

用于读书机器人的声音变换方法

邓杰, 赵群飞

(上海交通大学 图像处理与模式识别研究所, 上海 200240)

摘要: 目前拥有自动读书功能的机器主要通过安装语音库来实现自动朗读, 朗读的声音比较单调。若希望增加朗读声音的多样性就需要安装不同的语音库。这使得语音库容量过大, 浪费了存储资源。本文提出了一种声音变换技术, 通过对单一语音库中的声音进行变换来实现朗读声音的多样性。重点解决了声音变换中音调变换和音色变换相互干扰的问题。首先, 原始声音信号被分解成声音激励信号和声道滤波器信号。然后通过在频域修改声音激励信号的频谱, 利用短时傅立叶幅度谱重构激励信号的方法以及通过修改声道滤波器参数的方法来分别进行音调变换和音色变换。最后再将修改后的声音激励信号和声道滤波器信号重新合成回新的声音信号。实验表明, 利用本文提出的方法能有效地将基音参数和共振峰参数分开调整, 避免了音调变换和音色变换的混叠, 取得了良好的声音变换效果。

关键词: 自动读书机器人; 声音变换; 节奏变换; 音调变换; 音色变换

中图分类号: TP391.41

文献标识码: A

Voice modification for book reading robot

DENG Jie, ZHAO Qunfei

(Institute of Image Processing & Pattern Recognition, Shanghai Jiao tong University, Shanghai 200240, China)

Abstract: Current automatic book reading machines implement reading mainly through installing speech library. However, the voice is rather monotonous. Different speech libraries are needed to increase the diversity of the reading voice, which increases the capacity of the whole speech library, as a result, wastes the storage resources. In order to increase the diversity of the reading voice without increasing the size of speech library, a voice modification technology is proposed. The focus is on how to solve the mutual interference between pitch modification and timbre modification. Firstly, the voice is broken down into excitation and vocal tract filter. Secondly, excitation reconstruction through short-time Fourier transform magnitude and parameters modification of vocal tract filter are respectively used to achieve pitch modification and timbre modification. Finally, modified excitation and modified vocal tract filter are synthesized back to voice signals. Experiments using the proposed method show that pitch parameters and formants parameters are successfully separated. Modification interference between pitch and timbre is eliminated. A pretty good performance is achieved.

Keywords: automatic reading robot; voice modification; tempo modification; pitch modification; timbre modification

1 引言 (Introduction)

世界卫生组织估计, 目前全世界约有一亿八千万人患有视觉残疾, 其中约四千五百万人完全失明。如何帮助盲人获取纸质出版物上的信息具有重要意义。作者所在实验室开发了自动读书机器人来帮助盲人阅读各种出版物。该机器人如图 1 所示, 能自动翻阅书籍, 首先利用 OCR 功能自动识别出版物上的文字, 然后利用 TTS 语音合成功能将识别出的文字内容转换成声音朗诵给听众。语音合成功能依靠装载语音库, 利用微软的 TTS 语音合成引擎来实现。然而由于受到语音库的限制, 只能合成出语音库中存有的声音。这导致合成出的声音比较单调。如果通过安装不同的语音库来实现声音的多样性就会增加存储容量, 造成存储器的浪费。为了解决这一问题, 本文提出了一种声音变换方法, 通过对单一语音库中的声音进行变换来合成出具有多种特色的声音, 以满足读书机器人朗读声音的个性化需求。



图 1 自动读书机器人

Fig.1 Automatic book reading robot

声音变换包括四种类型: 节奏变换、音调变换、音色变换以及强度变换。声音节奏变换的难点在于如何只改变声音的播放速度而不改变音调和音色。声音音调变换的难点是压缩或扩展声音各次谐波间的空间距离而保持短时频谱包络以及声音节奏。声音音色变换的难点在于改变声音共振峰的位置和带宽的同时保持声音的节奏和基音频率。声音强度变换可以通过简单地把信号乘上一个强度因子来得到。研究人员提出了一系列声音变换的方法。例如, 同步叠加法^[1], 波形相似叠加法^[2], 声码器以及各种改进方法^[3-4], 峰值对齐叠加法^[5]。然而上述方法在

