DOI:10.19650/j. cnki. cjsi. J2412996

# 基于谱特征自适应估计的激光相干语音 探测信号增强方法

芮小博1,孔欣玥1,伍 洲2,3,张文喜2,3,曾周末1

(1.天津大学精密测试技术及仪器国家重点实验室 天津 300072; 2.中国科学院空天信息创新研究院 北京 100094;3.中国科学院大学 北京 100049)

摘 要:针对激光相干语音探测引入的缓变宽带背景噪声和测振目标造成的信道作用,本文提出了基于分析-重合成框架,针 对特定说话人的语音增强方法。该方法首先提取观测语音特征:基音频率、浊音概率、MCEP系数,其中,MCEP系数是能够表示 谱包络形状的谱包络特征。通过观测语音的谱包络特征和预训练的对应说话人语音谱包络特征 GMM,估计对应的纯净语音谱 包络特征,再与观测语音的基音频率和浊音概率一起重新合成语音信号,实现语音增强。噪声和信道参数的估计通过最大化观 测语音谱包络特征后验概率的自适应估计实现,然后通过 MMSE 估计得到纯净语音谱包络特征的估计值。合成信号实验和实 际信号采集实验检验了本文提出算法在激光相干语音探测场景下的去噪和均衡能力。

关键词:激光相干语音探测;语音增强;混合高斯模型;矢量泰勒级数

中图分类号: TH741 TN911.7 文献标识码: A 国家标准学科分类代码: 510.40

# Enhancement of speech detected by laser coherent detection method based on spectral feature adaptation

Rui Xiaobo<sup>1</sup>, Kong Xinyue<sup>1</sup>, Wu Zhou<sup>2,3</sup>, Zhang Wenxi<sup>2,3</sup>, Zeng Zhoumo<sup>1</sup>

(1. State Key Laboratory of Precision Measurement Technology and Instrument, Tianjin University, Tianjin 300072, China;
 2. Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China;

3. University of Chinese Academy of Sciences, Beijing 100049, China)

Abstract: To address the issue of slowly varying broadband background noise and channel effects caused by vibration of the target in laser-coherent speech detection, this paper proposes a speech enhancement method for specific speakers based on an analysis-resynthesis framework. This method first extracts the features from the observed signal: pitch, voiced speech probability, and MCEP coefficients, where MCEP coefficients represent the spectral envelope features which can capture the shape of the spectral envelope. A GMM trained by speech features of the corresponding speaker is used to help estimate the spectral envelope features of the clean speech from the spectral envelope features of the observed speech, and then the speech signal is resynthesized by combining it with pitch and voiced speech probability estimated from the observed speech to achieve speech enhancement. The estimation of noise and channel parameters is achieved by adaptation, which maximize the posterior probability of the observed speech's spectral envelope features, and then the estimation of the clean speech spectral envelope features is obtained by MMSE estimation. Both synthesized signal experiments and actual signal acquisition experiments verify the denoising and equalization capabilities of the algorithm in laser coherent speech detection scenarios.

Keywords: laser coherent speech detection; speech enhancement; Gaussian mixture model; vector Taylor series

## 0 引 言

麦克风通过测量振膜在空气声压作用下的振动以感知语音信号,由于声波能量随着距离增大快速衰减,其探测距离十分有限。激光语音探测方法通过激光光束测量 声源附近物体表面微振动以探测语音信号,可以实现超远距离的信号采集并且可深入人员难以到达的区域,在 军事安全等领域具有一定的应用价值。

激光相干语音探测指利用激光多普勒测振技术采集 语音信号。激光多普勒测振技术基于激光光波的多普勒 效应,用相干探测的方法探测多普勒频移,从而解算出测 量位置在光束方向的位移/速度/加速度,是实现激光语 音探测的一种较为成熟的方案<sup>[1]</sup>。然而在应激光相干语 音探测中,一些因素将会对采集语音信号的质量造成影 响:1)声源附近的外界噪声的影响;2)测量系统噪声和 光束受环境影响引入噪声;3)非合作测振目标在语音频 段的频响往往不是一条平稳的水平直线,这些影响具体 表现为引入宽带背景噪声、冲击噪声和信道作用(语音各 频段相对能量发生改变)。

因此,激光探测语音信号需要进行处理以提升信号 质量,语音增强方法能够一定程度上实现语音和噪声的 分离,具有在激光探测语音信号质量提升中进行应用的 潜力。单通道语音增强方法包括:谱减法、维纳滤波法、 最小均方差(minimum mean square error, MMSE)估计等 结合噪声功率谱估计实现的短时谱估计算法<sup>[2]</sup>;在其它 变换域进行语音和噪声分离的子空间法<sup>[3]</sup>和小波去噪 法<sup>[4]</sup>;将语音信号或特征建模为随机自回归过程使用卡 尔曼滤波<sup>[5]</sup>或粒子滤波<sup>[6]</sup>进行估计的方法;引入了语音 和噪声的先验模型的基于码本<sup>[7]</sup>、字典学习<sup>[8]</sup>、隐马尔可 夫模型<sup>[9]</sup>的语音增强方法;基于语音生成模型,通过带噪 语音特征估计纯净语音特征并重新合成语音的分析-重 合成语音增强方法<sup>[10]</sup>;基于深度学习的语音增强方法<sup>[11]</sup> 等。通常应用于数字语音通信设备、视听会议、自动语音 识别系统等。

目前,激光相干语音探测与语音增强方法结合的 相关应用研究较少,吕韬<sup>[12]</sup>采用最优对数谱(optimallymodified log-spectral amplitude,OMLSA)估计算法与相 位补偿相结合进行信号降噪处理,然后使用基于半波 整流的语音带宽扩展算法扩展相干侦测语音的带宽, 但测振目标造成的信道作用其实并不是简单的窄带宽 导致信号高频部分缺失;王亚慧等<sup>[13-14]</sup>分别提出拉 普拉斯分布下 MMSE 谱减语音增强算法和循环生成 对抗网络对采集语音进行增强,谱减法只能实现降低 噪声的影响,而激光探测信号中复杂多变的噪声干扰 使得基于深度学习的方法需要大量不同条件下的采集 数据。整体来看,相关研究有待进一步深入,缺乏对 侦听等应用场景下特定人语音信号增强等针对性的 研究。

