

# 基于DF-MAP说话人模型训练

王 军

CSLT, RIIT, THU

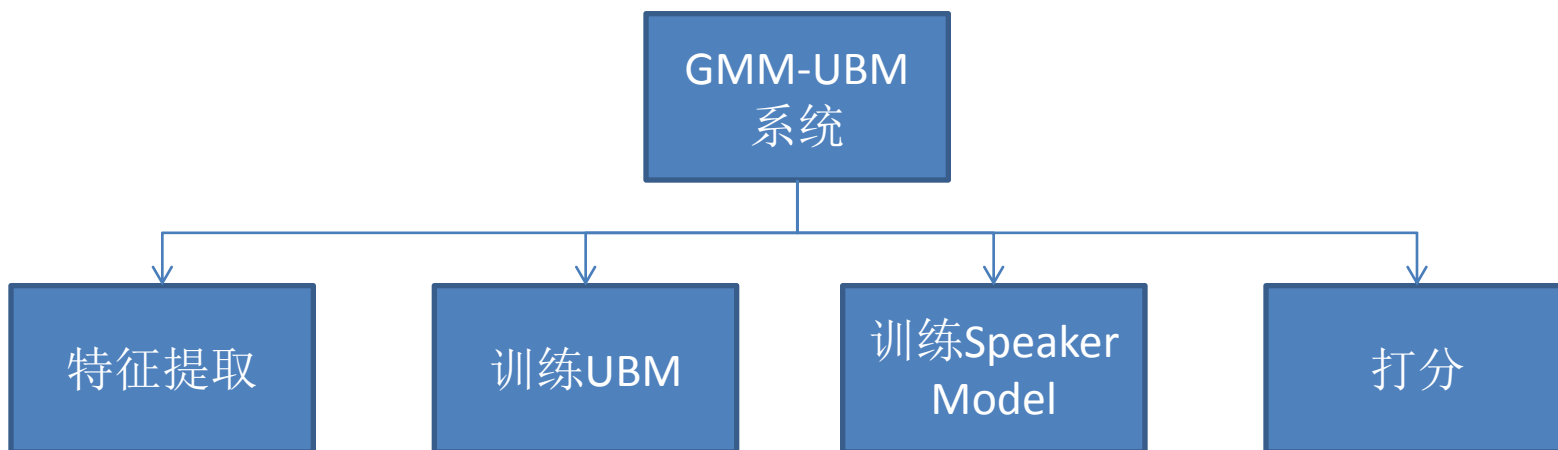
2013-01-06

# 目 录

- 1、基于GMM-UBM说话人识别
- 2、基于DF-MAP的说话人模型训练方法
- 3、实验
- 4、结论

# 1、基于GMM-UBM说话人识别

➤ 自NIST1999说话人确认评测以来，GMM-UBM由于其出色性能成为说话人确认的最主要方法。[1]



# 1、基于GMM-UBM说话人识别

## 特征提取:

分帧，预加重，提取特征，去除静音。

采用MFCC特征。DFT，三角带通滤波与离散谱卷积求对数能量，DCT。[1]

