施例对本发明进行详细描述。

[0058] 在本发明实施例中使用的术语是仅仅出于描述特定实施例的目的,而非旨在限制本发明。在本发明实施例和所附权利要求书中所使用的单数形式的“一种”、“所述”和“该”也旨在包括多数形式,除非上下文清楚地表示其他含义。

[0059] 应当理解,本文中使用的术语“和/或”仅仅是一种描述关联对象的关联关系,表示可以存在三种关系,例如,A和/或B,可以表示:单独存在A,同时存在A和B,单独存在B这三种情况。另外,本文中字符“/”,一般表示前后关联对象是一种“或”的关系。

[0060] 取决于语境,如在此所使用的词语“如果”可以被解释成为“在……时”或“当……时”或“响应于确定”或“响应于检测”。类似地,取决于语境,短语“如果确定”或“如果检测(陈述的条件或事件)”可以被解释成为“当确定时”或“响应于确定”或“当检测(陈述的条件或事件)时”或“响应于检测(陈述的条件或事件)”。

[0061] 本发明的核心思想在于,基于深度神经网络提取语音数据的高层声学特征,训练深度神经网络的目标是最大化相同用户的高层声学特征之间的相似度且最小化不同用户的高层声学特征之间的相似度,从而得到声学特征提取模型。该声学特征提取模型用于提取语音数据的高层声学特征。

[0062] 另外,在进行声学特征提取模型的训练时,需要首先对训练数据中的语音数据进行低层声学特征的提取,即进行预处理。其中该低层声学特征相对于高层声学特征粒度更粗,包含信息量也更粗;相反,经过声学特征提取模型处理后得到的高层声学特征相对于低层声学特征粒度更细,包含信息量也更细致,更适于建立声纹模型,以进行用户声纹的建立。在本发明实施例中,为了对这两种声学特征进行区分,将对语音数据进行预处理后得到的低层声学特征称为第一声学特征;将经过声学特征提取模型对低层声学特征进行处理后,得到的高层声学特征称为第二声学特征。

[0063] 在本发明中存在两个阶段:声学特征提取模型的建立阶段以及利用声学特征提取模型提取声学特征的阶段。其中两个阶段互相独立,声学特征提取模型的建立阶段可以是预先执行的阶段,但也可以在后续过程中不断对声学特征提取模型进行更新。下面结合实施例对这两个进行详细描述。

[0064] 图1为本发明实施例提供的建立声学特征提取模型的方法流程图,如图1所示,该方法可以包括以下步骤:

[0065] 在101中,从各用户标识对应的语音数据中分别提取第一声学特征,作为训练数据。

[0066] 可以预先采集已知用户的语音数据,在选择训练数据时可以对这些语音数据有一些质量要求,例如选取清晰度较好的语音数据,再例如删除长度过长或过短的语音数据,等等。

[0067] 对于采集到的语音数据首先进行预处理,从中提取各语音数据的第一声学特征。如前面所述的,该第一声学特征是低层的声学特征。在本发明实施例中可以采用FBank(Mel-scale Filter Bank,梅尔标度滤波器组)特征作为第一声学特征。例如,以25ms为一帧、10ms为步长提取语音数据的FBank特征。但本发明并不限于FBank特征,还可以采用其他特征作为第一声学特征。

[0068] 这样,就可以得到各用户标识对应的第一声学特征,从而构成训练数据。其中本发明并不限定用户标识的具体类型,可以是任意类型的标识,只要能够区分用户即可。在训练数据中可以包含同一用户对应的不同语音数据的第一声学特征,不同用户对应的语音数据的第一声学特征,等等。训练数据中各第一声学特征均具有对应的用户标识作为标签。

[0069] 在102中,利用各语音数据的第一声学特征,训练深度神经网络,得到声学特征提取模型;其中深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0070] 本步骤中,可以首先利用深度神经网络对各语音数据的第一声学特征进行学习,

输出各语音数据的第二声学特征。然后利用各语音数据的第二声学特征计算三元组损失，将三元组损失反馈给深度神经网络，以便调整深度神经网络的参数以最小化该三元组损失。

[0071] 为了方便对本发明实施例提供的声学特征提取模型的理解，对该声学特征提取模型的结构进行介绍。如图2所示，该声学特征提取模型可以包括深度神经网络层、池化和句子标准化层以及三元组损失层。

[0072] 其中，深度神经网络可以采用CNN、GRU (Gated Recurrent Unit, 门控递归单元) 等，当然也可以采用其他诸如RNN、LSTM等其他类型的深度神经网络。由于CNN相比较RNN、LSTM等而言，能够更加有效地减小频谱变化以及将频谱相关性在声学特征中进行体现，因此在本发明实施例中优选CNN这种类型的深度神经网络。

[0073] 然而，尽管深度神经网络具有很好地学习能力，但更难进行训练，在一定深度情况下准确性反而下滑。为了解决该问题，本发明可以基于CNN使用但不限于ResNet (Residual Net, 残差网络) 型CNN，或者采用GRU。