针对以上问题与不足,本文提出了基于谱特征自适 应估计的激光相干语音探测信号增强方法,引入了纯净 语音谱包络特征高斯混合模型,实现噪声和信道的自适 应估计,通过信号分析与重合成,一定程度上实现了特定 目标人语音激光相干探测信号的信号去噪,以及语音信 号在频率响应上的均衡。

# 1 方法

本文提出的语音增强算法基于分析-重合成框架,整体流程如图1所示。首先提取带噪语音的浊音概率、基 音频率、谱包络特征参数;然后根据提前训练好的特定说 话人的语音谱包络特征高斯混合模型(Gaussian mixture model, GMM),对噪声和信道参数进行预估计和自适应 估计,进而对纯净语音谱包络特征进行估计;最后,用带 噪语音的浊音概率、基音频率和估计纯净语音的谱包络 特征参数进行语音重合成,得到增强语音。



Fig. 1 Flow chart of the method proposed in this paper

#### 1.1 谐波加噪声语音生成模型

谐波加噪声模型(harmonic plus noise model, HNM) 是一类将语音建模成谐波成分与随机成分相叠加的语音 生成模型<sup>[15]</sup>。本文将语音信号表示为基音频率(声调)、 浊音(声带振动)概率、谱包络特征参数,然后使用 HNM 算法根据这些参数生成语音信号。理论上,带噪信号的 基音频率和浊音概率与纯净语音一致,因此可以将语音 增强问题转换为纯净语音信号谱包络参数的估计。

1)特征提取

通过在强噪声条件下具有鲁棒性的基音频率估计算 法 PEFAC (pitch estimation filter with amplitude compression, PEFAC)<sup>[16]</sup>估计经短时傅里叶变换(short time Fourier transform, STFT)后每帧信号的基音频率 $f_i$ 和浊 音概率 $pv_i$ ,  $i = 0, 1, \dots, I - 1$ 表示帧索引。谱包络特征参 数选择 MCEP(mel-cepstral)系数<sup>[17]</sup>,利用该系数可将一 帧语音信号的频谱  $S(e^{i\omega})$ 表示为 M 阶 MCEP 系数:

$$S(z) = \exp\sum_{m=0}^{M} c_{\alpha} [m] z_{\alpha}^{-m}$$
(1)

$$z_{\alpha}^{-1} = \psi(z) = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad |\alpha| < 1$$
(2)

式中:  $c_{\alpha}$  表示 MCEP 系数,  $z_{\alpha}^{-1}$  为全通传递函数,  $\alpha$  控制频 率的扭曲, 如果  $\alpha$  不为 0, 则系统具有梅尔尺度频率特 性。一帧语音信号的对数幅度谱可以表示为:

$$\log(|S(e^{j\omega})|) = \sum_{m=0}^{M} c_{\alpha}[m] real(\psi(e^{j\omega})^{m})$$
(3)

2)语音合成

语音合成基于1)特征提取所述的特征,计算谐波成 分与随机成分并进行叠加,谐波成分分布于浊音段,表 示为:

$$s_{h} \left[ \frac{1}{2} (l_{w} - 1) + i l_{h} + l \right] =$$

$$\left[ \sum_{m=1}^{M_{i}} \hat{A}_{i,m} [l] \cos(\hat{\theta}_{i,m} [l]), \quad pv_{i} > 0.5 \right]$$

$$0, \qquad pv_{i} \leq 0.5$$
(4)

 $pv_i > 0.5$ 时判定这一帧语音信号为有话段,在对应 帧的中间部分加入谐波成分。 $l_w$ 表示计算帧长, $A_{i,m}[l]$ ,  $\theta_{i,m}[l]$ 分别是第i帧m次谐波的幅度与相位,其中,l=0,  $1, \dots, l_h = 1, l_h$ 表示帧移。 $\hat{A}_{i,m}[l]$ 通过三角滤波器估计 幅度谱包络在对应谐波频率处的值得到,其它时间点处 的幅度经线性插值得到,谐波频率的估计与谐波幅度相 同,谐波相位经随机初始化,然后根据谐波频率和采样 时间计算得到。

$$s_{n}[n] = \begin{cases} c \sum_{i=0}^{l-1} w[n - il_{h}] \text{IDFT}\{\hat{E}_{i}[k] \text{HPF}\{D_{i}[k]\}\}, & pv_{i} > 0.5 \\ c \sum_{i=0}^{l-1} w[n - il_{h}] \text{IDFT}\{\hat{E}_{i}[k]D_{i}[k]\}, & pv_{i} \leq 0.5 \end{cases}$$

式中:  $\hat{E}_i[k]$  为谱包络估计,  $D_i[k]$  为能量归一化的高斯 白噪声的频谱, IDFT { • } 表示逆离散傅里叶变换操作,

(5)

HPF{·}表示高通滤波操作,常数 c 和窗 w 用于实现逆短时傅里叶变换。

合成语音信号最终表示为:  
$$s[n] = s_h[n] + s_n[n]$$
 (6)

#### 1.2 语音特征参数(MCEP 系数)增强方法

1) 观测信号特征模型

为了从观测语音谱包络特征中估计纯净语音谱包络特征的方法,本文采用了基于矢量泰勒级数(vector Taylor series, VTS)的特征增强方法<sup>[18]</sup>,首先需要建立观测信号模型。

将激光相干语音探测信号简单表示为纯净语音在 信道(卷积)和加性噪声的作用下得到观测信号,观测 信号经测量系统的采集和离散傅里叶变换转换到频域 可得:

$$Y_k = H_k X_k + N_k \tag{7}$$

 $X_k$ 、 $H_k$ 、 $N_k$ 和 $Y_k$ 分别表示纯净语音信号、信道、噪声、 观测信号的 DFT 系数, k代表频率采样点,由此可得观测 信号的功率谱:

$$|Y_{k}|^{2} = |H_{k}|^{2} |X_{k}|^{2} + |N_{k}|^{2} + 2\cos\theta_{k} |H_{k}| |X_{k}| |N_{k}|$$
(8)