## 训练Speaker Model:

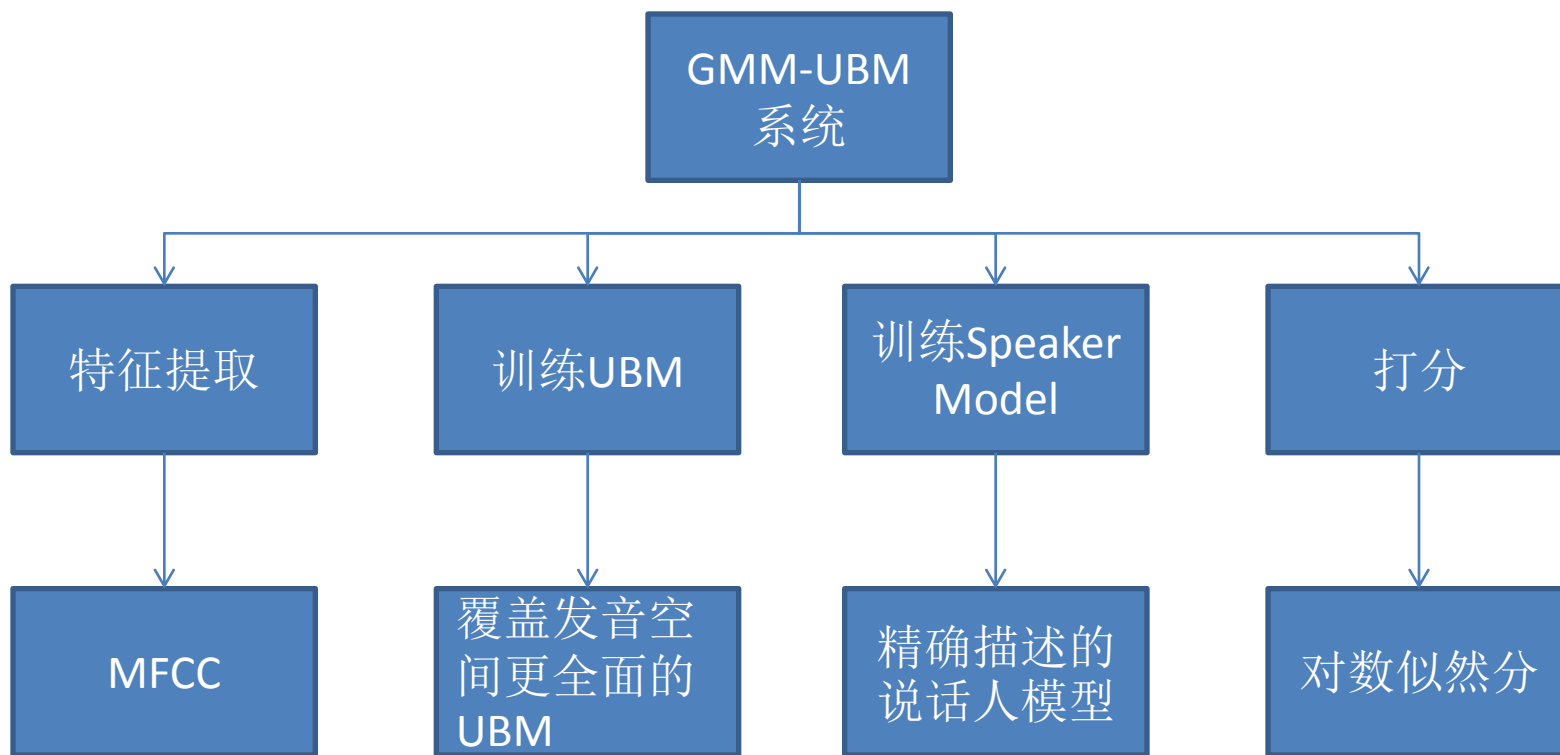
在GMM-UBM系统中，说话人模型通过贝叶斯自适应方法修改UBM中某些参数得到。[1][3]

## 训练UBM:

UBM是一个说话人无关、高阶的GMM。通常由数百人、男女均衡、数小时语音训练而成，用于表示说话人无关的特征分布。[1]