[0074] 首先对ResNET型CNN进行介绍。

[0075] ResNet可以用于简化CNN的训练。ResNet包括若干ResBlock (堆叠残差块)，各ResBlock包括低层输出和高层输入间的直接连接。如图3中所示，各ResBlock可以定义为：

[0076] $h = F(x, W_i) + x$

[0077] 其中， x 和 h 分别表示ResBlock的输入和输出。 F 表示堆叠的非线性层的映射函数。

[0078] 如图3所示，ResBlock可以包括两个卷积层和两个激活层。其中，两个卷积层可以包括诸如 3×3 的过滤器和 1×1 的stride (步幅)。每个ResBlock包括相同的结构，并且跳转连接是对 x 的相同映射。若通道的数量增加，则可以使用一个卷积层 (例如具有 5×5 的过滤器和 2×2 的stride)。因此，频率维度始终在卷积层中保持恒定。经过研究发现，语音识别在时间维度上对stride并不敏感。在本发明实施例中，可以使用如下ReLU (Rectified Linear Units, 修正线性) 函数作为所有激活层的非线性处理：

[0079] $\sigma(x) = \min\{\max\{x, 0\}, 20\}$

[0080] 下面对GRU进行介绍。

[0081] GRU相比较LSTM而言，训练速度更快且发散程度更小。本发明实施例中深度神经网络层可以采用多个GRU构成。例如，每个GRU可以包括一个 5×5 过滤器和 2×2 stride的卷积层，能够减少时域和频域的维度，从而允许GRU的计算速度更快。紧接着卷积层的是三个具有1024个单元的前向GRU层，在时间维度上进行循环。在GRU中也可以采用诸如ReLU进行激活。

[0082] 紧接着深度神经网络层的是池化和句子标准化层。池化和句子标准化层用来对深度神经网络层输出的帧级别的第二声学特征进行池化和句子标准化处理，从而得到句子级别的第二声学特征。

[0083] 具体地，如图2所示，池化和句子标准化层可以包括：池化层、仿射层和句子标准化层。

[0084] 其中池化层用于将帧级别的输入转变为句子级别的表示，即将帧级别的第二声学特征进行取平均，得到句子级别的第二声学特征。

[0085] 池化层的输出 h' 可以采用如下公式：

$$[0086] \quad h' = \frac{1}{T} \sum_{t=0}^{T-1} x'(t)$$

[0087] 其中, T为句子包含的帧数目, $x'(t)$ 为池化层的输入。

[0088] 经过池化 (Pooling) 层的处理, 使得本发明实施例提供的声学特征提取模型能够处理不同长度的语句, 解决了文本无关的情况。

[0089] 仿射层将句子级别的第二声学特征投射到预设的维度, 例如投射到512维度。

[0090] 长度标准化层将仿射层输出的句子级别的第二声学特征的长度进行规整, 使模为1。

[0091] 本发明实施例中, 三元损失层采用三元损失对深度神经网络层进行反馈训练, 以最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0092] 三元损失层可以采用三个样本作为输入: 锚样本, 包括一个用户的句子级别的第二声学特征; 正样本, 包括与锚样本同一用户的另一句子级别的第二声学特征; 负样本, 包括与锚样本不同用户的句子级别的第二声学特征。将上述样本构成一个三元组。

[0093] 三元损失层对深度神经网络层进行反馈, 以使得锚样本和正样本之间的余弦相似度 (在本发明实施例中样本之间的相似度采用余弦相似度体现, 但不排除其他相似度计算方式) 大于锚样本和负样本之间的余弦相似度。形式上,

$$[0094] \quad s_i^{ap} - \alpha > s_i^{an}$$

[0095] 其中, s_i^{ap} 为三元组 i 中锚样本 a 和正样本 p 之间的余弦相似度, s_i^{an} 为三元组 i 中锚样本 a 和正样本 n 之间的余弦相似度。训练目标是找到这些相似度中的最小边缘 α 。即计算三元组损失, 该三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。例如该三元组损失的计算函数 L 可以为:

$$[0096] \quad L = \sum_{i=0}^N [s_i^{an} - s_i^{ap} + \alpha]_+$$

[0097] 其中, N 为三元组的数目, 操作符 $[x]_+ = \max(x, 0)$ 。

[0098] 计算出的三元组损失反馈给深度神经网络层, 以不断调整深度神经网络层的参数, 从而逐渐训练深度神经网络, 最终最小化利用提取的第二声学特征计算的三元组损失。训练结束后, 得到声学特征提取模型, 此次训练过程结束。

[0099] 图4为本发明实施例提供的提取声学特征的方法流程图, 该流程基于如图1所示实施例建立的声学特征提取模型。如图4所示, 该方法可以包括以下步骤:

[0100] 在401中, 提取待处理语音数据的第一声学特征。

[0101] 本步骤是对待处理语音数据的预处理, 即从中提取第一声学特征, 该第一声学特征是低层的声学特征。此处提取的第一声学特征的类型和方式与图1所示实施例中步骤101中提取第一声学特征的类型和方式一致。在此不再赘述。