其中,  $\theta_k$  表示复变量( $H_k X_k$ ) 与  $N_k$  之间的夹角。

用向量形式表示信号的频谱  $Y = [Y_0, Y_1, \dots, Y_{k-1}]^T$ , 同理可得  $X \setminus N \setminus H$ , 根据式(3) 通过 MCEP 系数估计观测 信号的对数功率谱:

$$\log(|\mathbf{Y}|^2) \approx \mathbf{C}^{-1}\mathbf{y} \tag{9}$$

 $y = [y_0, y_1, \dots, y_m]^T$ 表示 MCEP 系数组成的向量,同 理可得 $x, h, n, C^{-1}$ 为转化矩阵,其伪逆表示为C,所以观 测语音信号的 MCEP 系数可估计为:

 $\alpha = [\cos\theta_0, \cos\theta_1, \cdots, \cos\theta_k, \cdots]^T$ ,本文将 $\alpha(0 \le \alpha_k \le 1)$ 设为定值, · 表示两个向量间的逐元素相乘,使用关于x, h, n的1阶 VTS 估计可得:

$$y \approx \boldsymbol{\mu}_{x} + \boldsymbol{\mu}_{h} + g(\boldsymbol{\mu}_{x}, \boldsymbol{\mu}_{h}, \boldsymbol{\mu}_{n}) + \boldsymbol{G}(\boldsymbol{x} - \boldsymbol{\mu}_{x}) + \boldsymbol{G}(\boldsymbol{h} - \boldsymbol{\mu}_{h}) + (\boldsymbol{I} - \boldsymbol{G})(\boldsymbol{n} - \boldsymbol{\mu}_{n})$$
(12)

其中, $\mu_x \mu_h \mu_n$ 分别是x h, n的均值矢量,I为单位 矩阵,

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\Big|_{\boldsymbol{\mu}_{x},\boldsymbol{\mu}_{n},\boldsymbol{\mu}_{h}} = \frac{\partial \mathbf{y}}{\partial \boldsymbol{h}}\Big|_{\boldsymbol{\mu}_{x},\boldsymbol{\mu}_{n},\boldsymbol{\mu}_{h}} = \boldsymbol{G}$$
(13)

$$\frac{\partial y}{\partial n} = I - G \tag{14}$$

$$\boldsymbol{G} = \boldsymbol{I} - \boldsymbol{C} \operatorname{diag} \left( \frac{\exp(\boldsymbol{C}^{-1}(\boldsymbol{\mu}_n - \boldsymbol{\mu}_x - \boldsymbol{\mu}_h)) + \boldsymbol{\alpha} \exp(\boldsymbol{C}^{-1}(\boldsymbol{\mu}_n - \boldsymbol{\mu}_x - \boldsymbol{\mu}_h)/2)}{1 + \exp(\boldsymbol{C}^{-1}(\boldsymbol{\mu}_n - \boldsymbol{\mu}_x - \boldsymbol{\mu}_h)) + 2\boldsymbol{\alpha} \exp(\boldsymbol{C}^{-1}(\boldsymbol{\mu}_n - \boldsymbol{\mu}_x - \boldsymbol{\mu}_h)/2)} \right) \boldsymbol{C}^{-1}$$
(15)

其中,diag(•)代表构造对角矩阵,其对角分量值分 别等于括号中向量的元素。

2) 纯净语音特征向量估计

线下训练阶段,用大量纯净语音的 MCEP 系数训练 GMM。GMM 用多个多维高斯分布来拟合多维随机变量 的概率分布,利用 GMM 建立纯净语音 MCEP 系数概率 分布的统计模型表示为:

$$p(\boldsymbol{x}) = \sum p(s) \operatorname{N}(\boldsymbol{x} | \boldsymbol{\mu}_{x,s}, \boldsymbol{\Sigma}_{x,s})$$
(16)

其中, s 代表 GMM 中的聚类, 对于聚类 s, p(s) 表示 其先验概率,  $N(x | \boldsymbol{\mu}_{x,s}, \boldsymbol{\Sigma}_{x,s})$  代表均值为 $\boldsymbol{\mu}_{x,s}$ , 协方差矩 阵为  $\boldsymbol{\Sigma}_{x,s}$  的高斯函数, 训练时令协方差矩阵为对角阵, 即  $\boldsymbol{\Sigma}_{s} = \text{diag}([\sigma_{s,1}^{2}, \sigma_{s,2}^{2}, \cdots \sigma_{s,M}^{2}])_{\circ}$ 

纯净语音特征向量估计阶段,根据观测信号的特征 向量y估计x的后验概率分布可得:

$$p(\mathbf{x} | \mathbf{y}) = \sum_{s} p(\mathbf{x} | \mathbf{y}, s) p(s | \mathbf{y})$$
(17)

因此,纯净语音特征的 MMSE 估计量为:

$$\hat{\boldsymbol{x}} = \mathbf{E}(\boldsymbol{x} | \boldsymbol{y}) = \sum p(s | \boldsymbol{y}) \mathbf{E}(\boldsymbol{x} | \boldsymbol{y}, s)$$
(18)

通过式 (12)将纯净语音特征 GMM 的每个聚类,映 射为观测语音 GMM 的对应聚类,所估计的观测语音特 征向量 GMM 各聚类的均值向量和协方差矩阵可表 示为:

$$\boldsymbol{\mu}_{y,s} \approx \boldsymbol{\mu}_{x,s} + \boldsymbol{\mu}_{h} + g(\boldsymbol{\mu}_{x,s}, \boldsymbol{\mu}_{h}, \boldsymbol{\mu}_{n})$$
(19)

$$\boldsymbol{\Sigma}_{\boldsymbol{y},\boldsymbol{s}} \approx \boldsymbol{G}_{\boldsymbol{s}} \boldsymbol{\Sigma}_{\boldsymbol{x},\boldsymbol{s}} \boldsymbol{G}_{\boldsymbol{s}}^{\mathrm{T}} + (\boldsymbol{I} - \boldsymbol{G}_{\boldsymbol{s}}) \boldsymbol{\Sigma}_{\boldsymbol{n}} (\boldsymbol{I} - \boldsymbol{G}_{\boldsymbol{s}})^{\mathrm{T}}$$
(20)