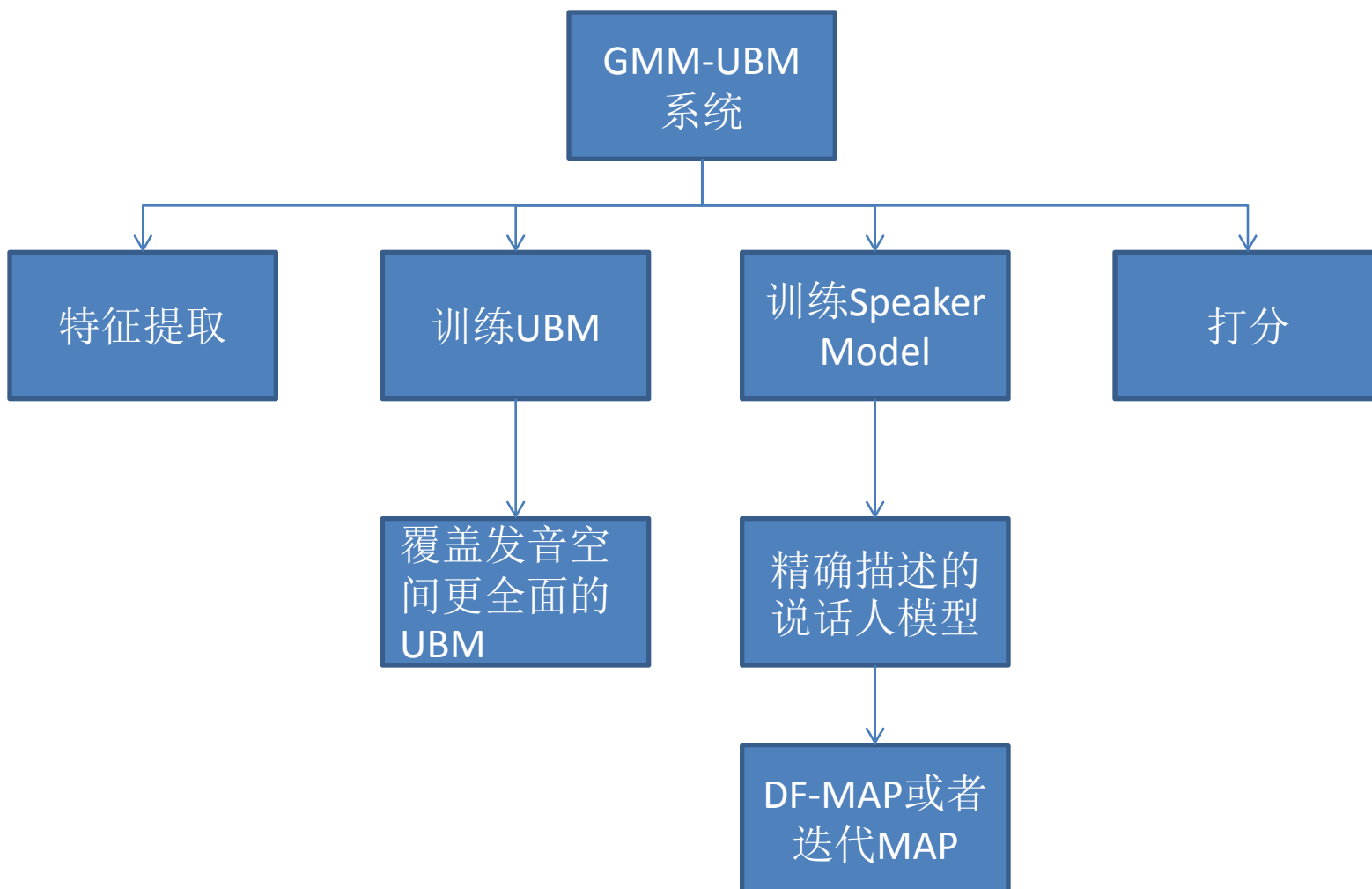
# 1、基于GMM-UBM说话人识别

## ➤ GMM-UBM系统



## 2、基于DF-MAP的说话人模型训练方法

### ➤ GMM-UBM系统性能



## 2、基于DF-MAP的说话人模型训练方法

说话人模型的训练： [1][3]

在GMM-UBM系统中，说话人模型通过贝叶斯自适应方法修改UBM中某些参数得到。

第一步是期望过程，计算训练数据在UBM各单高斯分布上的统计参数；

第二步，用新的统计参数与UBM的参数加权得到说话人模型的参数。

$$\tilde{\mu}_i = \alpha_i E_i(x) + (1 - \alpha_i) \mu_i$$

$$\alpha_i = \frac{n_i}{n_i + r}$$

## 2、基于DF-MAP的说话人模型训练方法

### DF-MAP说话人模型的训练： [1][3]

说话人训练语音覆盖到的发音情况，可以用自己的语音建模；未覆盖到的发音情况，可以用说话人无关的特征分布近似，从而减少训练语音和测试语音不同的影响。 [1]