基金项目: 项目 1 名称 (编号); 项目 2 名称 (编号)。

通讯作者: 姓名, 电子邮件

收稿/录用/修回: yyyy-mm-dd/yyyy-mm-dd/yyyy-mm-dd

双击此区域修改页脚, 日期部分由编辑部填写。基金一般只标注省级以上基金项目, 请注意核对基金名称。

改变声音基音频率的同时改变了声音共振峰的位置和带宽，反之亦然。结果使得在改变音调的时候却改变了音色。声音的个性特点遭到了破坏。Portnoff提出了用短时傅立叶变换来进行节奏变换^[6]。Griffin等人提出了利用修改短时傅立叶幅度谱来重构信号的方法处理节奏变换和音调变换^[7]。在Griffin的基础上Xinglei等人又提出了实时语谱迭代转换法和超前实时语谱迭代转换法^[8]。这些算法在信号处理的实时性方面有了较大改善，但依然没有把声音的基音频率参数和共振峰参数区分开来。导致在进行音调变换的时候改变了原始声音载有的个性特点，在进行音色变换时，声音的音调也发生了改变。虽然能在一定程度上进行声音变换，但由于音调和音色不能分开调整，使得变换的多样性受到限制。为了解决上述存在的问题，本文引入了源滤波模型，将声音信号的基音频率参数和共振峰参数区分开来，使得音调和音色能分开调整。声音信号首先通过源滤波模型分解成激励信号和声道滤波器信号，然后利用短时傅里叶变换幅度谱信号重构法修改激励信号，并根据需要同时修改声道滤波器参数。最后将修改后的激励信号经修改后的声道滤波器滤波生成变换后的声音信号。

2 声音信号的分解 (Decomposition of the voice signal)

源滤波模型认为声音信号是由声带振动产生的激励信号经过声道滤波产生的。因此，声音信号可以被分解成激励信号和声道滤波器两部分。激励信号携带了声音的基音频率，其大小决定着音调的高低。声道滤波器幅度谱的峰值被称为共振峰，其位置和带宽影响着声音的音色。图2显示了源滤波模型的原理。

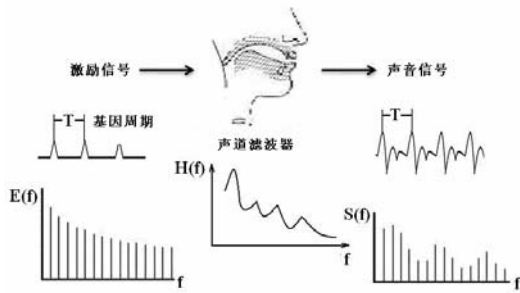


图2 源滤波模型原理
Fig.2 Principle of source filter model

源滤波模型的思想可以通过倒谱分析或线性预测分析来实现。本文采用线性预测分析来分解声音信号。假设 $s(n)$ 代表一段离散声音信号序列， $n=1,2,\dots$ 为序列号。根据线性预测的原理， $s(n)$ 可由之前的 p 个信号的加权值来预测。则 $s(n)$ 的预测值 $\hat{s}(n)$ 记为

$$\hat{s}(n) = \sum_{k=1}^p a_k s(n-k) \quad (1)$$

其中 $\hat{s}(n)$ 和 $s(n)$ 均为离散实数序列， a_k ($k=1,2,\dots,p$) 为加权参数，可由莱文森-杜宾算法求解得到。 p 为线性预测分析的阶数。理论上，当 p 趋近正无穷时，预测声音信号无限接近原始声音信号。通常情况下， p 取 10 到 12。相应的，原始声音信号和其预测信号的误差称为激励信号，定义为