通过得到的观测语音特征向量 GMM 来实现式(18) 中对于  $p(s|\mathbf{y})$  的估计。本文将信道作用视为固定值,因此式(20)中不考虑信道方差。

根据**y**关于**x**的一阶 VTS 估计计算 E(**x** | **y**, s)<sup>[19]</sup>, 聚 类 s 中, 将 **y** 估计为:

$$\mathbf{y} \approx \boldsymbol{\mu}_{x,s} + \boldsymbol{\mu}_{h} + g(\boldsymbol{\mu}_{x,s}, \boldsymbol{\mu}_{h}, \boldsymbol{\mu}_{n}) + \boldsymbol{G}_{s}(\boldsymbol{x} - \boldsymbol{\mu}_{x,s}) + \boldsymbol{G}_{s}(\boldsymbol{h} - \boldsymbol{\mu}_{h}) + (\boldsymbol{I} - \boldsymbol{G}_{s})(\boldsymbol{n} - \boldsymbol{\mu}_{n})$$
(21)

由于x 与 h、n不相关,聚类s中,x和y的联合分布可写作:

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}_{s} \sim N\left(\begin{bmatrix} \boldsymbol{\mu}_{x,s} \\ \boldsymbol{\mu}_{y,s} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{x,s} & \boldsymbol{\Sigma}_{xy,s} \\ \boldsymbol{\Sigma}_{yx,s} & \boldsymbol{\Sigma}_{y,s} \end{bmatrix}\right)$$
(22)

$$\Sigma_{xy,s} = \Sigma_{x,s} G_s^{\mathrm{T}}$$
(23)  
由此可得.

$$\boldsymbol{\mu}_{x,s} + \boldsymbol{\Sigma}_{x,s} \boldsymbol{G}_{s}^{\mathsf{T}} \boldsymbol{\Sigma}_{y,s}^{-1} (\boldsymbol{y} - \boldsymbol{\mu}_{y,s})$$

$$(24)$$

一阶 VTS 估计引入的偏差可能会导致 E(x|y,s) 与相差过大,因此将 E(x|y,s) 修正为:

$$E'(\boldsymbol{x} | \boldsymbol{y}, s) = \boldsymbol{\mu}_{x,s} + \boldsymbol{\sigma}'_{s} \odot \tanh\left(\frac{E(\boldsymbol{x} | \boldsymbol{y}, s) - \boldsymbol{\mu}_{x,s}}{\boldsymbol{\sigma}'_{s}}\right) \quad (25)$$

$$\boldsymbol{\sigma}_{s}' = \boldsymbol{\sigma}_{xs} \odot \frac{E(\boldsymbol{x} | \boldsymbol{y}, s) - \boldsymbol{\mu}_{x,s}}{\| E(\boldsymbol{x} | \boldsymbol{y}, s) - \boldsymbol{\mu}_{x,s} \|_{2}}$$
(26)

式(25)和(26)中, ③表示按位相乘,除号表示按位 相除,  $\sigma_{xs}$ 为对角阵  $\Sigma_{x,s}$ 中元素开方组成的向量  $\sigma_{xs} = [\sigma_{xs,1}, \sigma_{xs,2}, \dots, \sigma_{xs,M}]^{T}$ 。

3)噪声与信道特征参数的自适应估计

利用 EM(期望-最大化)算法对信道和噪声特征的 均值 $\mu_{h}$ , $\mu_{n}$ 进行自适应估计<sup>[20]</sup>。该方法在已知一段观 测语音的情况下,通过迭代的方式最大化其后验概率进 行参数估计。在迭代估计参数 $\mu_{h}$ 、 $\mu_{n}$ 时,将式(19)展 开为关于 $\mu_{h}$ 、 $\mu_{n}$ 的当前值 $\mu_{h_{0}}$ 、 $\mu_{n_{0}}$ 的一阶 VTS。

 $\mu_{y,s} \approx \mu_{x,s} + \mu_{h_0} + g(\mu_{x,s}, \mu_{h_0}, \mu_{n_0}) + G_s(\mu_h - \mu_{h_0}) + (I - G_s)(\mu_h - \mu_{n_0})$ (27) EM 算法的辅助函数可表示为:

$$Q(\boldsymbol{\lambda} | \bar{\boldsymbol{\lambda}}) = \sum_{t} \sum_{s} p(s_{t} = s | \boldsymbol{y}_{t}, \bar{\boldsymbol{\lambda}}) \log p(\boldsymbol{y}_{t} | s_{t} = s, \boldsymbol{\lambda})$$
(28)

其中,  $\lambda$  和  $\overline{\lambda}$  分别为噪声和信道的新旧参数, 用  $\gamma_{i}(s)$  表示 t 时刻聚类 s 的后验概率:

$$\boldsymbol{\gamma}_{t}(s) = p(s_{t} = s | \boldsymbol{y}_{t}, \bar{\boldsymbol{\lambda}})$$
(29)

EM 算法的步骤 M 通过最大化 Q 迭代求解 $\mu_h$ 和 $\mu_n$ 的更新值,加入惩罚项限制更新向量与原始向量之间的距离,避免一阶泰勒估计失效,以 $\mu_h$ 为例,更新值可写为:

$$\boldsymbol{\mu}_{h} = \operatorname{argmin}_{\boldsymbol{\mu}_{h}} \left( \left( -Q \right) + \boldsymbol{\eta}_{h} \| \boldsymbol{\mu}_{h} - \boldsymbol{\mu}_{h_{0}} \|_{2}^{2} \right), \ \boldsymbol{\eta}_{h} > 0$$
(30)

