



## (12) 发明专利申请

(10) 申请公布号 CN 104541324 A

(43) 申请公布日 2015. 04. 22

(21) 申请号 201380031695. 3

(74) 专利代理机构 北京银龙知识产权代理有限公司 11243

(22) 申请日 2013. 06. 26

代理人 曾贤伟 周捷

(30) 优先权数据

P. 403724 2013. 05. 01 PL

(51) Int. Cl.

G10L 15/14(2006. 01)

(85) PCT国际申请进入国家阶段日

2014. 12. 17

(86) PCT国际申请的申请数据

PCT/EP2013/063330 2013. 06. 26

(87) PCT国际申请的公布数据

W02014/177232 EN 2014. 11. 06

(71) 申请人 克拉科夫大学

地址 波兰克拉科夫

(72) 发明人 巴尔托什·焦尔科 托马什·贾奇克

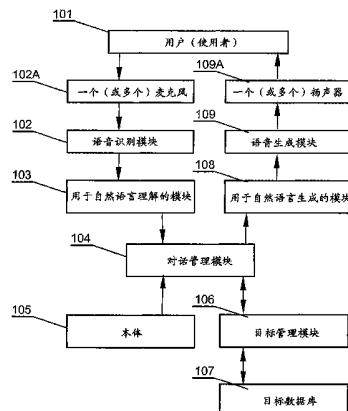
权利要求书1页 说明书8页 附图4页

### (54) 发明名称

一种使用动态贝叶斯网络模型的语音识别系统和方法

### (57) 摘要

一种用于语音识别的计算机实现的方法，包括以下步骤：通过输入设备（102A）的方式，记录（201）表示语音的电信号，并将该信号转换为频域或时-频域（202），基于动态贝叶斯网络在分析模块中分析信号（205），被配置为基于观察到的信号特征（OA, OV）生成单词（W）的假设和它们的概率，基于特定单词（W）假设和它们的概率，识别（209）出表示语音的电信号所对应的文本。该方法的特征在于，将观察到的信号特征（308-312）输入到分析模块（205）中，其中，所述观察到的信号特征是在至少两条并行信号处理线（204a, 204b, 204c, 204d, 201a）上，为频域或时-频域（202）中信号而确定的，其中每条线上的时间片段都不同，以及，在分析模块（205）中对至少两个不同的时间片段分析观察到的信号特征（308-312）之间的关系。



1. 一种用于语音识别的计算机实现的方法,包括以下步骤:

- 通过输入设备(102A),记录(201)表示语音的电信号,并将该信号转换为频域或时-频域(202),

- 基于动态贝叶斯网络在分析模块中分析信号(205),所述分析模块被配置为基于观察到的信号特征(OA, OV)生成单词(W)的假设和它们的概率,

- 基于特定单词(W)假设和它们的概率,识别(209)出表示语音的电信号所对应的文本,

该方法的特征在于,

- 将观察到的信号特征(308-312)输入到分析模块(205)中,该观察到的信号是对于在至少两条并行信号处理线(204a, 204b, 204c, 204d, 201a)上的频域或时-频域(202)中的信号而确定的,其中,每条线上的时间片段都不同,

- 以及在分析模块(205)中对至少两个不同的时间片段分析观察到的信号特征(308-312)之间的关系。

2. 如权利要求1所述的方法,其中时间片段具有预定的持续时长。

3. 如权利要求1或2所述的方法,其中时间片段取决于语音片段的内容,比如音素,音节,单词。

4. 如前述任一权利要求所述的方法,特征在于,在分析模块(205)中,定义描述模型的变量之间的确定性和概率性关系,而至少对观察到的信号特征与当前状态(Sti)的关联定义概率性关系。

5. 如前述任一权利要求所述的方法,其特征在于同时分析(205)不同的观察到的信号特征(OA, OV)。

6. 一种用于语音识别的、计算机实现的系统,包括:

- 用于记录代表语音的电信号的输入设备(102A),

- 用于将表示语音的所记录的电信号转换为频域或时-频域的模块(202),

- 基于动态贝叶斯网络的分析模块(205),被配置为分析表示语音的信号,并且,被配置为基于观察到的信号特征(OA, OV)生成单词(W)的假设和它们的概率,

- 用于基于已定义的单词(W)的假设与它们的概率,识别表示语音的电信号所对应的文本的模块(209),

该系统的特征在于进一步包括:

- 至少两个信号参数化模块(204a, 204b, 204c, 204d, 201a),用于在至少两条并行信号处理线上为分析模块(205)确定至少两个观察到的信号特征(308-312),其中每条线上的时间片段都不同,

- 其中分析模块(205)被配置为对至少两个不同的时间片段分析观察到的信号特征(308-312)之间的相关性。

7. 一种计算机程序,包括当在计算机上运行所述程序时,用于执行根据权利要求1-5的任何一个的计算机实现方法的所有步骤的计算机代码装置。

8. 一计算机可读介质,存储计算机可执行指令,当在计算机上执行时,所述指令执行根据权利要求1-5的任何一个的计算机实现的方法的所有步骤。

## 一种使用动态贝叶斯网络模型的语音识别系统和方法

### 技术领域

[0001] 本发明的目标是实现一种使用贝叶斯网络的语音识别系统和方法。特别地，涉及一种自动语音识别系统，其可以在用于广告和信息意图的对话系统中应用。对话系统的实施可以采用报亭或货摊的形式，其与顾客或观众开始一对对话，并且将呈现适当的多媒体内容。

### 背景技术

[0002] 语音识别系统在日常生活中变得越来越常见。比如，它们可以被用于信息电话中心，比如为公共交通所用。然而，这些系统仍然经常依赖于键盘和文本作为输入信息源，而不是使用语音作为输入信息源而运行。