$$e(n) = s(n) - \hat{s}(n) = s(n) - \sum_{k=1}^p a_k s(n-k) \quad (2)$$

对(2)式两边取 z 变换得到

$$E(z) = (1 - \sum_{k=1}^p a_k z^{-k}) S(z) = A(z) S(z) \quad (3)$$

于是在 z 域里，误差信号由原始声音信号 $S(z)$ 和传递函数 $A(z)$ 相乘得到。这里的 $A(z)$ 是一个全零点数字滤波器，代表了声道滤波器。通过调整 $A(z)$ 的零点，就可以调节声音共振峰的频率和带宽。误差信号 $E(z)$ 则携带了声音的基音频率信息。

3 激励信号和声道滤波器的修改 (Modification of excitation signal and vocal tract filter)

将上一节中的激励信号 $e(n)$ 转化到频率域得到频谱信号 $E(mS, \omega)$ 。于是，我们可以在频率域对 $E(mS, \omega)$ 进行修改，然后再将修改后的频谱信号 $E'(mS, \omega)$ 反变换回时域得到修改后的激励信号 $e'(n)$ 。然而频谱参数包括幅度谱和相位谱。而相位在实际操作中很不方便，所以很多时候我们需要直接利用幅度谱来重构信号。也就是说，需要首先得到激励信号 $e(n)$ 在频域的幅度谱 $|E(mS, \omega)|$ ，然后根据需要对幅度谱进行修改得到 $|E'(mS, \omega)|$ ，最后再利用修改后的幅度谱 $|E'(mS, \omega)|$ 来重构激励信号，使重构激励信号的短时傅里叶变换幅度谱尽可能地接近目标激励信号的短时傅里叶变换幅度谱。定义目标激励信号 $e(n)$ 与重构激励信号 $e'(n)$ 相似性距离 $D[e(n), e'(n)]$ 如(4)式所示

$$D[e(n), e'(n)] = \sum_{m=-\infty}^{\infty} \frac{1}{2\pi} \int_{-\pi}^{\pi} [|E(mS, \omega)| - |E'(mS, \omega)|]^2 d\omega \quad (4)$$

其中 $|E(mS, \omega)|$ 是目标激励信号的短时傅里叶变换幅度谱， $|E'(mS, \omega)|$ 为重构激励信号的短时傅里叶变换幅度谱。利用 $|E'(mS, \omega)|$ 取代 $|E(mS, \omega)|$ ，可推出如下的迭代方程^[7]

$$e^{i+1}(n) = \frac{\sum_{m=-\infty}^{\infty} w(mS-n) \frac{1}{2\pi} \int_{-\pi}^{\pi} \hat{E}^i(mS, \omega) e^{jom} d\omega}{\sum_{m=-\infty}^{\infty} w^2(mS-n)} \quad (5)$$

这里有

$$\hat{E}^i(mS, \omega) = E^i(mS, \omega) \frac{|E(mS, \omega)|}{|E^i(mS, \omega)|} \quad (6)$$

其中 $|E^i(mS, \omega)|$ 为 $e^i(n)$ 的短时傅里叶变换幅度谱, $w(mS - n)$ 为窗函数, 本文取汉明窗。从数学上可以证明距离测度 $D[e(n), e^i(n)]$ 会随着迭代次数 i 的增加而逐步减小^[7]。一般情况下迭代次数取 4 到 5 就能使得距离测度减小到可以接受的程度。在本文的声音变换算法中, 为了进一步减轻计算负担, 提高运算的实时性。我们采用了标准叠加公式(7)来取代迭代式(5)。

$$e^{i+1}(n) = \frac{\sum_{m=-\infty}^{\infty} \frac{1}{2\pi} \int_{-\pi}^{\pi} \hat{E}^i(mS, \omega) e^{j\omega n} d\omega}{\sum_{m=-\infty}^{\infty} w(mS - n)} \quad (7)$$