[0102] 在402中, 将提取出的第一声学特征输入声学特征提取模型, 得到待处理语音数据的第二声学特征。

[0103] 对于预先训练得到的声学特征提取模型, 由于其从训练数据中已经完成从第一声学特征到第二声学特征的自学习, 因此将步骤401中提取出的待处理语音数据的第一声学特征输入声学特征提取模型, 声学特征提取模型就能够输出待处理语音数据的第二声学特

征。该第二声学特征可以为句子级别的高层声学特征。

[0104] 在得到待处理语音数据的第二声学特征后,可以利用第二声学特征进行后续应用的处理,例如在403a中,利用待处理语音数据的第二声学特征,注册该待处理语音数据所对应用户标识的声纹模型,或者在403b中,将待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配,确定待处理语音数据对应的用户标识。

[0105] 在403a中,若待处理语音数据对应的用户标识已知,则可以利用提取的第二声学特征注册该用户标识对应的声纹模型。在注册声纹模型时,可以将提取的第二声学特征进行处理后,作为声纹信息存储于声纹模型库中。可以利用用户标识对应的一个或多个第二声学特征来进行声学模型的注册,具体注册过程本发明不做具体限制。

[0106] 在403b中,若待处理语音数据对应的用户标识未知,则可以利用提取的第二声学特征与声纹模型库中各已注册的声纹模型进行匹配,例如通过计算提取的第二声学特征与声纹模型库中各声纹模型之间相似度的方式进行匹配。若匹配到某个声纹模型,则可以确定该待处理语音数据对应该匹配到的声纹模型对应的用户标识。

[0107] 上述403a和403b是本发明实施例提供的两种在提取语音数据的第二声学特征后,对其的应用方式,当然除了这两种应用方式之外,还可以进行其他应用,本发明不做一一穷举。

[0108] 上述方法可以应用于语音识别系统中,执行主体可以为对应装置,该装置可以是位于用户设备的应用,或者还可以为位于用户设备的应用中的插件或软件开发工具包(Software Development Kit, SDK)等功能单元。其中,用户设备可以包括但不限于诸如:智能移动终端、智能家居设备、网络设备、可穿戴式设备、智能医疗设备、PC(个人计算机)等。其中智能移动设备可以包括诸如手机、平板电脑、笔记本电脑、PDA(个人数字助理)、互联网汽车等。智能家居设备可以包括智能家电设备,诸如智能电视、智能空调、智能热水器、智能冰箱、智能空气净化器等,智能家居设备还可以包括智能门锁、智能电灯、智能摄像头等。网络设备可以包括诸如交换机、无线AP、服务器等。可穿戴式设备可以包括诸如智能手表、智能眼镜、智能手环、虚拟现实设备、增强现实设备、混合现实设备(即可以支持虚拟现实和增强现实的设备)等等。智能医疗设备可以包括诸如智能体温计、智能血压仪、智能血糖仪等等。

[0109] 图5为本发明实施例提供的建立声学特征提取模型的装置结构图,该装置可以用于执行如图1中所示的操作。如图5所示,该装置可以包括:数据获取单元01和模型训练单元02。其中各组成单元的主要功能如下:

[0110] 数据获取单元01负责将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据。

[0111] 可以预先采集已知用户的语音数据,在选择训练数据时可以对这些语音数据有一些质量要求,例如选取清晰度较好的语音数据,再例如删除长度过长或过短的语音数据,等等。

[0112] 在本发明实施例中可以采用FBank特征作为第一声学特征。例如,以25ms为一帧、10ms为步长提取语音数据的FBank特征。但本发明并不限于FBank特征,还可以采用其他特征作为第一声学特征。

[0113] 模型训练单元02负责利用训练数据训练深度神经网络,得到声学特征提取模型;

其中深度神经网络的训练目标为：最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0114] 优选地，本发明实施例中采用的深度神经网络可以包括：CNN、ResCNN或者GRU。ResCNN和GRU的相关描述参见方法实施例中的记载，在此不再赘述。

[0115] 具体地，模型训练单元02可以首先利用深度神经网络对各语音数据的第一声学特征进行学习，输出各语音数据的第二声学特征；然后利用各语音数据的第二声学特征计算三元组损失，将三元组损失反馈给深度神经网络，以便调整深度神经网络的参数以最小化三元组损失；其中，三元组损失体现不同用户的第二声学特征之间的相似度与相同用户的第二声学特征之间的相似度的差值状况。

[0116] 更具体地，模型训练单元02在利用深度神经网络对各语音数据的第一声学特征进行学习，输出各语音数据的第二声学特征时，可以具体执行：利用深度神经网络对各语音数据的第一声学特征进行学习，输出帧级别的第二声学特征；对帧级别的第二声学特征进行池化和语句标准化处理，输出句子级别的第二声学特征。此时，在模型训练单元02计算三元组损失时利用的各语音数据的第二声学特征为各语音数据的句子级别的第二声学特征。