[0003] 已知各种类型的计算机化的交互报亭被用于与用户进行对话。比如，美国专利 US6256046 公开了一种在计算机化的报亭内的有源公共用户交互接口，其通过处理视觉数据、通过使用动作和色彩分析以检测表示用户出现的环境中的改变来感知用户。交互空间被定义，系统记录其环境的初始模型，该环境随着时间更新，以反映出不活动对象的添加或减去，并且补偿光的改变。该系统研发了针对移动对象的模型，因此当他们在交互空间的附近移动时，该系统能够跟踪用户。一立体摄像系统进一步增强了该系统感知位置和移动的能力。该报亭呈现出音频和视频的反馈来反映其“看到”了什么。

[0004] 美国专利申请 US20080204450 公开了一种用于提供虚拟世界的系统、方法和程序产品，其中主动提供的广告被嵌入在自动虚拟角色中。所提供的系统包括：用于将广告虚拟角色引入虚拟世界的注册系统；用于定向用户虚拟角色以实现广告虚拟角色所传递的广告内容的定向系统；用于定义广告虚拟角色如何在虚拟世界中移动的移动系统；以及用于定义广告虚拟角色如何将广告内容传递给用户虚拟角色的广告传递系统。

[0005] 诸如上述的已知的对话系统的缺陷包括，在与用户进行复杂对话时缺乏足够的语音识别能力。

[0006] 美国专利 US7203368 公开了一种模式识别程序，其使用 HMM（隐马尔科夫模型）和 CHMM（耦合隐马尔科夫模型）形成了分级的统计模型。分级的统计模型支持具有多个超节点的父层和具有与每一个父层的超节点相关联的多个节点的子层。经过训练之后，分级统计模型使用从数据集中提取的观察矢量来寻找基本的最优状态序列片段。对该过程进行改进是很有利的。

[0007] 一个比基于 HMM 的方案少一些限制的、更加通用的方案，是将贝叶斯网络用于语音识别。使用贝叶斯网络的方案包括动态贝叶斯网络（DBN），已经在以下出版物中被公开：

-M. Wester, J. Frankel, 以及 S. King 所著的：“Asynchronous articulatory feature recognition using dynamic Bayesian networks” (Proceedings of IEICI Beyond HMM Workshop, 2004) (“使用动态贝叶斯网络的异步分节特征识别”，公开于 2004 年 HMM 研讨会的 IEICI 会议录)；

-J. A. Bilmes 和 C. Bartels 所著的“Graphical model architectures for speech

recognition" , IEEE Signal Processing Magazine, vol. 22, pp. 89–100, 2005( "用于语音识别的图形模型构造", 公开于 IEEE 信号处理杂志, 2005 年, vol. 22, pp. 89–100) ;

-J. Frankel, M. Wester 和 S. King 所著的 " Articulator/feature recognition using dynamic Bayesian networks" , Computer speech and Language, vol. 21, no. 4, pp. 620–640, October 2007( "使用动态贝叶斯网络的发音器 / 特征识别", 公开于 2007 年 10 月, 计算机语音和语言 vol. 21, no. 4, pp. 620–640)。

使用贝叶斯网络的语音识别方法依据特征矢量对声音时长进行建模。在 DBN 中, 使用表示声音的变量替换表示时长的变量已经变得可能。然而, 所有的现有技术的方案都在预定的时间范围内进行语音分析。

[0008] 考虑到之前的现有技术, 有必要设计和实现一种允许提高人类和机器之间的对话效率的语音识别系统和方法。

## 发明内容

[0009] 本发明的目的在于提供一种用于自动语音识别的计算机实现的方法, 包括以下步骤: 通过输入设备记录表示语音的电信号, 并将该信号转换至频域或时 - 频域, 基于 DBN 在模块分析中分析信号, 被配置为基于观察到的信号特征 (OA, OV) 生成单词 (W) 的假设和它们的概率, 基于特定单词 (W) 假设和它们的概率识别出表示语音的电信号所对应的文本。该方法的特征在于, 将观察到的信号特征输入到分析模块中, 该观察到的信号是对于多个时间段、在至少两条并行信号处理线上的频域或时 - 频域中为信号而确定的, 其中在每条线上的时间片段都不同, 并且, 在分析模块中对至少两个不同的时间片段分析观察到的信号特征之间的关系。

[0010] 优选地, 时间片段具有预定的时长。

[0011] 优选地, 时间片段取决于语音片段的内容, 比如音素 (phonemes)、音节 (syllables)、单词 (words)。

[0012] 优选地, 该方法进一步包括在分析模块定义描述模型的变量之间的确定性和概率性关系, 而概率性关系至少被定义用于将观察到的信号特征与当前状态 (Sti) 进行关联。

[0013] 优选地, 该方法进一步包括同时分析不同的观察到的信号特征 (OA, OV)。

[0014] 本发明的另一个目标是实现用于语音识别的、计算机实现的系统, 包括用于将代表语音的电信号进行记录的输入设备, 用于将表示语音的记录的电信号转换为频域或时 - 频域的模块, 基于 DBN 的分析模块, 被配置为分析表示语音的信号, 并且, 被配置为基于观察到的信号特征 (OA, OV) 生成单词 (W) 的假设和它们的概率, 用于基于已定义的单词 (W) 的假设与它们的概率识别表示语音的电信号所对应的文本的模块。该系统进一步包括至少两个信号参数化模块, 用于为分析模块在至少两条并行信号处理线上为每条线上不同的时间片段确定至少两个观察到的信号特征, 其中分析模块被配置为分析在至少两个不同的时间片段上, 观察到的信号特征之间的相关性。

[0015] 本发明的目标是还提供一种计算机程序, 包括当所述程序在计算机上运行时, 用于执行根据本发明的计算机实现的方法的所有步骤的程序代码装置, 还有存储计算机可执行指令的计算机可读介质, 当在计算机上执行该指令时, 该指令执行根据本发明的计算机实现的方法的所有步骤。