利用(6)、(7)式, 取重构激励信号 $e'(n)$ 为 $e^4(n)$ 或 $e^5(n)$ 。

由上一节可知，声道滤波器 $A(z)$ 为一个全零点滤波器。可根据需要，来调整其零点的分布，从而改变其幅度谱的峰值。不同的滤波器具有不同的滤波特性，产生出不同的音色。

4 声音变换的步骤及实验结果 (Steps and results of voice modification)

图3描述了声音变换的具体步骤,将节奏变换,音调变换,音色变换统一到一个流程当中。图3中的实线框表示节奏变换、音调变换和音色变换都需要进行的处理步骤。虚线框表示三种变换在这些地方的处理有所不同。声音信号首先通过线性预测分析得到激励信号和声道滤波器参数。然后根据不同变换的需要利用短时傅里叶变换幅度谱重构信号法修改激励信号,并修改声道滤波器参数。最后将修改的激励信号和声道滤波器重新合成回时域信号就得到了变换后的声音信号。通过调整声音分析时的窗移长度和声音合成时的窗移长度的比值可实现声音节奏变换,通过调整重采样率可实现音

调变换,通过调整每一帧声音的声道滤波器 A_m 的零点参数可实现音色变换。当然也可以同时调整这些参数以实现特定的混合变换效果。

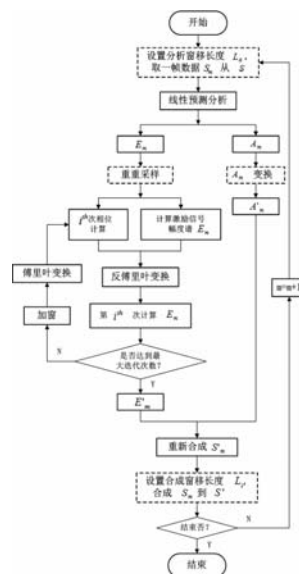


图 3 声音变换流程图
Fig.3 Chart of voice modification

4.1 节奏变换

节奏变换的原理如图 4 所示。其中 L_a 为声音分析时的窗移长度, L_s 为声音合成时的窗移长度, L 为窗口长度。本文中的窗函数取汉明窗。通过改变 L_a/L_s 的比值来调节声音的节奏。当 $L_a/L_s > 1$ 可以加快原声音的节奏, $L_a/L_s < 1$ 则放慢原声音的节奏。

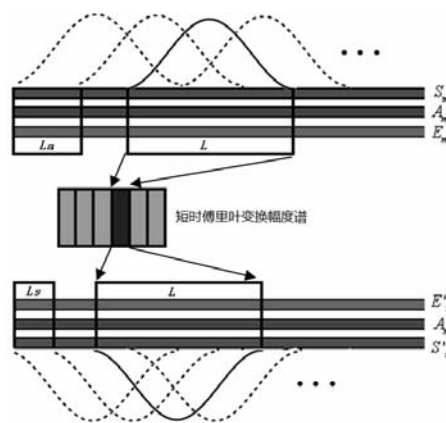


图4 节奏变换原理图
Fig.4 Principle of tempo modification

图 5a 显示了英文句子“*We were away a year ago.*”的短时傅里叶变换幅度谱, 图 5b 是用本文提出的方法进行节奏加快处理后的声音的短时傅里叶变换幅度谱。从图 5 中可以看出, 声音的基音频率、