[0117] 对于声学特征提取模型的具体架构以及该架构中各层次所执行的具体处理可以参见方法实施例中的相关描述，在此不再赘述。

[0118] 图6为本发明实施例提供的提取声学特征的装置结构图，如图6所示，该装置可以包括：预处理单元11和特征提取单元12。其中各组成单元的主要功能如下：

[0119] 预处理单元11负责提取待处理语音数据的第一声学特征。该第一声学特征的类型和提取方式与图5中数据获取单元01获取训练数据时所采用的第一声学特征的类型和提取方式一致。例如，第一声学特征可以采用FBank特征。

[0120] 特征提取单元12负责将第一声学特征输入声学特征提取模型，得到待处理语音数据的第二声学特征。

[0121] 在得到待处理语音数据的第二声学特征后，可以利用第二声学特征进行后续应用的处理，例如该装置还可以包括：

[0122] 声纹注册单元(图中未示出)，负责利用待处理语音数据的第二声学特征，注册该待处理语音数据所对应用户标识的声纹模型。

[0123] 再例如，该装置还可以包括：

[0124] 声纹匹配单元(图中未示出)，负责将待处理语音数据的第二声学特征与已注册的各用户标识的声纹模型进行匹配，确定待处理语音数据对应的用户标识。

[0125] 图7示出了适于用来实现本发明实施方式的示例性计算机系统/服务器012的框图。图7显示的计算机系统/服务器012仅仅是一个示例，不应对本发明实施例的功能和使用范围带来任何限制。

[0126] 如图7所示，计算机系统/服务器012以通用计算设备的形式表现。计算机系统/服务器012的组件可以包括但不限于：一个或者多个处理器或者处理单元016，系统存储器028，连接不同系统组件(包括系统存储器028和处理单元016)的总线018。

[0127] 总线018表示几类总线结构中的一种或多种，包括存储器总线或者存储器控制器，外围总线，图形加速端口，处理器或者使用多种总线结构中的任意总线结构的局域总线。举例来说，这些体系结构包括但不限于工业标准体系结构(ISA)总线，微通道体系结构(MAC)

总线,增强型ISA总线、视频电子标准协会(VESA)局域总线以及外围组件互连(PCI)总线。

[0128] 计算机系统/服务器012典型地包括多种计算机系统可读介质。这些介质可以是任何能够被计算机系统/服务器012访问的可用介质,包括易失性和非易失性介质,可移动的和不可移动的介质。

[0129] 系统存储器028可以包括易失性存储器形式的计算机系统可读介质,例如随机存取存储器(RAM)030和/或高速缓存存储器032。计算机系统/服务器012可以进一步包括其它可移动/不可移动的、易失性/非易失性计算机系统存储介质。仅作为举例,存储系统034可以用于读写不可移动的、非易失性磁介质(图7未显示,通常称为“硬盘驱动器”)。尽管图7中未示出,可以提供用于对可移动非易失性磁盘(例如“软盘”)读写的磁盘驱动器,以及对可移动非易失性光盘(例如CD-ROM,DVD-ROM或者其它光介质)读写的光盘驱动器。在这些情况下,每个驱动器可以通过一个或者多个数据介质接口与总线018相连。存储器028可以包括至少一个程序产品,该程序产品具有一组(例如至少一个)程序模块,这些程序模块被配置以执行本发明各实施例的功能。

[0130] 具有一组(至少一个)程序模块042的程序/实用工具040,可以存储在例如存储器028中,这样的程序模块042包括——但不限于——操作系统、一个或者多个应用程序、其它程序模块以及程序数据,这些示例中的每一个或某种组合中可能包括网络环境的实现。程序模块042通常执行本发明所描述的实施例中的功能和/或方法。

[0131] 计算机系统/服务器012也可以与一个或多个外部设备014(例如键盘、指向设备、显示器024等)通信,在本发明中,计算机系统/服务器012与外部雷达设备进行通信,还可与一个或者多个使得用户能与该计算机系统/服务器012交互的设备通信,和/或与使得该计算机系统/服务器012能与一个或多个其它计算设备进行通信的任何设备(例如网卡,调制解调器等等)通信。这种通信可以通过输入/输出(I/O)接口022进行。并且,计算机系统/服务器012还可以通过网络适配器020与一个或者多个网络(例如局域网(LAN),广域网(WAN)和/或公共网络,例如因特网)通信。如图所示,网络适配器020通过总线018与计算机系统/服务器012的其它模块通信。应当明白,尽管图7中未示出,可以结合计算机系统/服务器012使用其它硬件和/或软件模块,包括但不限于:微代码、设备驱动器、冗余处理单元、外部磁盘驱动阵列、RAID系统、磁带驱动器以及数据备份存储系统等。

[0132] 处理单元016通过运行存储在系统存储器028中的程序,从而执行各种功能应用以及数据处理,例如实现一种建立声学特征提取模型的方法,可以包括:

[0133] 将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;

[0134] 利用所述训练数据训练深度神经网络,得到声学特征提取模型;

[0135] 其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0136] 再例如,实现一种提取声学特征的方法,可以包括:

[0137] 提取待处理语音数据的第一声学特征;

[0138] 将所述第一声学特征输入声学特征提取模型,得到待处理语音数据的第二声学特征。

[0139] 上述的计算机程序可以设置于计算机存储介质中,即该计算机存储介质被编码有计算机程序,该程序在被一个或多个计算机执行时,使得一个或多个计算机执行本发明上

述实施例中所示的方法流程和/或装置操作。例如,被上述一个或多个处理器执行的方法流程,可以包括:

[0140] 将从各用户标识对应的语音数据中分别提取的第一声学特征,作为训练数据;

[0141] 利用所述训练数据训练深度神经网络,得到声学特征提取模型;

[0142] 其中所述深度神经网络的训练目标为:最大化相同用户的第二声学特征之间的相似度且最小化不同用户的第二声学特征之间的相似度。

[0143] 再例如,被上述一个或多个处理器执行的方法流程,可以包括:

[0144] 提取待处理语音数据的第一声学特征;

[0145] 将所述第一声学特征输入声学特征提取模型,得到待处理语音数据的第二声学特征。

[0146] 随着时间、技术的发展,介质含义越来越广泛,计算机程序的传播途径不再受限于有形介质,还可以直接从网络下载等。可以采用一个或多个计算机可读的介质的任意组合。计算机可读介质可以是计算机可读信号介质或者计算机可读存储介质。计算机可读存储介质例如可以是一——但不限于——电、磁、光、电磁、红外线、或半导体的系统、装置或器件,或者任意以上的组合。计算机可读存储介质的更具体的例子(非穷举的列表)包括:具有一个或多个导线的电连接、便携式计算机磁盘、硬盘、随机存取存储器(RAM)、只读存储器(ROM)、可擦式可编程只读存储器(EPROM或闪存)、光纤、便携式紧凑磁盘只读存储器(CD-ROM)、光存储器件、磁存储器件、或者上述的任意合适的组合。在本文件中,计算机可读存储介质可以是任何包含或存储程序的有形介质,该程序可以被指令执行系统、装置或者器件使用或者与其结合使用。

[0147] 计算机可读的信号介质可以包括在基带中或者作为载波一部分传播的数据信号,其中承载了计算机可读的程序代码。这种传播的数据信号可以采用多种形式,包括——但不限于——电磁信号、光信号或上述的任意合适的组合。计算机可读的信号介质还可以是计算机可读存储介质以外的任何计算机可读介质,该计算机可读介质可以发送、传播或者传输用于由指令执行系统、装置或者器件使用或者与其结合使用的程序。

[0148] 计算机可读介质上包含的程序代码可以用任何适当的介质传输,包括——但不限于——无线、电线、光缆、RF等等,或者上述的任意合适的组合。

[0149] 可以以一种或多种程序设计语言或其组合来编写用于执行本发明操作的计算机程序代码,所述程序设计语言包括面向对象的程序设计语言——诸如Java、Smalltalk、C++,还包括常规的过程式程序设计语言——诸如“C”语言或类似的程序设计语言。程序代码可以完全地在用户计算机上执行、部分地在用户计算机上执行、作为一个独立的软件包执行、部分在用户计算机上部分在远程计算机上执行、或者完全在远程计算机或服务器上执行。在涉及远程计算机的情形中,远程计算机可以通过任意种类的网络——包括局域网(LAN)或广域网(WAN)连接到用户计算机,或者,可以连接到外部计算机(例如利用因特网服务提供商来通过因特网连接)。

[0150] 由以上描述可以看出,本发明提供的上述方法、装置、设备和计算机存储介质可以具备以下优点:

[0151] 1) 本发明的声学特征提取模型能够自学习到达到训练目标的最优声学特征。相比较现有预设特征类型和变换方式的声学特征提取方式,实现更加灵活,准确性更高。

[0152] 2) 本发明中优选ResCNN或GRU类型的深度神经网络,从而在采用较高层级深度的神经网络情况下,也能够保证特征提取的准确性,且提高深度神经网络的训练速度。

[0153] 3) 本发明在训练声学特征提取模型的过程中,对深度神经网络的输出进行池化和句子标准化处理,使得该模型除了能够对文本相关的语音数据进行特征提取之外,也能够对文本无关的语音数据进行很好地特征提取。

[0154] 4) 经过试验后发现,本发明能够更好的处理大规模的语音数据并且能够很好地适应不同语言的处理。

[0155] 在本发明所提供的几个实施例中,应该理解到,所揭露的方法、装置和设备,可以通过其它的方式实现。例如,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。

[0156] 所述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。

[0157] 另外,在本发明各个实施例中的各功能单元可以集成在一个处理单元中,也可以是各个单元单独物理存在,也可以两个或两个以上单元集成在一个单元中。上述集成的单元既可以采用硬件的形式实现,也可以采用硬件加软件功能单元的形式实现。