### 附图简要说明

- [0016] 已经在附图中的示例性实施例中公开了本发明的目标，其中：
- [0017] 附图 1 示出了依据本发明的系统的方框图；
- [0018] 附图 2 示出了自动语音识别过程的方框图；
- [0019] 附图 3 示出了在不同长度的并行时间周期内使用 DBN 对语音进行模型化；
- [0020] 附图 4 描述了使用与附图 3 中示出的 DBN 相似的 DBN 进行单词序列解码的例子（为了示例性目的，已经被简化的版本）。

### 具体实施方式

[0021] 附图 1 示出了依据本发明的系统的方框图。该系统可以被用于交互性广告或其它提供信息的对话系统中。对话尽可能地接近现实中的对话。由于使用诸如模式识别、语义分析的技术，使用语音合成所伴随的本体认知和自然语言生成，这种假设可以被实现。

[0022] 在可能使用本发明的对话系统中，可以包括高质量显示器或图像投影器。在优选的实施例中，对话系统还配备有用户出现检测器，或者在更加高级的例子中，配备有诸如生物检测器、面部识别模块等的用户特征检测器。对话系统还可以包括方向性麦克风，以能够更加有效地采集语音。

[0023] 输出信息被调整为对话的情景并且根据用户爱好而被确定。对话系统优选地还输出与用户对话的视觉虚拟角色或用户的图像。

[0024] 采用语音识别的对话系统与用户或多个用户 101 交互地交谈。用户 101 通过说话将问题输入至声音输入模块，比如输入至麦克风 102A 中。被麦克风记录的声音通过语音识别模块 102 进行处理，之后传送至用于识别自然语言的模块 103。

[0025] 用于理解的模块 103 负责在预期的响应情景中解释用户 101 的状态识别假设，其采用使得它们可以被机器所理解并且可以容易地并且快速地被处理的方式。比如，若在旅游信息现场实施该系统，基于伴随着语音假设概率的语音假设列表的用于理解的模块 103 需要确定说话者是否在寻找一个具体的地点，如果是的话，这个地点是什么，或者服务、公共交通运行的时间信息等。在最简单的版本中，为了实现该意图，该模块使用关键词，但是它们也可以使用基于句法模型（例如，句子解析器）和 / 或语义（比如，单词网或语音 HMM）的更加高级的方案，这些内容在 D. Jurafsky, J. H. Martin, "Speech and Language Processing", Second Edition, Pearson Education, Prentice Hall, 2009 中被公开。

[0026] 被自然语言理解模块 103 处理之后，句子或句子假设被输入至对话管理模块 104（比如，诸如在 D. Jurafsky, J. H. Martin, "speech and Language Processing", Second Edition, Pearson Education, Prentice Hall, 2009 中所描述的），其通过适当地询问本体模块 105，与目标管理模块 106 和目标数据库 107 合作，决定将要对用户询问呈现的响应。

[0027] 本体模块 105 包括关于世界的有规则的知识，比如，在某些特定类型中可以购买哪些产品、人们将与所选择的产品一起购买哪种产品等等的信息。本体模块可以包括来自社会服务的额外的不同类型的数据，例如，核实正在与用户进行对话的朋友是否在该用户访问的城市中等等。本体模块还可以包括任何其它的实际知识，使用一种计算机或其它机器可以处理它的方式将其系统化。

[0028] 目标管理模块 106 被用于在计算机中实现商业、广告、谈判等的已知规则，其用于指引专业人员（比如，商业雇员），其职责被依据本发明的系统所执行。

[0029] 确定响应的内容之后，自然语言中的响应在生成自然语言的模块 108 中产生，紧接着被送至语音生成模块 109 中。通过安装在系统内的扬声器或其它输出设备 109A 将以语音形式生成的响应输出至用户 101。

[0030] 本发明中所使用的关键元件是使用贝叶斯网络进行分析的计算机实现的模块。贝叶斯网络能够对复杂的现象进行模型化，其中分开的元素可以彼此依赖。将基本模型产生为有向的、非循环的图形，其中节点表示模型的分开的元素（随机变量），而边表示这些元素之间的相关性。

[0031] 附加地，边被分配有概率值，该概率值指定在另一个事件假定了一特定值的情况下，发生多个事件中的一个事件。通过使用贝叶斯法则，复杂的条件概率可以在贝叶斯网络的特定路径中计算。这些概率可以被用于推断网络中的各个元素将要采用的值。

[0032] 每一个网络变量必须有条件地独立于不和它连接的其它变量。使用这种方式生成的图形可以被解释为事件的紧密表示 (compact representation)，考虑到图形的节点之间的条件独立性，这些事件和假定发生的累积概率。

[0033] DBN 可被用于语音识别。之后，节点不仅仅表示单个随机变量，还表示变量序列。这些可以被解释成时间序列，其允许根据时间的经过来对语音模型化。因此，多个连续的观察状态给出对于到达最终状态的明确路径的判断。

[0034] 使用标准贝叶斯网络基于声音的期望时长，依赖于发音者特征的矢量。网络具有用于每个特征的单个离散变量和用于声音时长的单个连续变量。网络描述了特征之间的关系。表示特征的节点的值依赖于输入到网络中的值，并且可选地，依赖于其它特征。表示时长的节点值是隐层（如在 HMM 中一样），仅仅直接依赖于从其它节点接收的值。

[0035] 引入 DBN 允许用代表声音的变量代替代表时长的变量。具有特征之间的关系的整个网络使用这样的方式被复制，这种方式为：多个网络中的一个网络表示在时刻 t-1 处所分析的信号，并且下一个网络表示在时刻 t 所分析的信号。这两个网络通过边相连，这些边具有可能随着时间变化的状态之间的转移的概率值。