(1) 训练语音较短时，一次MAP可能导致说话人语音训练不充分。

(2) 训练语音中说话人语音分布不均匀，导致对于某些发音训练充分，而另一些发音训练不充分。



# 3、实验

## 实验环境配置

➤时变声纹数据库<sup>[2]</sup>（方法针对建模，不考虑时变因素）

➤VPR5.0

➤UBM

UBM1: CCB\_Adapt\_110720.ubm

UBM2: UBM\_WFCC\_Reading\_32.ubm

UBM3: 男45人女38人,共83人的语音训练

UBM4: 利用训练语音对UBM3自适应

### 3、实验

表1 不同MAP factor对说话人识别结果的影响

	CCB_Adapt_11 0720	UBM_WFCC_Reading_32
<b>Baseline MAP=16</b>	15.271610%	6.107486%
<b>MAP=0</b>	----	7.116808%
<b>MAP=1</b>	<b>13.433333%</b>	6.553955%
<b>MAP=2</b>	13.641808%	5.940254%
<b>MAP=4</b>	14.052542%	<b>5.707910%</b>
<b>MAP=8</b>	14.633333%	<b>5.705508%</b>

# 3、实验

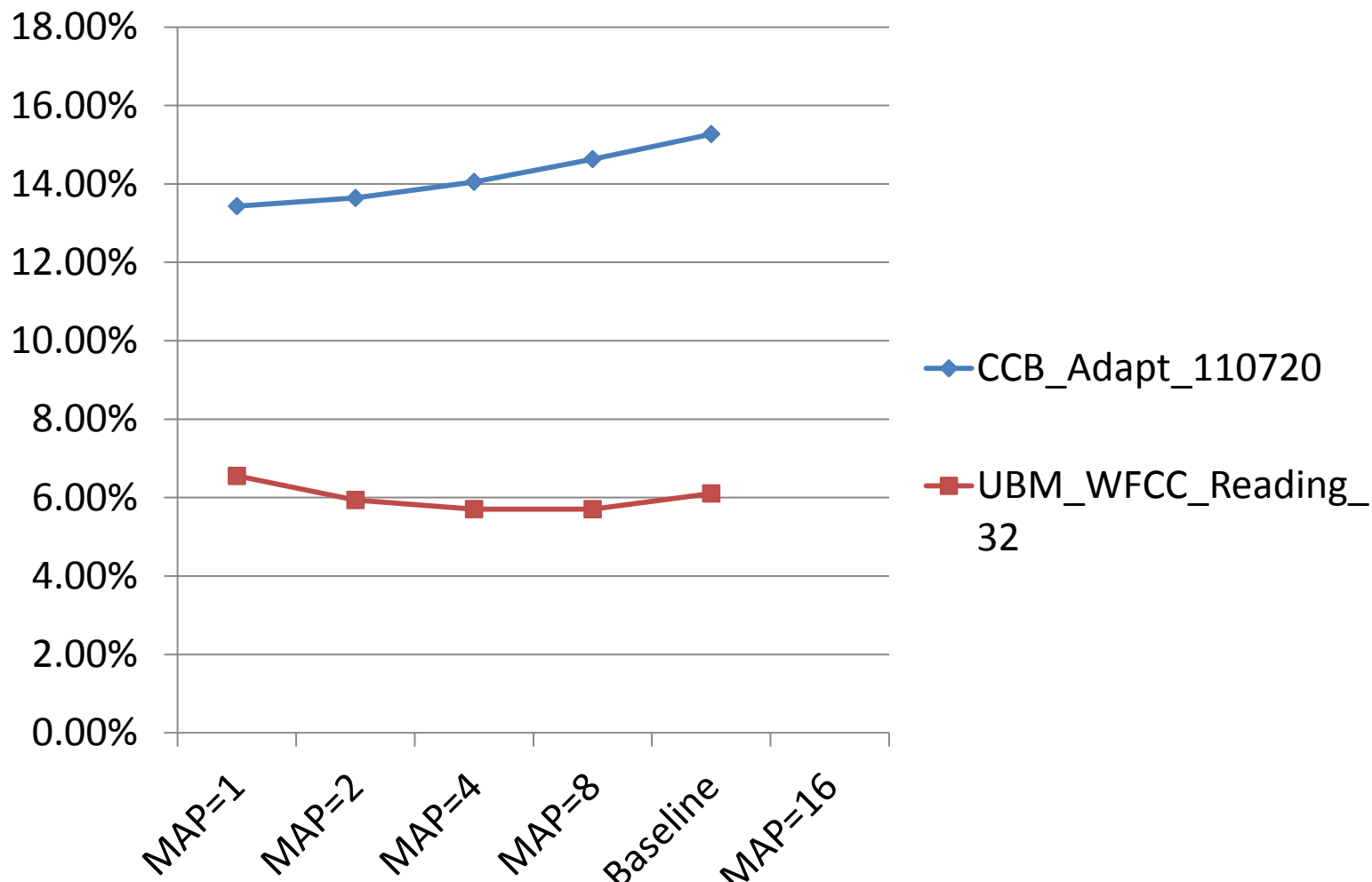


图1 不同MAP factor对说话人识别结果的影响

### 3、实验

表2 MAP迭代算法对说话人识别结果影响

	CCB_Adapt_1107 20/ utterance	UBM_WFCC_Reading_3 2/ utterance
<b>baseline</b>	15.271610%	6.107486%
<b>迭代1</b>	14.120198%	<b>5.734887%</b>
<b>迭代2</b>	13.284463%	5.867090%
<b>迭代3</b>	12.817090%	6.334322%
<b>迭代4</b>	12.605650%	7.502542%
<b>迭代5</b>	12.587429%	10.266667%
<b>迭代6</b>	<b>12.556497%</b>	14.002684%
<b>迭代7</b>	12.600000%	17.966667%
<b>迭代8</b>	12.687429%	22.433333%
<b>迭代9</b>	12.816808%	25.916667%
<b>迭代10</b>	12.933333%	29.153672%