[0158] 以上所述仅为本发明的较佳实施例而已,并不用以限制本发明,凡在本发明的精神和原则之内,所做的任何修改、等同替换、改进等,均应包含在本发明保护的范围之内。

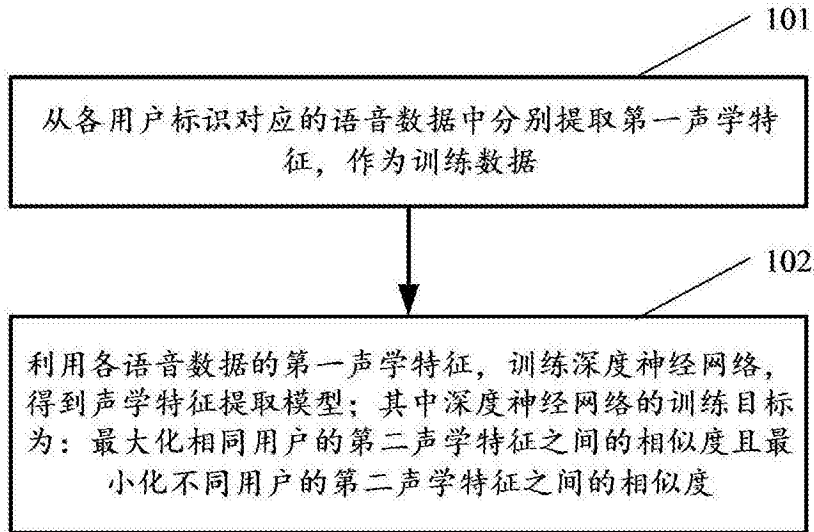


图1