[0036] 应注意到，本发明并不只限制于具有两个子网络的实施例。可以有更多的子网络，每一个子网络用于后续的时刻。典型地，可以有数百或数千个网络。可以将这种结构多次复制至多个后续的时刻。附加地，在某些情况下，这种局部贝叶斯网络结构可以在离散时间之间修改其本身。

[0037] DBN 模型还可以被用于将来自不同源的信号的相关信息联系起来，比如声学特征和视觉特征（比如嘴唇运动）。这种类型的系统对于在恶劣的声学条件的地方的应用特别有用。在例如街道、飞机场、工厂等地方，低信噪比 (SNR) 值使得仅仅使用源自声学路径的信息极大地降低所获取的结果的质量。增加从对相同类型的噪声不敏感的其它信号类型获取的信息，就会消除这些困难，并且使得在这些地方也可以使用语音识别系统。

[0038] 发明人已经注意到，当用于语音分析时，与 HMM 方法相比，贝叶斯网络具有更少的限制条件。

[0039] 附图 2 示出了语音识别过程的方框图。下面的描述也会参考附图 3 中的某些特征，附图 3 示出了在不同长度的时间周期上使用 DBN 对语音进行模型化。

[0040] 此处通过这样的方法使用 DBN 对语音进行模型化, 即分开的观察表示不同的时长——如图 3 中所示。这些不同的时长可以是预定长度的片段, 比如 5ms, 20ms, 60ms, 这依赖于语音片段的内容, 比如音素、音节、单词或者两种类型的结合, 比如, 5ms, 20ms, 音素, 单词。

[0041] 所呈现的方法允许提取不同的信息类型, 并且对所获取的特征进行直接融合, 这是由于使用 DBN 模型对状态概率进行了估计 (附图 3 中的 St1 至 St6)。

[0042] 基于那些描述模型的变量之间的两种类型的关系, 在 DBN 中进行推断: 确定性关系 (在图 3 中用直箭头表示) 和概率性关系 (在图 3 中用波浪形箭头表示)。

[0043] 基于那些公知事实定义确定性关系, 比如, 当分析一个给定单词  $W_{ti}$  时, 已知第一个音素  $P_{ti}$  的类型和位置  $W_{ps}$ 。之后, 通过知晓是否已经发生了从该音素到下一个音素的转移  $P_{tr}$ , 可以确定当前音素在单词中的位置: 若还没有出现音素转移, 那么  $t+1$  时刻的  $W_{ps}$  等于  $t$  时刻的  $W_{ps}$ , 若已经观察到了前述的转移, 那么,  $t+1$  时刻的  $W_{ps}$  等于  $W_{ps}+1$ 。

[0044] 可以使用相似的方式获取到从一个单词到另一个单词的转移  $W_{tr}$  有关的信息。出现从用音标标出的单词的最后一个音素进行的转移, 暗示着有必要对另一个单词  $W_{ti}$  进行分析。

[0045] 另一种类型的关系是概率性关系。为了基于其之间存在概率性关系的变量进行推断, 有必要确定限定发生这些事件的概率 (概率密度函数—PDF) 的函数。该类型的关系用于将信号的观察到的特征和当前状态  $St_i$  进行关联。优选的 PDF 函数是高斯混合模型—GMM。

[0046] 某些关系可以既是确定性也是概率性的, 比如连续单词  $W_{ti}$ 。一旦没有发生从一个单词到另一个单词的转移, 该关系就是确定性的——该单词与  $t-1$  时刻的单词一样。一旦发生转移, 使用从语言模型中得到的知识, 以概率性方式确定下一个单词  $W_{ti+1}$ 。

[0047] 基于声学特征的观察值在 DBN 中进行推断。然而, 任何观察都可能受到测量误差 ( $s$ ) 的影响。引入属于相同组 (图 3 中的 OA11, OA23 和 OA33 或者 OV11 和 OV23) 的、相关的、随时间变化的观察量之间的概率关系同样可以减小这种误差。

[0048] 状态  $S_{ti}$  和之前的状态  $St_{i-1}$  用于估计观察量是说出给定音素 (附图 3 中的 Pt1 至 Pt6) 的结果的概率。

[0049] 给定音素的出现还和转移状态  $P_{tr}$  在一定概率上相关。音素  $P_{ti}$ 、音素转移率  $P_{tr}$ 、音素在单词  $W_{ps}$  中的位置和从单词  $W_{tr}$  进行的转移率允许估计假设记录的声音中包含单词  $W$  的正确性。

[0050] 语音具有某些频率特征和能量特征在短时期内几乎保持恒定的特性。然而, 在长时期内它们显著变化。然而, 并没有定义第一和第二种情况出现时的特殊时刻, 因而, 使用 DBN 模型非常有利。不同片段中的观察量之间的关系可以呈现, 也可以不呈现。

[0051] 比如, 对于具有四个时间长度的配置的变形而言, 它们假定 5ms、20ms、音素和单词用于并行分析。可能存在不同的模型配置, 例如, 存在全部四个范围之间的关系, 还可能仅仅在 5ms 和 20ms 层与音素层之间存在关系, 仅仅在 20ms 层和音素层之间存在关系, 或者仅仅在音素层和单词层之间存在关系。

[0052] 此外, 每一个范围可以具有一些观察量类型, 其与不同种类的语音特征相关。比如, 它们中的一个可以是频率特征矢量, 另一个可以是能量, 在一个可以是视觉特征矢量。这些同样可以是相同种类的声学特征, 但是通过不同的方法获得 (比如 WFT (小波傅里叶变

换)、MFCC(梅尔倒谱系数)),还可以是在不同的时间范围内使用同样的方法所获取的声学特征,比如在 20ms 的滑动窗中、或是在 50ms 的滑动窗中,两者均是每隔 10ms 提取一次。