### 3、实验

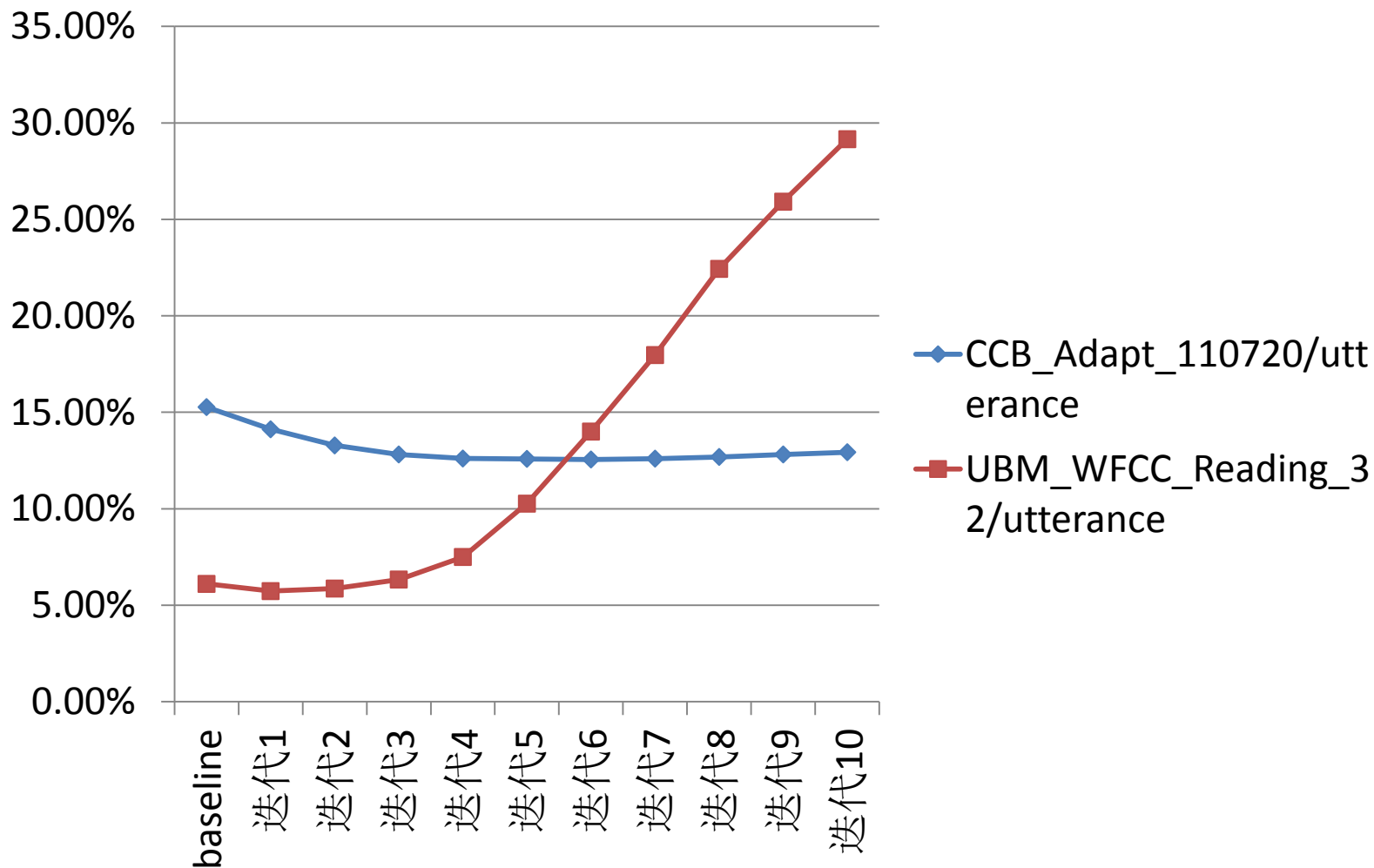


图2 MAP迭代算法对说话人识别结果影响

### 3、实验

表3 DF-MAP算法对说话人识别结果影响a

	数据1	数据3	数据4	数据5	数据6	数据7
<b>baseline</b>	11.969 492%	12.258 399%	14.125 108%	13.373 025%	13.432 591%	14.570 763%
<b>DF-MAP</b>	10.200 000%	10.352 573%	12.344 298%	11.8043 67%	11.7668 44%	12.800 000%
<b>Relative improve ment</b>	14.78 %	15.54%	12.61%	11.73%	12.40%	12.15%