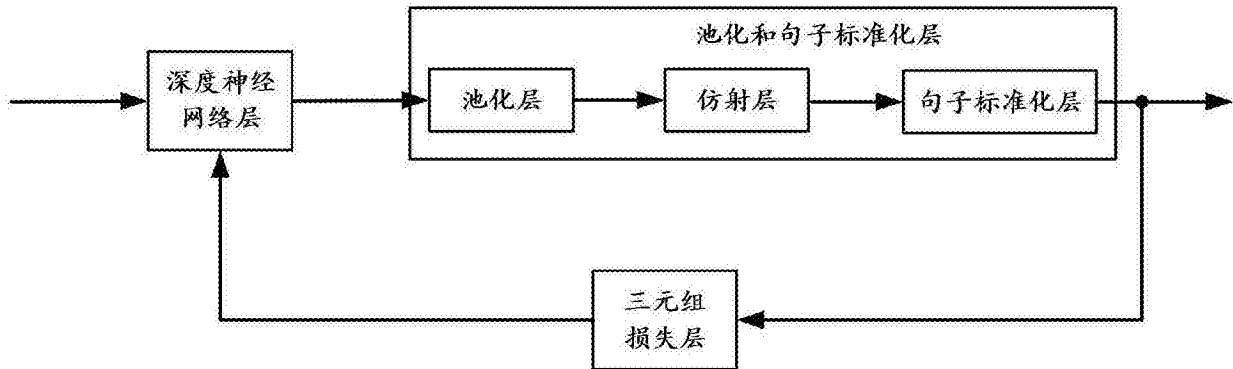


图2

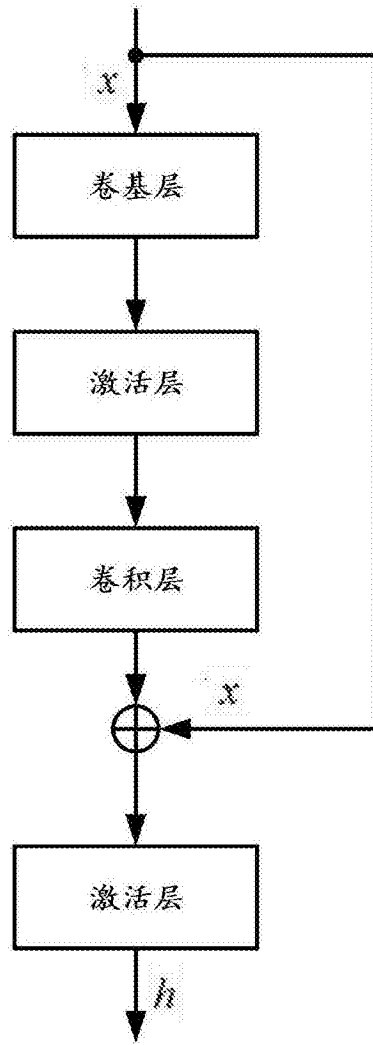


图3

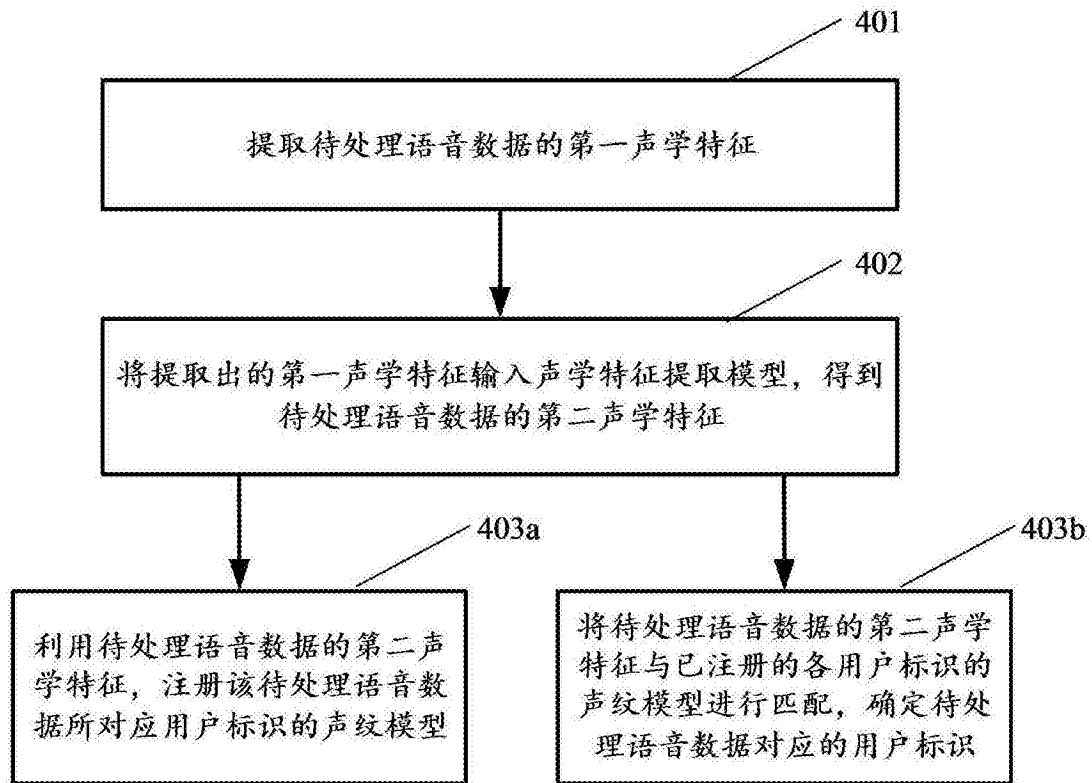


图4