共振峰的位置和带宽都几乎没有改变，因此声音的音调和音色得到了很好的保持。原始语音的节奏则加快了1.5倍。

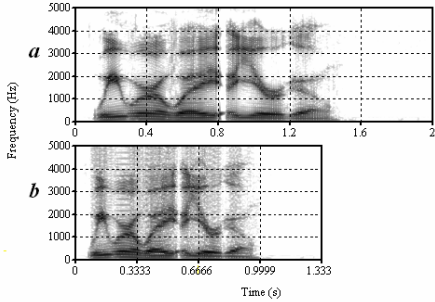


图5 节奏变换实验结果
Fig.5 Experiment result of tempo modification

4.2 音调变换

音调变换原理如图6所示。其中 L_a 为声音分析时的窗移长度， L_s 为声音合成时的窗移长度， L 为声音分析时的窗口长度， L' 为声音合成时的窗口长度。声音信号首先通过线性预测分析被分解成激励信号和声道滤波器。然后对激励信号进行重采样后变换到频率域进行处理。处理完成后再将其通过声道滤波器滤波生成目标声音。在音调变换时， L_a 等于 L_s ，通过改变 L/L' 的比值来调节音调。当 $L > L'$ 时可以提高原声音的音调，而 $L < L'$ 则降低原声音的音调。

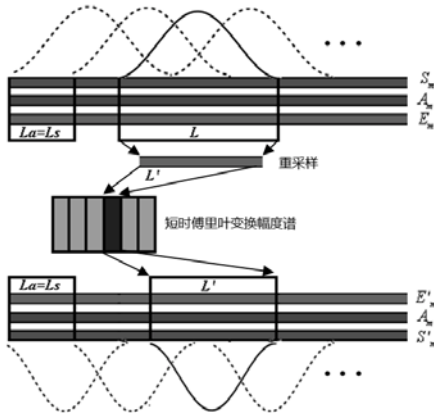


图6 声调变换原理图
Fig.6 Principle of pitch modification

图7a显示了英文句子“We were away a year ago.”的短时傅里叶变换幅度谱，图7b是利用Griffin提出的幅度谱重构信号法进行音调升高处理后的声音的短时傅里叶变换幅度谱，图7c是用本文提出的方法进行音调升高处理后的声音的短时傅里叶变换幅度谱。从图7中可以看出，Griffin的方法虽然升

高了音调，却同时改变了声音共振峰的位置和带宽，在改变音调的同时损坏了原声音的音色。这是因为Griffin的方法^[7]没有将基音频率参数和共振峰参数区分开来造成的。本文提出的方法解决了Griffin方法中的这一不足，有效地将基音频率参数和共振峰参数区分开来，因此不会在升高音调的同时改变声音的音色。

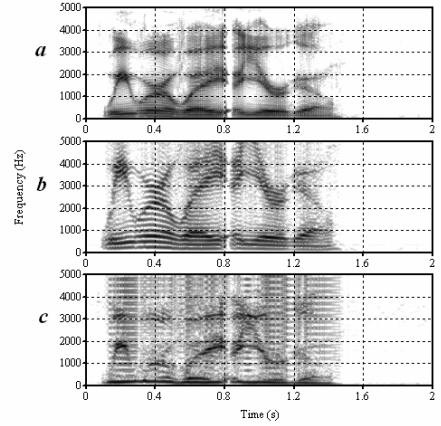


图7 声调变换实验结果
Fig.7 Experiment result of pitch modification

4.3 音色变换

音色变换的原理如图8所示。其中 L_a 为声音分析时的窗移长度， L_s 为声音合成时的窗移长度， L 为窗口长度。在音色变换时， L_a 等于 L_s ，声音分析时的窗口长度和声音合成时的窗口长度均等于 L 。通过线性预测分析得到声道滤波器后，可对其零点进行修改来改变滤波器的滤波特性。

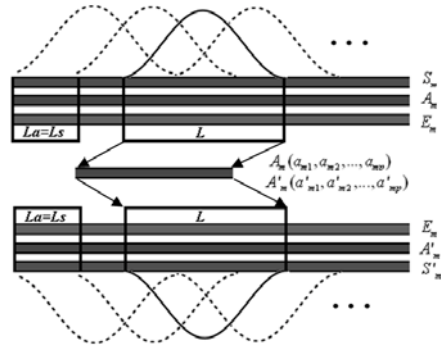


图8 音色变换原理图
Fig.8 Principle of timbre modification

图9a显示了英文句子“We were away a year ago.”的短时傅里叶变换幅度谱，图9b是利用本文方法对声音音色进行修改后的声音的短时傅里叶变换幅度谱。从图9中可以看出，利用本文的方法在对声音共振峰进行修改后并没有改变声音的基音频