[0053] 另外,某些范围可以仅仅在特别类型的特征分析中出现,并且,在其它类型中并不出现(附图 3——声学特征 1 观察量 (308) 持续 60ms、声学特征 2 观察量 (310) 持续 20ms、以及视觉特征 1 观察量 (309) 持续 30ms)。

[0054] 可能存在在分析过程中使用(也是同时使用)的更多类型的描述信号的特征,比如,声音的有声/无声描述、共振峰频率或基音频率。

[0055] 附图 2 中公开的方法开始于步骤 201,获取语音信号。下一个步骤 202 通过诸如 WFT 的方法或使用短时傅里叶变换(STFT)进行的时频变换将信号处理成频率域。可以应用其它变换以允许对包括在不同时刻上的不同频率子带内的信息(比如信号能量)进行定量描述。

[0056] 接着,在步骤 203,时-频谱被分割成恒定的帧,比如 5ms、20ms、60ms 等等,或者依据预定的算法进行分割,诸如,在例如如下中所公开的:

-P. Cardinal, G. Boulian 和 M. Comeau 所著的 "Segmentation of recordings based on partial transcriptions", Proceedings of Interspeech, pp. 3345-3348, 2005; 或者

-K. Demuynck 和 T. Laureys 所著的 "A comparison of different approaches to automatic speech segmentation", Proceedings of the 5th International Conference on Text, Speech and Dialogue, pp. 277-284, 2002; 或者

-Subramanya, J. Bilmes 和 C. P. Chen 所著的 "Focused word segmentation for ASR", Proceedings of Interspeech 2005, PP. 393-396, 2005.

[0057] 分割模块 (203) 将频谱分析过程分成多个线程,独立地对其进行参数化。

[0058] 线程数目可以不是前面所述的 4 个。附图 2 中的例子使用四个分开的线程,其具有 5ms 的帧——204a、20ms 的帧——204b、音素——204c 和单词——204d,而在方框 204a 至 204d 中从每个线程中提取的特征,表示了在特定时刻的语音。这些参数化方框使用诸如 MFCC、感知线性预测 (PLP) 的处理算法,或者是其它的,比如:

-H. Misra, S. Ikbal, H. Bourlard 和 H. Hermansky 所著的 "Spectral entropy based feature for robust ASR", Proceedings of ICASSP, pp. 1-193-196, 2004; 和 / 或

-L. Deng, J. Wu, J. Droppo 和 A. Acero 所著的 "Analysis and comparison of two speech feature extraction/compensation algorithms", IEEE Signal Processing Letters, vol. 12, no. 6, pp. 477-480, 2005; 和 / 或

-D. Zhu 和 K. K. Paliwal 所著的 "Product of power spectrum and group delay function for speech recognition", Proceedings of ICASSP, pp. 1-125-128, 2004.

[0059] 从模块 204a-204d 获取的特征与诸如信号能量和视觉特征矢量的观察量 201a 一起被传至 DBN205。DBN 模型,使用对 BN 的近似推测的其自己的嵌入的算法(比如变分消息传播 (variational message passing)、期望传播 (Expectation Propagation) 和 / 或吉布斯采样 (Gibbs Sampling)),利用语音识别中所使用的动态编程算法(比如维特比解码和 / 或 Baum-Welch),并基于词典 206 和语言模型 207 的内容(比如单词的二元语法)来确定单词假设并且计算它们的概率。在大多数情况下,这些假设会部分重叠,因为 DBN 可以在相同的时间长度上呈现不同的假设。在进一步的语言模型 208 中(优选地,比 DBN 中使用的第

一个语言模型更加高级)可以继续处理这些假设,以便获取识别的语音文本 209。

[0060] 附图 3 公开了示例性 DBN 结构。术语 W301 表示单词, Wtr302 表示单词转移率, Wps303 表示在特定单词中的音素的位置, Ptr304 表示音素转移率, Pt 305 表示音素, Spt 306 表示之前的状态, S 307 表示状态, OA1 308 表示在 60ms 的时间窗中, 观察到的第一种类型的声学特征, OV1 309 表示在 30ms 时间窗中观察到的第一种类型的视觉特征, OA2 310 表示在 20ms 时间窗中观察到的第二种类型的声学特征, OA3 311 表示在 10ms 时间窗中观察到的第三种类型的声学特征, 同时, OV2 312 表示在 10ms 时间窗中观察到的第二种类型的视觉特征。

[0061] 箭头表示变量之间的相关度(依赖度),如前面所描述的一样。通过条件概率分布(CPD) 定义转移率,其是基于训练数据,在贝叶斯网络的训练过程中被计算出的。

