

[19] 中华人民共和国国家知识产权局

[51] Int. Cl.

G10L 15/20 (2006.01)

G10L 15/26 (2006.01)



[12] 发明专利申请公开说明书

[21] 申请号 200380106508.X

[43] 公开日 2006年1月25日

[11] 公开号 CN 1726532A

[22] 申请日 2003.10.31

[21] 申请号 200380106508.X

[30] 优先权

[32] 2002.12.20 [33] EP [31] 02102875.8

[86] 国际申请 PCT/EP2003/012168 2003.10.31

[87] 国际公布 WO2004/057574 英 2004.7.8

[85] 进入国家阶段日期 2005.6.17

[71] 申请人 国际商业机器公司

地址 美国纽约

[72] 发明人 沃尔克·菲舍尔

谢格弗里德·昆兹曼

[74] 专利代理机构 中国国际贸易促进委员会专利商
标事务所
代理人 康建忠

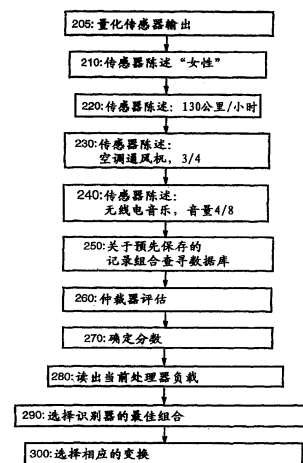
权利要求书 2 页 说明书 11 页 附图 2 页

[54] 发明名称

基于传感器的语音识别器选择、自适应和组合

[57] 摘要

本发明涉及一种操作语音识别系统的方法和相应的系统，其中多个识别器程序是可访问的，以便被激活进行语音识别，并且根据需要被组合，以便有效地改进单个识别器完成的语音识别的结果。为了适应各种工作环境的动态变化的声学条件，以及只具有有限的可用计算能力的嵌入式系统，提出用传感器装置收集(210、220、230、240)表征语音识别边界条件的选择基础数据，例如讲话人和环境噪声等，b)利用(260)程序控制的仲裁装置评估收集的数据，例如包括软件机构和物理传感器的判定引擎，从而从多个可用的识别器中选择(290)最适合的识别器或其组合。



1、一种操作语音识别系统的方法，其中程序控制的识别器(1)执行下述步骤：

- 5 把语音信号分成多帧，并计算每帧的任意类型的特征向量，
用字符或字符组标记所述帧，每个音素产生多个标记，
根据预定的声学模型对所述标记解码，构成一个或多个单词或者
一个单词的多个片段，

10 在所述方法中，多个识别器是可访问的，以便被激活进行语音识别，并且被组合以平衡由单个语音识别器进行的语音识别的结果，其特征在于下述步骤：

- a)用传感器装置(5)收集(210、220、230、240)表征语音识别边界条件的选择基础数据，
b)利用(260)程序控制的判优装置(6)评估收集的数据，
15 c)根据所述评估，从多个可用的识别器中选择(290)最适合的识别器或其组合。

2、按照权利要求 1 所述的方法，其中所述传感器装置(5)是下述一个或多个：

判定逻辑，包括软件程序，物理传感器或者它们的组合。

20 3、按照权利要求 1 所述的方法，还包括下述步骤：

- a)在实现下述一个或多个的判定逻辑中处理(260)物理传感器(5)输出：统计检验，判定树，模糊隶属关系函数，
b)从所述处理返回(270)将用在传感器选择/组合判定中的置信度值。

25 4、按照权利要求 1 所述的方法，其中导致识别器选择判定的所述选择基础数据被保存在数据库中以便反复快速访问(250)，从而获得识别器的快速选择。

5、按照权利要求 1 所述的方法，还包括下述步骤：

根据(280)当前的处理器负载，选择(290)识别器的数目和/或组合。

6、按照权利要求 1 所述的方法，还包括下述步骤：

保存一个识别模型如何被转换成另一识别模型的变换规则(7)，而不是保存多个模型本身。

7、一种具有执行根据前述权利要求 1-6 之一所述的方法的步骤
5 的装置的计算机系统。

8、一种在数据处理系统中执行的计算机程序，包括当在计算机上执行时，完成根据前述权利要求 1-6 任意之一所述的方法的相应步骤的计算机程序代码部分。

9、一种保存在计算机可用介质上的计算机程序产品，包括当所
10 述计算机程序产品在计算机上执行时，使计算机执行根据权利要求 1-6 任意之一所述的方法的计算机可读程序单元。