### 3、实验

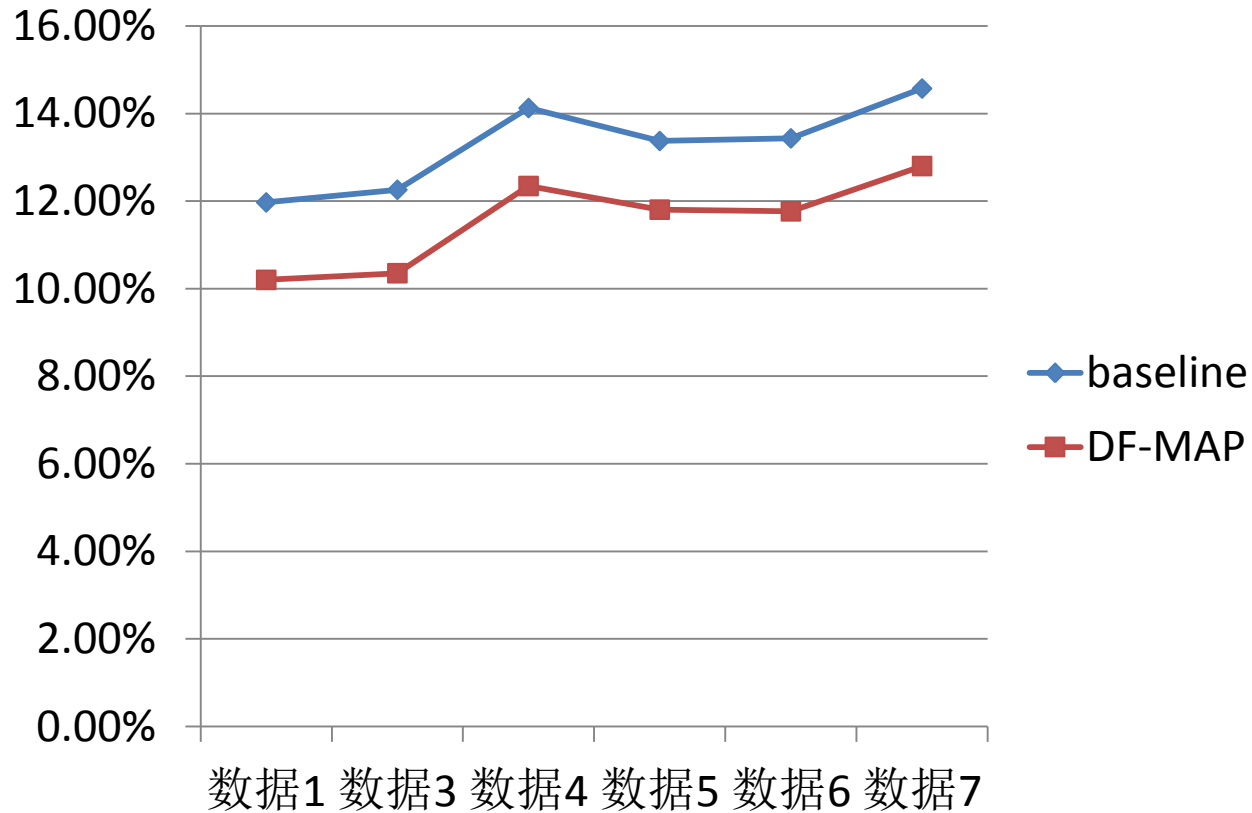


图3 DF-MAP迭代算法对说话人识别结果影响a

### 3、实验

表4 DF-MAP算法对说话人识别结果影响b

	数据1	数据2	数据4	数据5	数据6	数据7
<b>baseline</b>	9.2730 23%	6.10748 6%	11.1739 44%	9.88422 9%	8.80131 7%	11.7405 37%
<b>DF-MAP</b>	8.5833 33%	5.58827 7%	10.7627 12%	9.67796 6%	8.52831 5%	11.1000 00%
<b>Relative improve ment</b>	7.44%	8.5%	3.68%	2.09%	3.1%	5.46%