率，因此在改变音色时候不会改变音调。而在Griffin的方法中，由于声音的基音频率参数和共振峰参数没有区分开来，音调和音色是混合在一起的，在改变音调的同时必然改变音色，反之亦然，图7b清楚地说明了这一点。

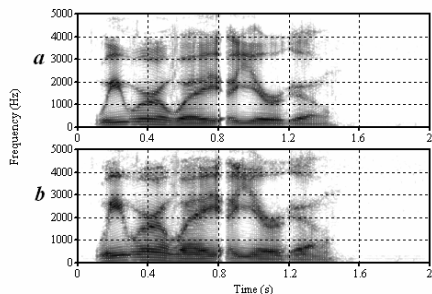


图9 音调变换原理图

Fig.9 Experiment result of timbre modification

5 结论 (Conclusion)

本文提出了一种基于源滤波模型和短时傅立叶变换幅度谱信号重构的声音变换技术。成功解决了读书机器人语音库不能满足多样性需求的问题。利用该方法可以很好地处理声音的节奏变换、音调变换和音色变换。该法将声音的基音频率参数和共振峰参数区分开来，使得音调变换和音色变换不至相互影响。因此能够在改变音调的时候不改变音色，这就保持了原声音的个性特征；在改变音色时候不改变音调，这就保持了原声音的音调特点。音调和音色能在互不影响的情况下各自调整就使得声音变换更加灵活，在一定意义上又增加了声音的多样性。当然也可以利用本文提出的方法将声音的节奏变换、音调变换和音色变换结合在一起以获得一个具有特定效果的混合变换声音。

[1] Wayman, J L, Reinke, R E, Wilson, D L. High quality speech expansion, compression, and noise filtering using the sola method of time scale modification[C]//Twenty-Third Asilomar Conference on Signals, Systems and Computers: Vol.2. San Jose, CA, USA: Maple Press, 1989: 714-717

[2] Verhelst W, Roelands M. An overlap-add technique based on waveform similarity (WSOLA) for high quality time-scale modification of speech[C]//Proceedings of ICASSP '93: Vol.2. New York, NJ, USA: IEEE, 1993:554-557

[3] Laroche J, Dolson M. Improved phase vocoder time- scale modification of audio[J]. IEEE Transactions on Speech and Audio Processing, 1999, 7(3): 323-332.

[4] Dolson M. The Phase Vocoder: A Tutorial Computer Music Journal[J], 1986, 10(4): 14-27.

[5] Dorrán D, Lawlor R, Coyle E. High quality time-scale modification of speech using a peak alignment overlap-add algorithm (PAOLA)[C]//Proceedings of International Conference on Acoustics, Speech and Signal Processing: Vol.1. Piscataway, NJ, USA: IEEE, 2003:700-703

[6] Portnoff M. Time-scale modification of speech based on short-time Fourier analysis[J]. IEEE Transactions on Acoustics, Speech and Signal Processing, 1981, 29(3): 374-390.

[7] Griffin D, Jae L. Signal estimation from modified short-time Fourier transform[J]. IEEE Transactions on Acoustics, Speech and Signal Processing, 1984, 32(2):236-243

[8] Xinglei Z, Beauregard G, Wyse L. Real-Time Signal Estimation From Modified Short-Time Fourier Transform Magnitude Spectra[J]. IEEE Transactions on Audio, Speech, and Language Processing, 2007, 15(5): 1645-1653.

作者简介:

邓杰 (1986--), 男, 硕士。研究领域: 智能机器人控制, 声音信号处理。

赵群飞 (1960--), 男, 理学博士, 教授。研究领域: 特种机器人智能控制, 机器视觉等。