基于传感器的语音识别器选择、自适应和组合

5 技术领域

本发明涉及计算机化语音识别的领域。

背景技术

特别地，本发明涉及操作大词汇量语音识别系统的方法，其中程
10 控识别器执行下述步骤：

1. 把语音信号分解成长度不必相等的短的时间间隔，即帧，得到每帧的抽取的特征向量，例如包括谱系数，

2. 用字符或字符组标记帧，每帧产生多个标记，

3. 对所述标记解码，从而构成一个或多个单词或者一个单词的
15 多个片段，

4. 在该方法中，多个识别器是可访问的，以便被激活进行语音识别，并且所述多个识别器基于请求被组合，以便改进单个识别器的语音识别结果。

更特别地，上述这种连续语音识别器通过把依赖于上下文的子字
20 单元，比如音子或三音子模拟成基本的隐马可夫模型(也称为“HMM”)，捕捉语声的许多变化。这些模型的统计参数一般由数百小时的被标记训练数据估计得到。虽然如果训练数据与应用场景的声学特性充分相符，那么这提供高的识别精度，但是可以看出如果语音识别器不得不应付具有显著不同，并且可能高度动态变化的特性的声学
25 环境，那么识别准确性显著降低。

在线和(无)监督的批次自适应技术通过重新估计声学模型参数解决该问题，但是如果只存在很少量的数据和/或计算资源稀少，那么它们都是不可行的，或者-在批次自适应的情况下-不能正确地处理声学环境中的动态变化。

目前的大词汇量连续语音识别器采用隐马可夫模型(Hidden Markov Models(HMM))来根据语音信号,计算具有最大后验概率的单词序列 w 。

隐马可夫模型是处理状态的有限集 $S=\{S_1, \dots, S_N\}$, 并且为状态被占用的每个时间 $t(t=1,2, \dots, T)$ 的输出的观察创造条件的随机自动机 $A=(\pi, A, B)$ 。

初始状态向量

$$\pi = [\pi_i] = [P(s(1)=s_i)], \quad 1 \leq i \leq N \quad (1)$$

给出在时间 $t=1$ 时, HMM 处于状态 s_i 的概率, 转换矩阵

$$A = [a_{ij}] = [P(s(t+1)=s_j | s(t)=s_i)], \quad 1 \leq i, j \leq N \quad (2)$$

保持描述从状态 s_i 到 s_j 的转换的一阶时间不变性过程的概率。观测值是从语音信号得到的连续取值的特征向量 $x \in R$, 输出概率由一组概率密度函数(这里也称为 pdfs)定义:

$$B: [b_i] = [P(x|s(t)=s_i)], \quad 1 \leq i \leq N \quad (3)$$

对于任意给定的 HMM 状态 s_i , 未知分布 $p(x|s_i)$ 通常由基本高斯 pdfs 的混合物近似

$$\begin{aligned} p(x|s_i) &= \sum_{j \in M_i} (w_{ji} \cdot N(x | \mu_{ji}, \Gamma_{ji})) \\ &= \sum_{j \in M_i} \left(w_{ji} \cdot |2\pi\Gamma_{ji}|^{-1/2} \cdot \exp\left(-\frac{(x - \mu_{ji})^T \Gamma_{ji}^{-1} (x - \mu_{ji})}{2}\right) \right) \end{aligned} \quad (4)$$

其中 M_i 是与状态 s_i 相关的高斯函数的集合。此外, x 表示观测的特征向量, w_{ji} 是第 i 个输出分布的第 j 个混合分量权重, μ_{ji} 和 Γ_{ji} 是状态 s_i 下的第 j 个高斯函数的平均矩阵和协方差矩阵。要注意为了符号的简单性, 从等式 4 省略了均值向量的状态和混合分量下标。

现有技术的语音识别器通常由下述组件组成:

·计算允许信号的短小部分(帧)的参数的特征抽取。
·频繁使用的特征是通常由能量值和它们的时间导数富集(enrich)的谱参数或 Mel 频标倒谱系数(MFCC)。