[0062] 附图 4 描述了使用附图 3 中所示的 DBN 用于解码单词序列的一个例子。与图 3 中的区别在于,为了实现语音识别,对于不同长度的两帧,这里使用了一种类型的信号声学特征。网络呈现出解码短语:“Cat is black”的过程——语音转换是:/kæt iz blæk/。音素状态依赖于两种类型的观察量 01 和 02。时刻 t 处的前一状态 306 是时刻 t-1 处的状态 307 的完美复制。依赖于在单词 303 中的当前位置、音素转移至另一个音素 304 的出现率、音素的状态 306 和前一状态 307,对单词 301 的后续音素进行分析。若转移概率值大于 0.5,则出现音素转移。附图 3 中的贝叶斯网络的分开的节点的符号被这些状态的值取代。对于 302 和 304 而言,它们是这样的值:T(真)/F(假) 分别表示后续单词或后续音素之间出现转移或未出现转移。对于音素在单词 303 中的位置而言,它是当前所分析的音素的索引(1-3 指代单词‘cat’、1-2 指代单词‘is’、1-4 指代单词‘black’)。只有在时刻 t-1 的前述时刻中,获取到音素 304 的转移值为‘T’时,才会出现音素索引的变化。此外,只有当在特定单词中的上一个索引的音素转移 304 的时刻所获得的单词转移 302 发生的情况下,单词 301 才会变化。在这种情况下,作为使用语言模型的结果,连续的单词之间的关系从确定性改变成概率性。在附图上部的表中示出二元语法语言模型(应用多对单词的模型)的示例性值。此外,还公开了示例性的语言模型中初始单词概率值。通过同时处理不同时长内的片段和各种类型的特征所达到的技术效果是提高了语音识别质量,因为在一种类型的时间片段下会更好地识别以各种方式说出的一种类型的音素,并且其它的需要不同类型的片段,但是,为每一种音素类型确定合适的分析时间窗是复杂的。此外,考虑到在更局部的时间片段上精确地提取信息,某些特征表示稳定特性,同时其它的特征需要更全局的时间片段。使用如图 3 所示的结构,可以一次提取两种类型的特征。在传统的系统中,只是通过局部特征或只是通过全局特征来提取所用的一些信息。此外,比如,视觉特征可以具有与声学特征不同的持续时间,即,比如,对要说出声音的嘴唇的观察量可以持续长于或短于特定声音的时间。

[0063] 本领域技术人员能够容易地认识到,前述的语音识别方法可以被一个或多个计算机程序执行或控制。这些计算机程序典型地使用诸如个人电脑、个人数字助理、移动电话、数字电视的接收器和解码器、信息亭或类似的计算设备中的计算资源来被执行。应用被存储于非易失性存储器中,比如闪存,或者被存储于易失性存储器中,比如 RAM,并且通过处理器执行应用。这些存储器是用于存储计算机程序的示例性存储介质,所述计算机程序包括执行依据本文所公开的技术概念的计算机执行的方法的所有步骤的计算机指令。

[0064] 尽管已经参考特定的优选实施例限定、描述和描写了本文公开的发明,在前述具

体方式中的这些参考和实施的例子并不对本发明作任何限制。然而，在不违背技术概念的更广范围的条件下，很显然可以对本发明做出各种修改和改变。所公开的优选实施例仅仅是示例性的，并非本文公开的技术概念的全部范围。