辅助函数加上惩罚项之后分别对 $\mu_n$ 和 $\mu_n$ 求导取 0 求最值得:

$$\sum_{t} \sum_{s} \boldsymbol{\gamma}_{t}(s) \boldsymbol{G}_{s}^{\mathrm{T}} \boldsymbol{\Sigma}_{\boldsymbol{y},s}^{-1}(\boldsymbol{y} - \boldsymbol{\mu}_{\boldsymbol{y},s}) + \boldsymbol{\eta}_{h}(\boldsymbol{\mu}_{h} - \boldsymbol{\mu}_{h_{0}}) = 0$$
(31)
$$\sum_{t} \sum_{s} \boldsymbol{\gamma}_{t}(s) (\boldsymbol{I} - \boldsymbol{G}_{s})^{\mathrm{T}} \boldsymbol{\Sigma}_{\boldsymbol{y},s}^{-1}(\boldsymbol{y} - \boldsymbol{\mu}_{\boldsymbol{y},s}) + \boldsymbol{\eta}_{n}(\boldsymbol{\mu}_{n} - \boldsymbol{\mu}_{n_{0}}) = 0$$
(32)

对两个参数分别进行重估计,可分别得到 $\mu_h$ 和 $\mu_n$ 的更新值:

$$\boldsymbol{\mu}_{h} = \boldsymbol{\mu}_{h_{0}} + \left(\sum_{t}\sum_{s}\gamma_{t}(s)\boldsymbol{G}_{s}^{\mathsf{T}}\boldsymbol{\Sigma}_{y,s}^{-1}\boldsymbol{G}_{s} + \boldsymbol{\eta}_{h}\boldsymbol{I}\right)^{-1}$$

$$\left(\sum_{t}\sum_{s}\gamma_{t}(s)\boldsymbol{G}_{s}^{\mathsf{T}}\boldsymbol{\Sigma}_{y,s}^{-1}(\boldsymbol{y}_{t} - \boldsymbol{\mu}_{x,s} - \boldsymbol{\mu}_{h_{0}} - g(\boldsymbol{\mu}_{x,s}, \boldsymbol{\mu}_{h_{0}}, \boldsymbol{\mu}_{n_{0}}))\right)$$

$$(33)$$

 $\boldsymbol{\mu}_{n} = \boldsymbol{\mu}_{n_{0}} + \left(\sum_{t} \sum_{s} \boldsymbol{\gamma}_{t}(s) (\boldsymbol{I} - \boldsymbol{G}_{s})^{\mathrm{T}} \boldsymbol{\Sigma}_{\boldsymbol{y},s}^{-1} (\boldsymbol{I} - \boldsymbol{G}_{s}) + \boldsymbol{\eta}_{n} \boldsymbol{I}\right)^{-1}$   $\left(\sum_{t} \sum_{s} \boldsymbol{\gamma}_{t}(s) (\boldsymbol{I} - \boldsymbol{G}_{s})^{\mathrm{T}} \boldsymbol{\Sigma}_{\boldsymbol{y},s}^{-1} (\boldsymbol{y}_{t} - \boldsymbol{\mu}_{\boldsymbol{x},s} - \boldsymbol{\mu}_{h_{0}} - g(\boldsymbol{\mu}_{\boldsymbol{x},s}, \boldsymbol{\mu}_{h_{0}}, \boldsymbol{\mu}_{h_{0}}, \boldsymbol{\mu}_{h_{0}})\right)\right)$  (34)

#### 1.3 参数初始值估计

1) 噪声参数初始值估计

假设噪声在一个处理时段内保持平稳,在这一个处 理时段内估计噪声的特征向量与对角协方差矩阵。假设 平稳噪声的短时谱系数符合复高斯分布<sup>[21]</sup>,且短时谱各 频带相互独立,可得短时谱系数的概率密度函数为:

$$p_{N_k}(d) = \frac{1}{2\pi\sigma_k^2} \exp\left(-\frac{|d|^2}{2\sigma_k^2}\right)$$
(35)

式中: $N_k$ 表示短时谱第k个频率采样点处的系数,可以推导得出其幅值 $|N_k|$ 符合瑞利分布:

$$p_{\parallel N_k \parallel}(d) = \frac{d}{\sigma_k^2} \exp\left(-\frac{d^2}{2\sigma_k^2}\right)$$
(36)

因此  $\log |N_k|$ 符合对数瑞利分布<sup>[22]</sup>,其均值与方差 分别表示为 $\mu_k = \sigma_k + \ln 2/2 - \gamma/2, \nu_k = \pi^2/24, \gamma$ 表示欧 拉常数。因此噪声特征的均值初始值和方差可估计为:

$$\boldsymbol{\mu}_{n} = 2\boldsymbol{C}[\boldsymbol{\mu}_{0}, \boldsymbol{\mu}_{1}, \cdots, \boldsymbol{\mu}_{k}, \cdots]^{\mathrm{T}}$$
(37)

$$\boldsymbol{\Sigma}_{n} = \frac{\pi^{2}}{6} \boldsymbol{C} \boldsymbol{I} \boldsymbol{C}^{\mathrm{T}}$$
(38)