“打标记器”用表示可能有意义的子字单元, 例如依赖于上下文的音子(phone)或子音子的许多标记标识每个特征向量。常见的特征向量的分类技术包括利用高斯混合密度的统计分类或者使用神经网络

络的分类。

·“解码器”截取每个标记作为 HMM 的输出，并计算最大后验概率的单词序列。为了有效地处理来自标记步骤的可选择结果，采用搜索策略和修剪技术。流行的例子是异步栈解码和时间同步 Viterbi (维特比) 解码或集束搜索。

最近已证明通过组合来自并行运行的几个基本识别器的(中间)结果，能够显著降低错字率。可以分出三种主要方法：

·计算特征的不同集合，并把它们组成为被传递给打标记器的单个特征向量的特征组合方法。

·似然组合方法还计算不同的特征向量，但是独立地对它们分类。源于不同的标记步骤的结果根据其证据被组合，对于每一帧，备选标记的单一向量被传送给解码器。

·ROVER(识别器输出表决错误减少)是一种使用动态编程技术把来自几个解码器传递的输出合并到单字假设网络中的后处理方法。在组合网络的每个分支点，后续的表决机构为最终的抄录选择分数最高的单词。

这里提出的发明的主要目的是克服与这些方法相关的一些问题，同时保持增大的识别准确性。

现有技术中已知如果用在未用训练数据正确表示的声学环境中，那么语音识别器的识别准确性显著降低。在诸如桌面口述之类的应用中，通过允许最终用户在不同的环境中登记到该系统中，能够容易地解决该问题，也可考虑输入特征向量的归一化的方法。但是，面对语音作为普遍计算中的输入媒介的重要作用，不允许提前的自适应步骤的应用的数目日益增大。此外，如果识别器不得不处理可能大量的动态变化的声学环境，那么由于缺少足够数量的在线自适应数据，或者由于计算资源有限，自适应方法可能变得不可行。

具有极大量的参数的更准确的声学模型有助于克服这种情况，但是在这里报告的发明中针对的典型应用中是不可行的。除了其它许多应用之外，这些应用是诸如交互式语音响应解决方案，消费设备(移动

电话机、PDA、家用电器)用话音驱动接口,和汽车中的资源短缺语音识别之类的应用。

文献中已证明和单个基本识别器相比,上面提及的组合方法能够在嘈杂环境中产生明显更好的准确性。但是,这些方法对 CPU 施加了不断增大的计算负载,还需要数量增大的存储器来存储几个声学模型和中间结果;于是,它们不适合于资源短缺的语音识别器。

发明内容

于是,本发明的目的是提供一种语音识别方法和系统,它适合于说话者的环境中的动态改变的噪声,以及在由于资源有限,因此只具有有限的计算能力的(嵌入式)系统中运行的特定要求。

本发明的目的由在公开的独立权利要求中陈述的特征实现。在各个从属权利要求中陈述了本发明的其它有利方案和实施例。现在应参考附加的权利要求。

根据本发明的基本方面,提出在语音识别系统内执行下述步骤:

a)用传感器装置收集表征语音识别边界条件的选择基础数据,例如说话人,环境噪声,

b)利用程序控制的传感装置评估收集的数据,即判定引擎,包括软件机构,物理传感器,它们的组合等,

c)根据所述评估,从多个可用的识别器中选择最适合的识别器或者它们的组合。

这样,在具有不断变化的噪声水平,并且其中已存在多个“检测装置”的环境中能够获得显著的优点。从而,传感器装置要被非常广泛地理解为能够提供所述选择基础数据的物理的或者呈逻辑程序形式的任意结构,所述选择基础数据可在存在或者不存在额外的用户输入的情况下被评估,以便由增加的知识将增大识别率的思想所启发,增加定义当前的讲话情景的细节的知识。从而,有利的是,传感器装置可以是判定逻辑,包括软件程序,它解释可由任何物理传感器,比如可检测以特定速度行驶,在特定车型中安装冬季和/或夏季轮胎(pneus)