[0065] 因此，保护范围不被限定为说明书中描述的优选实施例所，而仅仅通过随后的权利要求所限定。

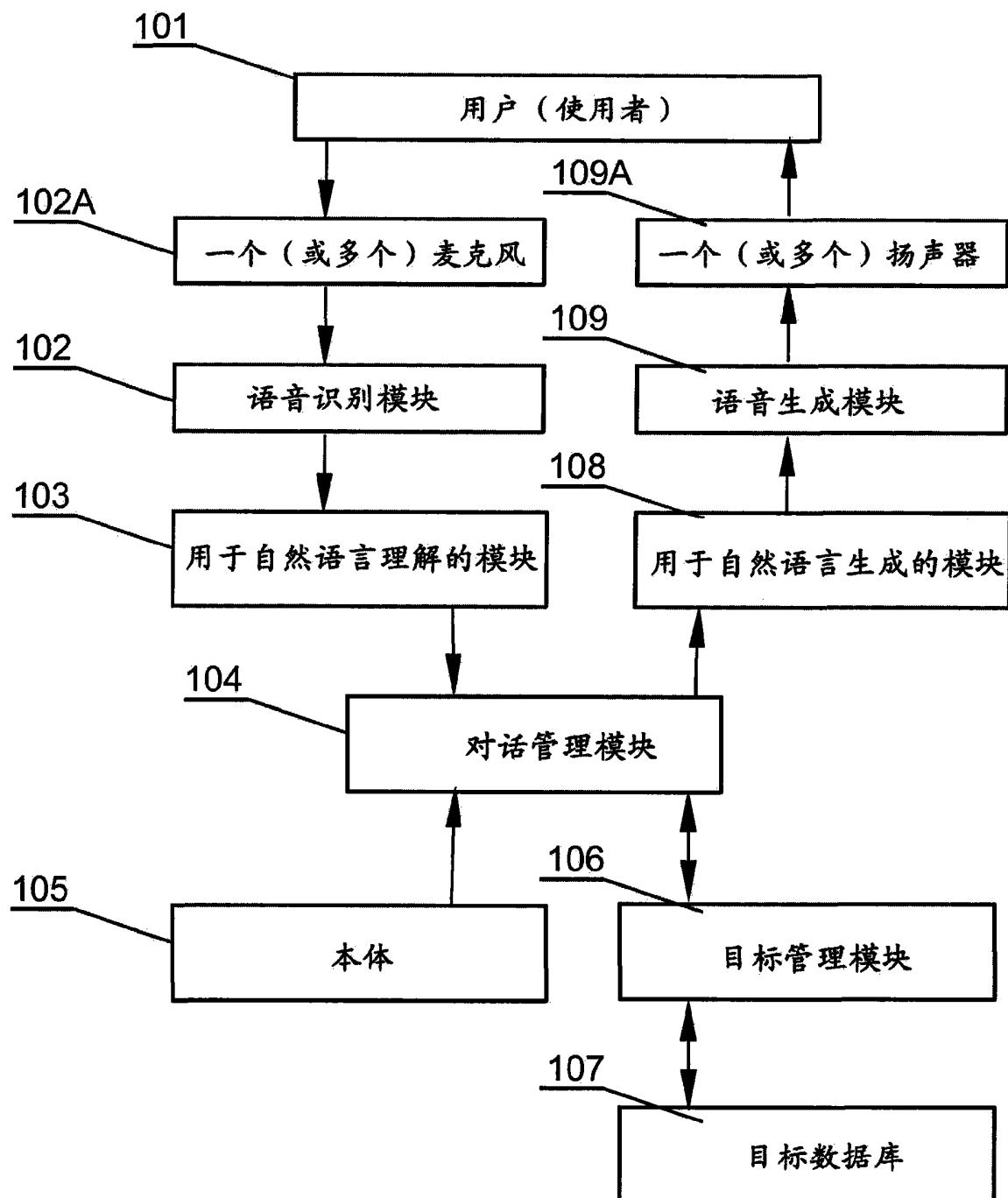


图 1

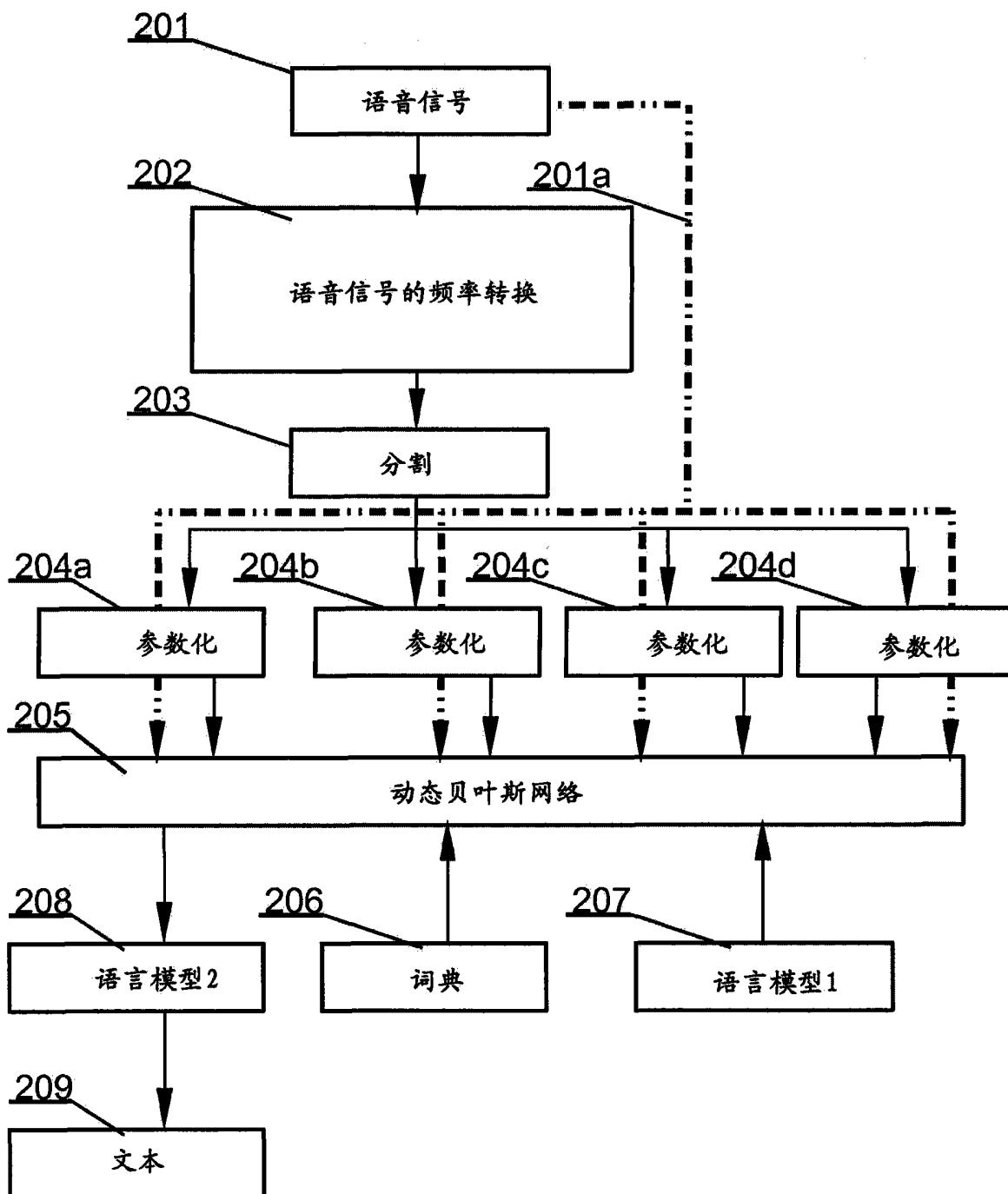