因为语音的短时幅度谱具有稀疏性,对于一段带噪 语音的短时幅度谱的每个频带来说,所有时间点的取值 大部分都是该频带的噪声幅值,因此每个频带处带噪语 音分布与噪声分布相近。如图 2 所示是噪声为白噪声, SNR 为 5 dB 时,采样频率为 16 kHz 的带噪语音信号和 对应噪声信号分帧加窗后经 640 点的离散傅里叶变换得 到一系列频谱中第 16 个和第 6 个频率采样点处的幅值 的分布直方图,图中曲线表示均值滤波后得到的分布曲 线。根据式 (36),噪声幅度谱概率分布的最大值在  $\sigma_k$ 处取得,因此  $\sigma_k$  的估计通过直方图统计带噪语音第 k 个 频率采样点短时幅度谱分布,经均值滤波后取分布曲线 最大值处的横坐标实现。



corresponding noisy speech

### 2)信道参数估计

根据式(7)可得:

$$\begin{split} |H_{k}|^{2} |X_{k}|^{2} &= |Y_{k}|^{2} + |N_{k}|^{2} - 2\cos\beta_{k} |Y_{k}| |N_{k}| (40) \\ \text{式中}: \beta_{k} \, \text{表示复变量} \, Y_{k} \, \text{与} \, N_{k} \, \text{之间的夹角}, \text{为了方便起} \\ \text{见}, \text{用定值} \, \beta(0 < \beta < 1) \, \text{代替}, \, \text{所以将} |H_{k}|^{2} \, \text{估计为}: \\ & \text{E}(|H_{k}|^{2}) \approx \end{split}$$

$$\frac{\mathrm{E}(|Y_k|^2) + \mathrm{E}(|N_k|^2) - 2\mathrm{cos}\beta\mathrm{E}(|Y_k|)\mathrm{E}(|N_k|)}{\mathrm{E}(|X_k|^2)}$$
(41)

由于 *X* 的概率密度分布经 GMM 建模,根据式(9)、 式(16)和对数高斯分布的期望值代数表达式可得:

$$\mathbf{E}(|\boldsymbol{X}|^{2}) = \sum_{s} p(s) (\boldsymbol{C}^{-1} \boldsymbol{\mu}_{x,s} + 0.5\boldsymbol{C}^{-1} \boldsymbol{\Sigma}_{x,s} (\boldsymbol{C}^{-1})^{\mathrm{T}})$$
(42)

根据式 (35) 可得 E( $|N_k|^2$ ) =  $2\sigma_k^2$ , E( $|Y_k|^2$ ) 的值 由带噪语音功率谱的均值估计, 最终得到信道特征参数 的初始估计值为:

$$\boldsymbol{\mu}_{h} = \boldsymbol{C}\log(\operatorname{E}(|\boldsymbol{H}|^{2}))$$
(43)

#### 2 实验方案

#### 2.1 算法的实现与评估

本文训练纯净语音 GMM 的训练数据集来自 THCHS-30数据库<sup>[23]</sup>,该数据库语音采样频率为 16 kHz。 用特定说话人的时长 30 min 的语音进行训练,对这段语 音进行分帧加窗并计算 25 阶 MCEP 系数,选取帧长为 512,帧移为 64,窗函数为汉明窗,计算 MCEP 系数时选 取  $\alpha = 0.5$ 。同时,用 PEFAC 算法计算语音帧的浊音概 率,将语音分为浊音部分和清音/无语音段部分,用 Kmeans 算法分别分为 160 个和 40 个簇,然后以这些簇 的中心为 GMM 中每个高斯成分的均值的初始值,借助 Scikit-leam 库训练协方差矩阵为对角阵,高斯分布的聚 类数量为 200 的 GMM。

选用语音信号评价中的客观指标:对数似然比(loglikelihood ratio, LLR)和语音质量感知评估(perceptual evaluation of speech quality, PESQ)<sup>[24]</sup>对算法效果进行评 估。LLR 是一种语音失真的客观评价方法,通过评估纯 净语音信号和处理后语音信号的全极点模型之间的差异 实现,数值越低代表语音质量越好。PESQ 在 2001 年成 为国际电信联盟推出的 P. 862 标准,PESQ 得分与主观 评测相关度较高,能较好得预测语音的主观视听感受,分 值区间为-0.5~4.5,得分越高说明语音质量越好。

#### 2.2 合成信号实验

合成信号实验用向纯净语音添加信道作用和特定信 噪比的噪声形成的信号对算法进行测试。信号合成方法 如图 3 所示。

以不在训练集中但属于同一说话人的语音作为测试



集的纯净语音信号,信道作用的添加方法如下:首先用白 噪声激励,激光测振仪采集的方式测量目标物上某一测 点的频响,将一段测试集纯净语音信号进行傅里叶变换 后,乘以该点的幅频响应,再进行逆傅里叶变换。

用激光多普勒测振仪在测振目标上进行无语音信号的 空采时,在外界环境与测振目标稳定的条件下,观察到采集 到的信号存在散粒噪声——不同强度的白噪声(测量位移 信号时)或幅度谱随频率线性上升的噪声(测量速度信号 时)<sup>[25]</sup>,此外,低频段还明显存在可能来源于目标物振动与 仪器的其它噪声。静止纸盒上一点的空采速度信号和位 移信号如图 4 所示。为了实验方便选取白噪声和白噪声 差分噪声(幅度谱随频率线性上升的噪声)用于合成信号。



Fig. 4 Measurement on a paper box when no speech exists

分别随机选取纸盒、水瓶、窗帘上一个测点的频响添加到随机截取的长度为15s同一说话人未出现在训练集中的语音信号上,再添加白噪声、白噪声差分噪声进行实验,并使得信噪比分别为0、5、10dB,图5所示为添加水瓶频响和白噪声的信噪比为5dB的合成信号与对应纯净语音的对比。用2.1节所述指标进行评价,选用吕韬使用的相位补偿OMLSA算法<sup>[12]</sup>作为对比算法。





#### 2.3 真实实验

采用 Polytec LDV(PSV-500)激光多普勒测振仪实测 扬声器播放的不在训练集中但属于同一说话人的语音信 号引起的不同目标物表面微振动,以此对算法性能进行 验证,如图 6 所示。通过调节测振仪焦距改变回光强度 和调节扬声器音量,可以模拟不同的测振距离和激励强 度得到信噪比不同的观测信号。可以用 PESQ 指标对算 法处理效果进行评估,此外将算法处理后语音与原语音 进行互相关运算,可以实现信号波形对齐,进行波形图、 语谱图和 LLR 指标的对比。