### 3、实验

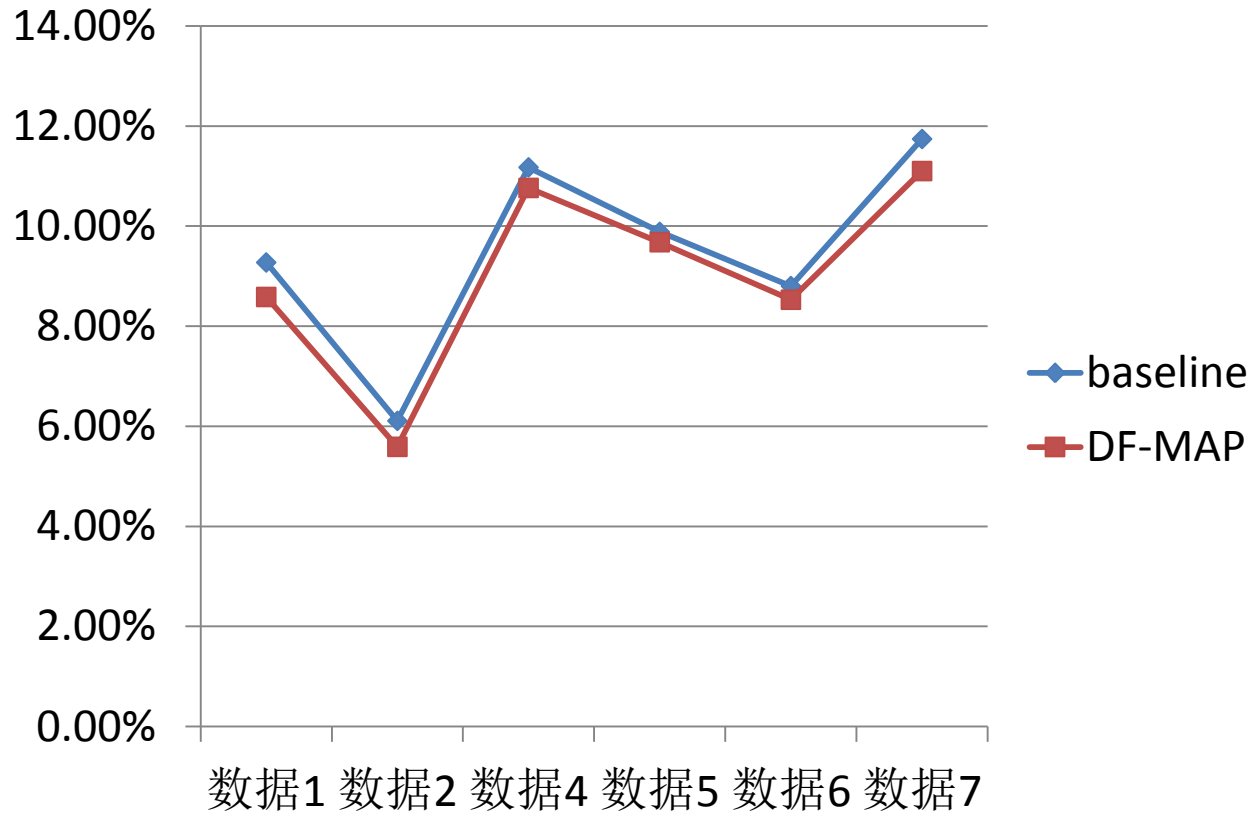


图4 DF-MAP对说话人识别结果影响b

## 4、结论

- 1、通过**DF-MAP**和**MAP**迭代方法可以对训练语音进行均匀、充分的训练，从而得到更精确的说话人模型，提高说话人识别性能。
- 2、最佳**MAP**迭代次数与**UBM**相关。不同的**UBM**会导致不同收敛速度，可依据开发集确定迭代次数。
- 3、最佳**DF-MAP**因子和训练语音相关。不同的训练语音分布导致不同的**DF-MAP**因子。可依据开发集确定**DF-MAP**因子。

。

## 参考文献

- [1]熊振宇,大规模、开集、文本无关说话人辨认研究,:[博士学位论文].北京:清华大学计算机科学与技术系,2005.
- [2]Linlin Wang and Thomas Fang Zheng, “Creation of Time-Varying Voiceprint Database,” Technical Session-6 (Oral), Oriental-COCOSDA, Nov. 24-25, 2010, Kathmandu, Nepal.
- [3]Douglas Reynolds, “Gaussian Mixture Models,” Encyclopedia of Biometric Recognition, 2008.

# UBM训练方法

UBM是一个说话人无关、高阶的GMM。通常由数百人、男女均衡、数小时语音训练而成，用于表示说话人无关的特征分布。

- (1) UBM训练需要较多人数、同信道语音，数据难找。
- (2) UBM训练费时较长。

利用训练语音自适应可以得到一个在训练集上局部最优的UBM，且耗时大幅降低，性能大幅提高。