等而产生的噪声的麦克风，照相机，可从其它可用数据评估的噪声产生设备(例如通风机，音响设备)的 ON/OFF 位置检测的，或者可向用户请求的一些基础数据。当然，也可使用它们的组合。从而，收集的检测数据的一些处理被认为包括在传感器装置内。

15 此外，对于有效的评估来说，最好增加下述步骤：

a) 在实现下述一个或多个的判定逻辑中处理物理传感器输出：统计检验，判定树，模糊隶属关系函数，

b) 从所述处理返回将用在传感器选择/组合判定中的置信度值。

此外，通过增加表达在根据上面提及的处理定义的一组条件下，
10 语音识别有多“好”的评级标准，例如基于数标的标准或者“优”、“中”、“差”任意之一等，用户也可对该过程产生影响。

此外，导致识别器选择判定的所述选择基础数据最好被保存在数据库中，以便于识别器的反复快速选择。这能够主要基于数据库中的查寻，可能还有一些额外的似真性检验来进行识别器选择判定，而不是运行完成的选择判定逻辑。从而，能够节约计算资源。
15

此外，根据本发明的优选方面，提出根据当前的系统负载选择识别器的数目。在具有有限计算资源的嵌入式系统，例如部署在汽车中的嵌入式系统中，这是有利的。

此外，根据本发明的另一优选方面，提出对所考虑的应用特有的
20 各种条件提供模型变换的提前估计。这最好通过只保存一个识别模型如何被变换成另一识别模型的变换规则，而不是保存多个模型本身来实现。这有助于节省存储空间，并且在语音识别系统的运行时间期间，能够在传输中计算不同的模型。

从而，提供选择最适合于当前声学环境中的操作的一个或多个变
25 换的机构，并且提出识别器的动态组合的方法，所述方法在随着时间相当频繁地改变的嘈杂环境中得到改进的识别准确性。

本发明的体系结构提供不得不处理高度变化的声学环境的语音识别应用的改进准确性，此外，通过限制组合的识别器的数目，它还在计算资源可变的情况下，提供可缩放的识别准确性。

这里介绍的发明目的在于在不利的声学环境中，增大通用的基于HMM的语音识别器的稳健性。通过把基于传感器的方法用于声学模型的动态创建以及它们的组合，本发明解决了在上面的背景技术中描述的问题。

- 5 通过把一个或多个模型变换应用于初始的声学模型，动态创建特定于环境的识别器。和在线自适应技术不同，适合的变换不是在运行时间期间计算的，而是在提前的训练步骤中确定的。通用的声学模型和特定于环境的变换与相关的指标函数一起被保存，所述指标函数允许运行时间期间，变换的基于传感器的选择。这确保最匹配当前声学环境的特征的模型的创建和使用。由于在识别过程的组合中，不使用未被传感器识别的模型变换，因此在不必要地增大计算资源的情况下，得到更好的准确性。此外，和自适应模型的存储相比，保存预先计算的变换需要少得多的存储器。
- 10

- 根据本发明，提出借助一个或多个外部存储器取回表征语音识别器工作的声学环境的信息，并把该信息和于一个或多个声学模型的动态创建和组合。
- 15

- 模型的加权组合的方法不在本发明的范围中。但是，通过利用特定于环境的，预先计算的模型变换来创建这些模型是这里描述的发明的一个独创思想。除了已提及的需要较小存储容量的优点之外，这还避免不同特征向量的计算，不同特征向量的计算是基于子带的方法中的一个计算费用高的步骤。
- 20

附图说明

- 附图中举例说明了本发明，但是本发明并不受附图的限制，其中：
- 25 图1是表示根据本发明的一个优选实施例，说明发明原理的概述的示意块图，

图2是表示在汽车中的嵌入式系统中应用的远程信息处理领域中的例证应用的发明基本原理的概述的示意块图。

具体实施方式

现在参考附图，尤其参考图 1，更详细地说明本发明的方法和系统的优选实施例。

通用基准语音识别器 1 被用于从为某一应用特有的各种声学环境 E_j 收集训练语音数据 y -附图标记 2。特定于环境的训练数据 y 被监督地或者不受监督地收集，并被用于所考虑的每个工作环境的声学模型变换的计算，参见块 3。下面，给出举例说明利用预存储变换的特征的两个例子。