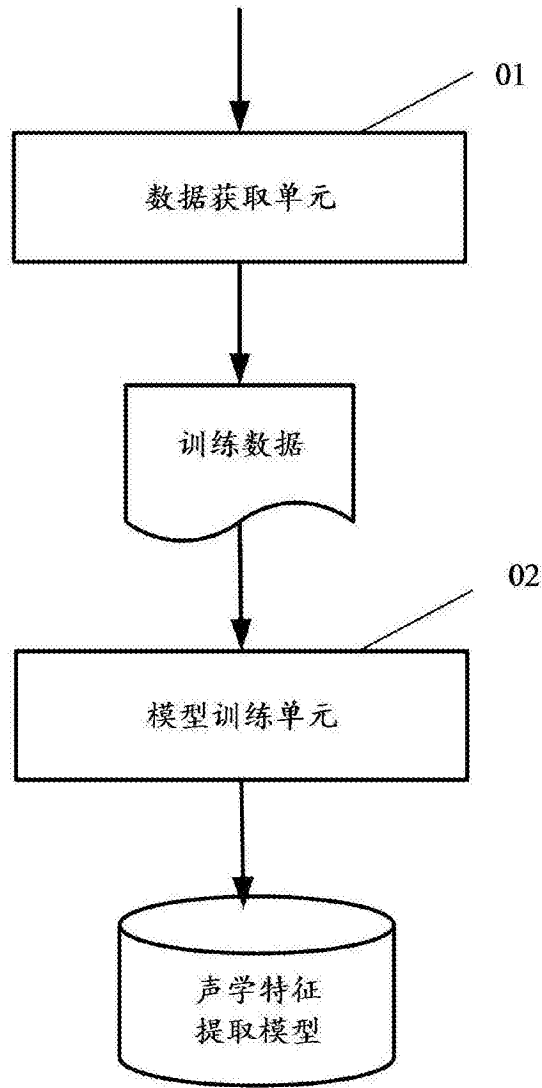


图5

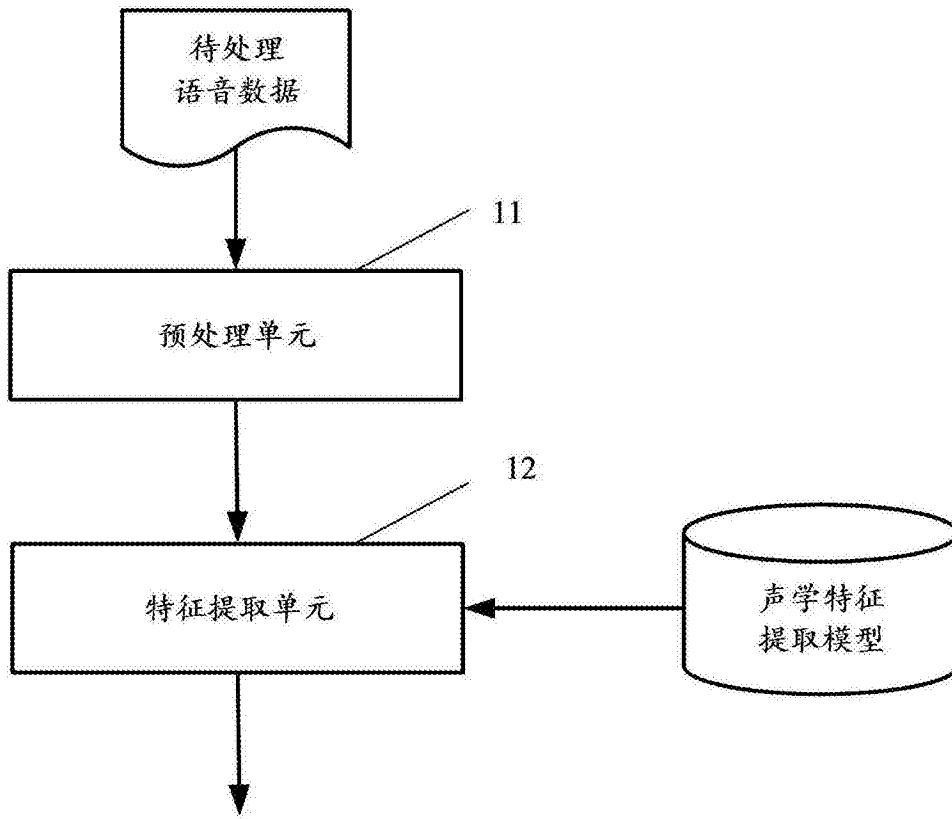


图6

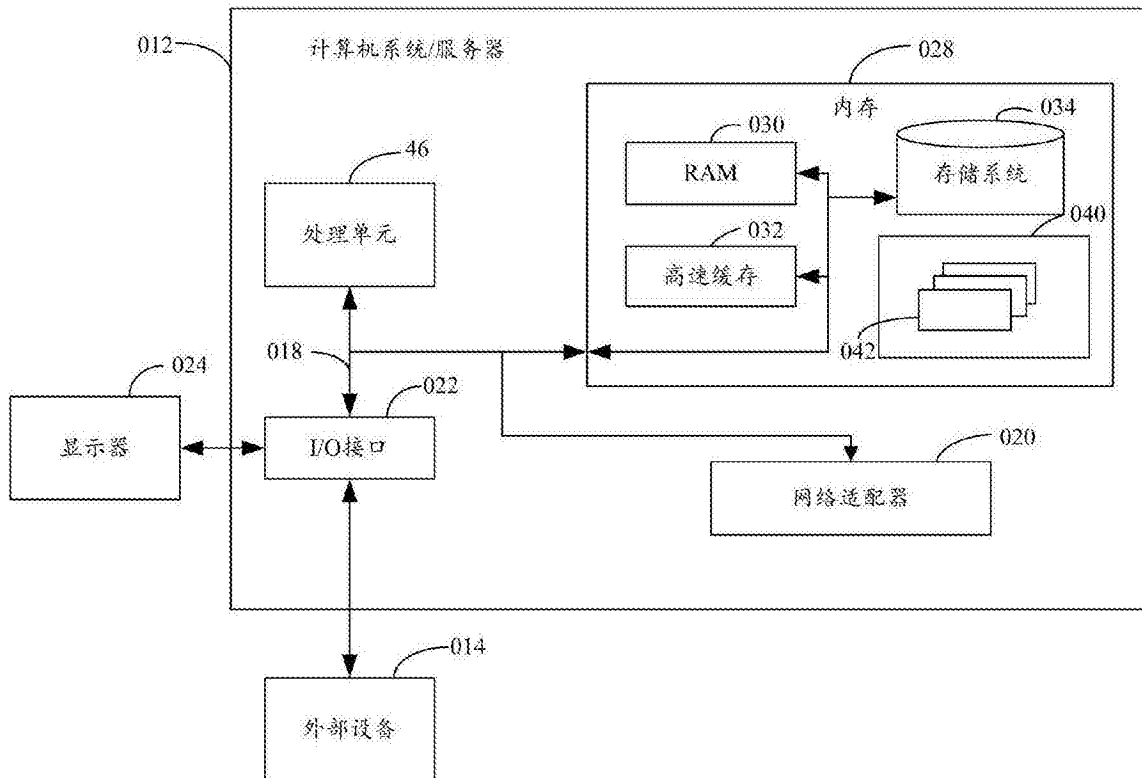


图7