图 2

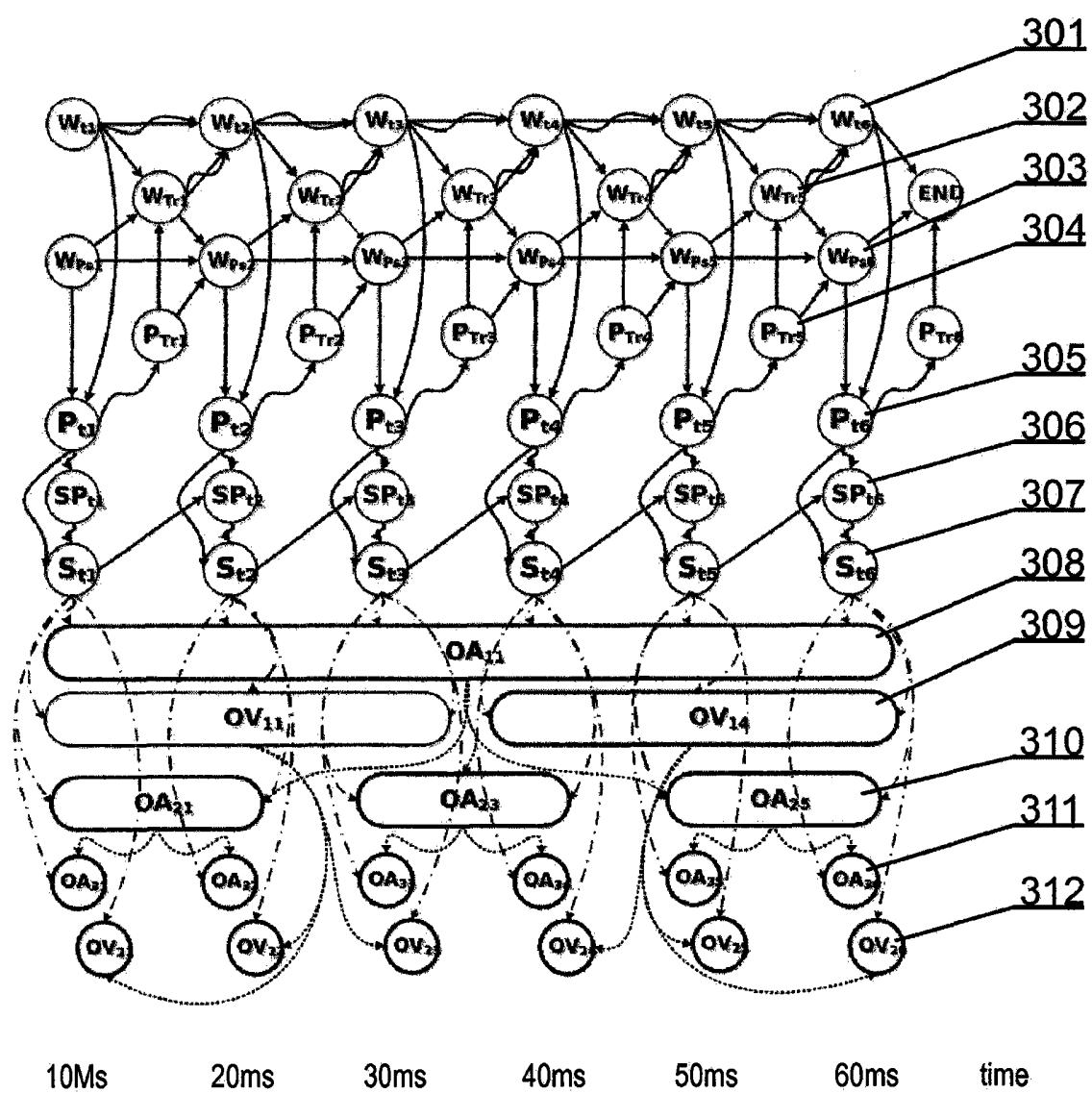


图 3

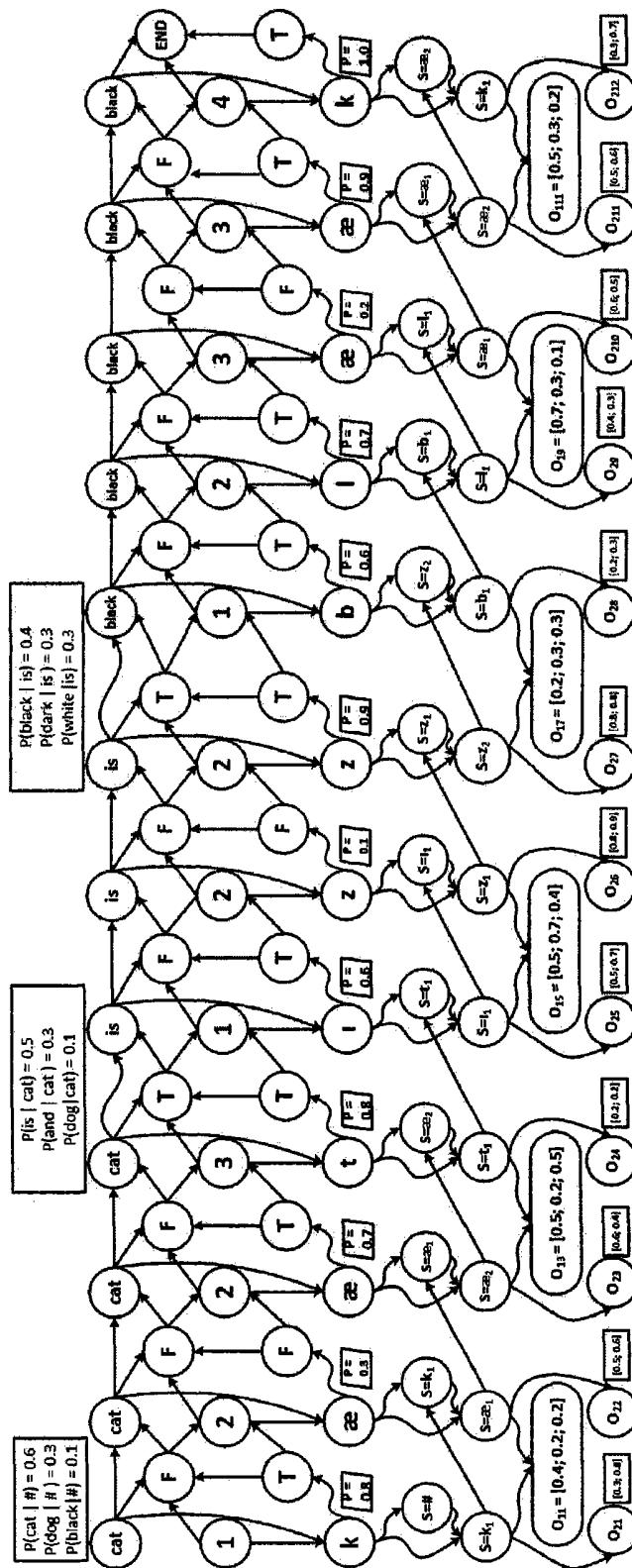


图 4