# 3 结果分析

#### 3.1 特征提取方法的验证

用激光多普勒测振仪实际采集的语音信号验证 PEFAC 算法对于浊音概率和基音频率的估计。如图 7 所示为从窗帘上探测到的一段语音信号的浊音概率与基 音频率的估计结果,最上方为浊音概率 pv 的估计;中间 为观测信号波形图,浊音概率大于 0.5 和小于等于 0.5 的部分分别用深色和浅色标注,最下方为观测信号 0~1 000 Hz 频段的语谱图,浊音概率大于 0.5 的部分所 估计的基音频率用实线在图中标出。总得来说,PEFAC 算法在 50~1 000 Hz 的低频噪声能量不过于高的情况下 可以较好地实现激光相干探测语音的浊音概率与基音频 率的估计。

332





观测信号分帧加窗后每帧信号计算 25 阶 MCEP 系数,形成 26 维特征向量,用特征向量恢复该帧信号的对数幅度谱包络与原语音对数幅度谱的对比如图 8 所示, 其中,频率扭曲系数 α 使得 MCEP 系数对于对数幅度谱 低频部分的形状刻画更为精细。



recovered by the MCEP coefficients

#### 3.2 合成信号实验结果分析

合成信号实验由于信道和噪声已知,可以将对信道 和噪声的初始估计和自适应估计迭代后的结果与实际值 进行对比。图9所示为添加纸盒上一点处频响和白噪声 差分噪声,SNR 为5 dB 时,信道和噪声的短时对数功率 谱均值、初始估计μ,和μ,恢复的对数功率谱、噪声和信 道参数经自适应估计分别经迭代 2 次后 μ<sub>h</sub> 和 μ<sub>n</sub> 恢复的 对数功率谱的对比结果。可以看出信道和噪声参数的自 适应估计可以使信道和噪声的估计在一定程度上更接近 于真实值。由于语音能量集中在 30~3 000 Hz,语音能量 过小的频段会存在 1.3 中第 2)部分信道估计方法得到的 估计值偏差过大,经迭代修正也无法很好地估计的情况。



信道和噪声参数经自适应估计后,用1.2节中第2) 部分中所述 MMSE 估计方法估计纯净语音 MCEP 系数。 图 10 展示了一帧带噪语音、其对应纯净语音和本文算法 根据带噪语音估计的 MCEP 系数恢复的对数幅度谱包 络,可见本文算法通过 MCEP 系数的估计能得到与纯净 语音相近的对数幅度谱包络估计。

每种添加了不同信道和噪声的信号取 10 组用本文



算法进行处理,并用 2.1 所述指标进行评价,将评价值取 平均,结果如表 1 所示,每种合成信号每种指标的 3 行数 据从上到下分别为原始合成信号、相位补偿 OMLSA 算法 处理后、本文算法处理后的结果,其中较好的结果经加粗 标注。本文提出算法处理信号的 LLR 指标在所有情况 下较好,这是由于这个算法的机制使得处理后的信号的 短时对数幅度谱包络与原语音信号较为接近。在本文的

Table 1

测试范围内,所提出算法处理信号的 PESQ 指标显示,在 有些低信噪比的情况下,相较相位补偿 OMLSA 算法本文 提出算法有时能更好地提高语音质量;然而,在高信噪比 情况下本文提出算法则效果欠佳,分析原因在于语音合 成算法比较粗糙,且算法各模块存在太多近似环节。因 此,在噪声不至于影响基音频率检测的情况下,本文算法 更适用于较低信噪比时的语音信号质量增强。

评价指标	输入						
	SNR/dB	纸盒+白	纸盒+差分	水瓶+白	水瓶+差分	窗帘+白	窗帘+差分
		1.18	1.53	1.03	1.25	1.32	1.66
	0	1.39	1.98	1.47	1.72	1.60	2.05
		1. 53	1.82	1.48	1.80	1.77	2.03
		1.43	1.72	1.32	1.56	1.64	1.91
PESQ	5	1.93	2. 29	1.73	2.17	2.10	2. 52
		1.75	1.95	1.83	1.90	1.99	2.16
		1.75	1.93	1.56	1.87	1.90	2.17
	10	2.37	2.67	2.14	2.63	2.49	2.91
		1.95	2.07	1.96	2.05	2.09	2.29
LLR		1.94	1.96	1.87	1.87	1.89	1.88
	0	1.92	1.93	1.88	1.89	1.90	1.90
		1. 23	1.03	1.10	1.03	1.09	1.01
		1.90	1.94	1.84	1.86	1.86	1.84
	5	1.88	1.92	1.86	1.88	1.87	1.86
		1.07	1.01	1.10	1.07	1.02	0. 98
		1.88	1.89	1.81	1.85	1.82	1.83
	10	1.85	1.86	1.82	1.87	1.83	1.85
		1.01	1.06	1.04	1.01	1.01	0.97

表1 合成信号实验客观评价指标对比

Objective measures of the proposed method against OMLSA with phase compensation

#### 3.3 真实实验结果

图 11 是测振目标物为纸盒时的一组实际采集信号的波形图和语谱图,将纯净语音信号、观测信号、相位补偿 OLMSA 算法增强后的信号和本文算法增强后的信号





#### 图 11 语音采集实验结果



333

对比,可以看出本文算法的去噪和均衡效果。该组观测 信号、相位补偿 OLMSA 算法增强后的信号和本文算法增 强后的信号 PESQ 分别 1.39、1.63 和 1.79,LLR 分别为 1.88、1.96 和 1.15。

# 4 结 论

为了提高激光相干语音探测的信号质量,本文创新 性地提出了一种基于分析-重合成框架和噪声与信道参 数自适应估计的语音增强算法。通过对观测语音进行特 征提取,结合预训练的目标人物语音谱包络特征参数 GMM 和观测语音中包含的信道与噪声信息,估计观测信 号谱包络特征所对应的纯净语音谱包络特征,进而用纯 净语音谱包络特征的估计值进行语音重合成。合成信号 实验与实际实验证明了本文提出的算法可以在一定程度 上实现激光相干探测语音信号的去噪和均衡。然而该算 法由于语音合成模块比较粗糙,且算法各模块存在太多 近似环节,导致其几乎不可能得到质量很高的语音信号, 且实验表明本文算法与传统语音增强算法一样,并不能 实现语音可懂度的提高,作者将把改善这些问题作为之 后的工作目标。