• MLLR(最大似然线性回归)自适应通过使用线性变换更新 HMM 均值向量(参见等式 4)。

$$\mu^{(adapt)} = W\mu^{(base)} + \omega,$$

这里变换参数 W 和 ω 被确定，以使自适应数据 y 的似然性达到最大。应注意为使符号简单起见，从等式 4 省略了均值向量的状态和混合分量下标(index)。不同的变换可被应用于属于不同(音子或音位变体)类别的均值向量；例如，把语音和静默均值向量的具体变换看作一个简单例子。在任何情况下，对于每个环境 E_j ，这导致一组变换参数

$$T_j = \{W_i, \omega_i | i=1, \dots, n_j\}$$

• 并行模型组合(PMC)估计“噪声”HMM $\lambda_j^{(noise)} = (\pi, A, B)_j$ 的参数，参考等式 1-3，“噪声”HMM 模拟环境 E_j 的影响，并且与基准识别器的“干净”(或者与环境无关的)HMM 组合。于是，变换参数由“噪声”HMM 的参数给出，即：

$$T_j = \{(p, A, B)_i\}$$

运行时间期间预计算的特定于环境的变换的应用和所得到的声学模型要求识别器训练期间和运行时间期间的声学环境的表征。对于根据本发明实施例的用途，使用一个传感器，所述传感器可被看作计算在本发明的范围中有意义的量值的外部(物理)设备或者计算机程序(软件)或它们的组合。

在块 6 中执行的应用于基准模型的一个或多个模型变换的运行时

间选择以连续监视环境的相关参数的一组传感器 d_k 5 提供的输出为基础。为此，传感器输出经过可采用诸如统计检验，(二元)判定树，或者模糊隶属关系函数之类的方法的判定逻辑，并且对于所考虑的每个环境，返回置信度分数 χ_j ， $1 \leq j \leq n$ 。应注意用于这些检验的参数最好在模型变换估计的自适应数据的处理期间获得。同样，作为描述如何确定环境 E_j 的模糊隶属关系函数的参数的例子，举例说明该原理：

·在识别器训练期间，自适应数据 y 被传送给一组传感器 5，该组传感器 5 可测量源于语音信号本身的任何特征，或者有用的任何外部量值，以便描述自适应数据的环境的声学。

10 ·传感器输出 $z = d_k(y)$ 被量化并以直方图形式保存，所述直方图给出在环境 E_j 中观察 z 的相对频率。随后，直方图可由(多变量)概率密度函数近似，或者可被用于在运行时间期间充当置信度量度的相对频率的直接查找。

·用于传感器 d_k 和环境 E_j 的模糊隶属关系函数 χ_{jk} 可通过特征 z 内分段线性函数的定义的选择，由直方图构成：

$\chi_{jk}(z) = 0$ ，如果 z 小于或等于 z_1 ，或者 z 大于或等于 z_4 ；

$\chi_{jk}(z) = z / (z_2 - z_1)$ ，如果 z_1 小于 z ，并且 z 小于 z_2 ；

$\chi_{jk}(z) = 1$ ，如果 z_2 小于或等于 z ，并且 z 小于或等于 z_3 ；

$\chi_{jk}(z) = 1 - z / (z_4 - z_3)$ ，如果 z_2 小于或等于 z ，并且 z 小于或等于 z_3 ；

20 这里特征值 z_i ， $i \leq 4$ 被选择成使 $p(z \leq z_i) = q_i$ 。概率 q_i 一般被选择成识别 z 的非常罕见的值(例如 $q_1 = 0.05$ ， $q_2 = 0.20$ ， $q_3 = 0.85$ ，和 $q_4 = 0.95$)。同样，这应被理解为只是一种例证定义。

·如果几个传感器被用于监视环境，那么它们各自的置信分数 χ_{jk} 被组合，以便获得特定环境 E_j 的最终分数；例如在通过采用最小值的模糊分数的情况下

$$\chi_j = \min_k \{ \chi_{jk} \},$$

它对应于逻辑“与”运算。当然，也可使用关于模糊集合定义的任何其它运算。

此外，环境(或变换)选择的特征可利用除语音识别器使用的帧速

率之外的帧速率来计算，并且一般将在某一时间间隔内被求平均数，以便获得防止离群值的稳健性。它们可由语音信号本身或者已知的影响声学环境的任意其它量值计算得到。虽然信噪比(SNR)可被看作将从语音信号本身计算的最重要参数之一，不过也可考虑诸如移动汽车的实际速度或路面之类的特征，或者关于说话人的性别或语速的知识5 的利用。于是，对于关联参数的计算和抽取，我们主张全自动方法和需要用户交互作用的方法的使用。

只要置信度分数不显著改变，那么当前的 HMM 声学模型 7 被识别器用于输入的语音信号 8 的解码。如果在 6 中检测到一个或多个新环境，那么应用与这些环境相关的变换 T_j ，变换后的声学模型被用于10 解码。为此，置信度分数被分级，只有 M 个最佳得分的环境的变化被考虑用于未来的处理。重要的是注意考虑中的环境的数目 M 可变化：

- 如果置信度分数不允许环境的明确识别，那么 M 可能较大。

- 如果设备或(远程)识别服务器的工作负载-其计算和分布在现有技术中已知，并且存在于任何现代操作系统中-已分别较高，那么 M 15 将较小，以便实现可接受的响应时间(以识别准确性为代价)。

此外，获得的置信度分数还在识别器组合 8 期间被使用，识别器组合 8 可被用于获得更好的识别准确性。如上所述，现有技术的语音识别器包括三个主要的处理阶段：特征抽取，语音帧的标记和解码。20 而在本发明中，提出单个特征向量的使用，组合可在图 1 中的打标记器 8a 或解码器 8b 中进行。在第一种情况下，归一化的置信度分数被用于增大等式 4 中的 HMM 输出概率：

$$\hat{p}(x|S_i) = \chi_{jk}(z) \cdot p(x_k|S_i)$$

在单词假设的组合的情况下，置信度量度可被用于解析结(tie)，25 如果每个识别器对于指定的语音信号范围产生不同的结果，那么会发生结。这种情况下，提出把从最佳得分的识别器获得的副本(transcription)分配给所考虑的该部分语音信号。

另外参见图 2，以前述实施例在应用于汽车中的嵌入式系统中的远程信息处理领域中的例证应用，给出了发明基本原理的概述。

在第一块 205 中，传感器数据-来自四个传感器设备的选择基础数据从物理设备被读出并被量化，以致数据可用于程序评估。

从而，收集的选择基础数据表示下述可评估的陈述：

1. “驾驶员是女性”，来自具有封闭的图像识别器工具的照相机，
5 -210，
2. “车速为 130 公里/小时”； -220，
3. “空调打开，并且通风机以 75%功率运转”， 230。
4. 无线电打开，并且音量为 8 级中的 4 级，并且播放古典音乐，
-240。

10 随后在步骤 250 中，在数据库中进行查寻，得到其中满足 4 个条件中的 3 个的数据集被保存的判断。从而，与该数据集相关的模型组合被保留为最可能的识别器组件之一。

随后在步骤 260 中，本发明提供的程序控制的仲裁器被用于评估收集的数据，在步骤 270 中，对本例中有意义的多个模型组合确定分
15 数。随后在步骤 280 中，确定当前可用的计算负载。结果可能得到最多 2 模型组合被允许用于语音识别，不过三个最佳得分的提议建议 4 模型的组合。由于其它两个活动的优先权高于语音识别，因此这种限制可被采取。

从而在下一步骤 290 中，选择只具有两个模型的最适合的识别器
20 组合。这需要新的评分过程。

随后在步骤 300 中，选择变换，以便计算选择的最佳两个模型。其它步骤根据上面的说明进行。

可用硬件，软件，或硬件和软件的组合来实现本发明。可集中地
25 在一个计算机系统中实现根据本发明的工具，或者按照分布式方式实现本发明的工具，在这种情况下，不同的部件被散布在数个互连的计算机系统中。适合于实现这里描述的方法的任意类型的计算机系统或其它设备都是适合的。硬件和软件的典型组合可以是具有计算机程序的通用计算机系统，当被加载和执行时，所述计算机程序控制计算机系统执行这里描述的方法。

本发明也能嵌入计算机程序产品中，所述计算机程序产品包含能够实现这里描述的方法的全部特征，并且当被装入计算机系统时，能够实现这些方法。

- 5 本文中的计算机程序意味着一组指令的用任意语言、代码或符号表示的任意表述，所述一组指令意图使具有信息处理能力的系统直接地，或者在下述任一或下述二者之后执行特定的功能：a) 转换成另一种语言，代码或符号；b) 用不同的材料形式再现。

图1

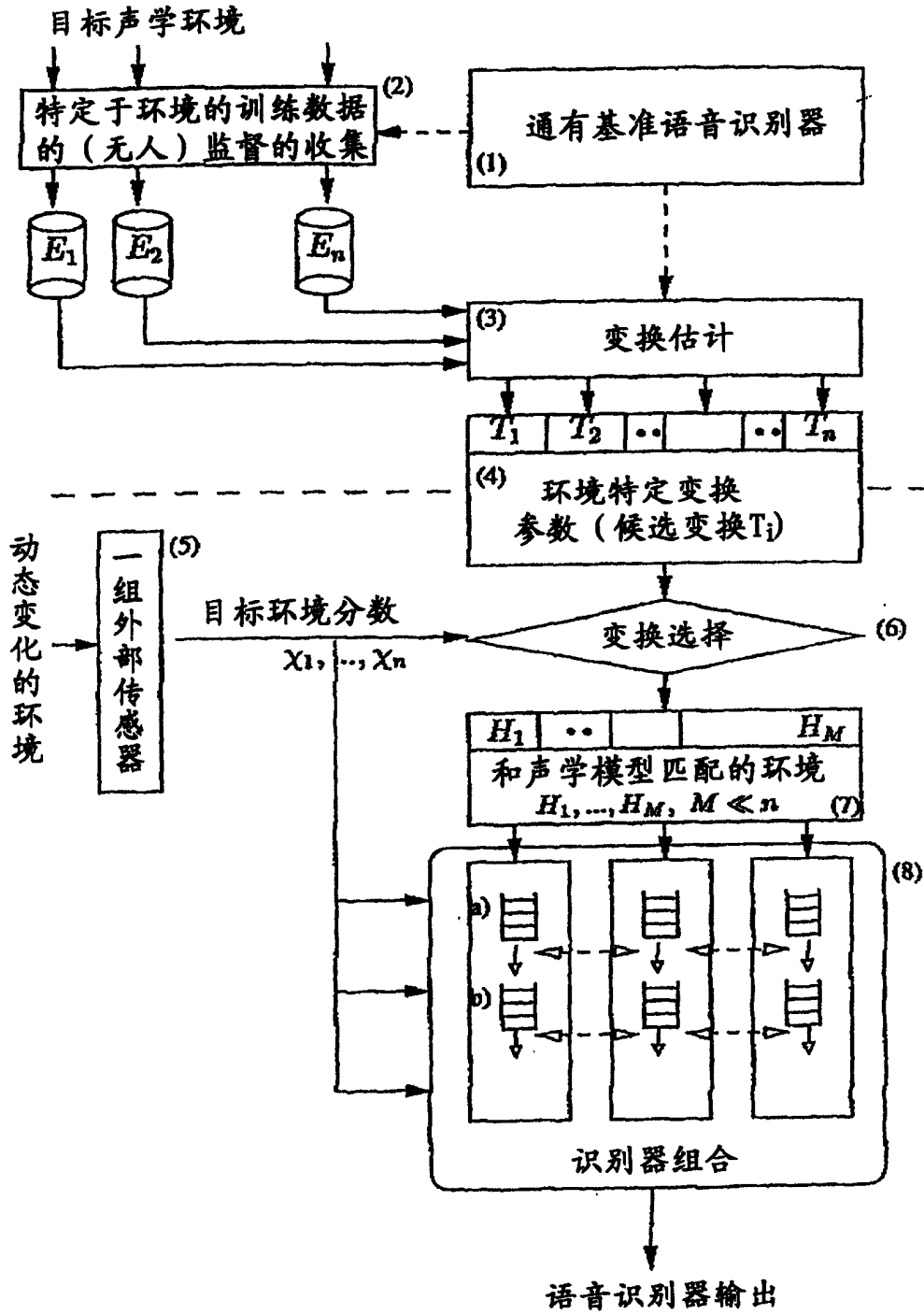


图 2