#### 参考文献

- [1] LYU T, GUO J, ZHANG H Y, et al. Acquirement and enhancement of remote speech signals [J]. Optoelectronics Letters, 2017, 13(4):275-278.
- ZHENG CH SH, ZHANG H Y, LIU W ZH, et al. Sixty years of frequency-domain monaural speech enhancement: From traditional to deep learning methods [J]. Trends in Hearing, 2023, 27:23312165231209913.
- GHORPADE K, KHAPARDE A. Single-channel speech enhancement using single dimension change accelerated particle swarm optimization for subspace partitioning[J]. Circuits, Systems, and Signal Processing, 2023, 42(7): 4343-4361.
- [4] 李如玮,鲍长春,窦慧晶. 基于双正交小波包分解的自适应阈值语音增强[J]. 仪器仪表学报,2008(10): 2135-2140.

LI R W, BAO CH CH, DOU H J. Speech enhancement using adaptive threshold based on bi-orthogonal wavelet packet decomposition [J]. Chinese Journal of Scientific Instrument, 2008(10):2135-2140.

- [5] XUE W, MOORE A H, BROOKES M, et al. Speech enhancement based on modulation-domain parametric multichannel Kalman filtering [J]. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2020, 29: 393-405.
- [6] 朱志宇. 基于高斯粒子滤波器和 TVAR 模型的语音增

强技术[J]. 仪器仪表学报,2008(9):1903-1907.

ZHU ZH Y. Speech enhancement technique based on gaussian particle filter and time-varying autoregressive model [ J ]. Chinese Journal of Scientific Instrument, 2008(9):1903-1907.

- [7] SONG Y J, MADHU N. Investigations on the optimal estimation of speech envelopes for the two-stage speech enhancement[J]. Sensors, 2023, 23(14):23146438.
- [8] KWON K, SHIN J W, KIM N S. NMF-based speech enhancement using bases update [J]. IEEE Signal Processing Letters, 2014, 22(4):450-454.
- [9] XIANG Y, SHI L M, HOJVANG J L, et al. A novel NMF-HMM speech enhancement algorithm based on poisson mixture model [C]. 2021 IEEE International Conference on Acoustics, Speech and Signal Processing. 2021:721-725.
- [10] LIU B, TAO J H, WEN ZH Q, et al. Speech enhancement based on analysis-synthesis framework with improved parameter domain enhancement[J]. Journal of Signal Processing Systems, 2016, 82(2):141-150.
- [11] 袁文浩,屈庆洋,梁春燕,等. 基于感知条件网络的可控语音增强模型[J]. 仪器仪表学报,2023,44(5):53-60.
  YUAN W H, QU Q Y, LIANG CH Y, et al.

Controllable speech enhancement model based on perceptual conditional network [J]. Chinese Journal of Scientific Instrument, 2023,44(5):53-60.

[12] 吕韬. 远距离激光相干语音信号侦测技术研究[D]. 中国科学院大学(中国科学院长春光学精密机械与物 理研究所), 2019.

> LYU T. Research on the remote laser coherent speech signal detection technology [D]. Changchun: University of Chinese Academy of Sciences (Changchun Institute of Optics, Fine Mechanics and Physics, Chinese Academy of Sciences), 2019.

[13] 王永彪,张文喜,王亚慧,等. 拉普拉斯分布下的 MMSE 谱减语音增强算法[J]. 计算机应用,2020, 40(3):878-882.
WANG Y B, ZHANG W X, WANG Y H, et al. Speech enhancement algorithm based on MMSE spectral subtraction with Laplacian distribution [J]. Journal of Computer Applications, 2020,40(3):878-882.

[14] 王亚慧. 远距离激光相干语音探测系统解调及增强技术研究[D]. 北京邮电大学, 2022.
 WANG Y H. Research on demodulation and enhancement technology of long-distance laser coherent speech detection system [D]. Beijing University of Posts and Telecommunications, 2022.

- [15] ERRO D, SAINZ I, NAVAS E, et al. Harmonics plus noise model based vocoder for statistical parametric speech synthesis[J]. IEEE Journal of Selected Topics in Signal Processing, 2013, 8(2):184-194.
- [16] GONZALES S, BROOKES M. PEFAC-a pitch estimation algorithm robust to high levels of noise[J]. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2014, 22(2):518-530.
- YOSHIMURA T, TAKAKI S, NAKAMURA K, et al. Embedding a differentiable mel-cepstral synthesis filter to a neural speech synthesis system [C]. 2023 IEEE International Conference on Acoustics, Speech and Signal Processing, 2023:11816-11820.
- [18] LOWEIMI E, BARKER J, HAIN T. Channel compensation in the generalised vector Taylor series approach to robust ASR [C]. Interspeech. 2017:2466-2470.
- [19] LI J Y, SELTZER M L, GONG Y F. Improvements to VTS feature enhancement [C]. 2012 IEEE International Conference on Acoustics, Speech and Signal Processing, 2012:4677-4680.
- [20] LI J Y, DENG L, YU D, et al. A unified framework of HMM adaptation with joint compensation of additive and convolutive distortions [J]. Computer Speech & Language, 2009, 23(3):389-405.
- [21] NIELSEN J K, CHRISTENSEN M G, BOLDT J B. An analysis of traditional noise power spectral density estimators based on the Gaussian stochastic volatility model[J]. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2023, 31:2299-2313.
- [22] LEE H, LEE M H, YOUN S, et al. Speckle reduction via deep content-aware image prior for precise breast tumor segmentation in an ultrasound image [J]. IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control. 2022, 69(9):2638-2650.
- [23] WANG D, ZHANG X W. THCHS-30: A free Chinese speech corpus[J]. ArXiv preprint arXiv: 1512.01882, 2015.

- [24] RIX A W, BEERENDS J G, HOLLIER M P, et al. Perceptual evaluation of speech quality (PESQ) a new method for speech quality assessment of telephone networks and codecs [C]. 2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings, 2001, 2:749-752.
- [25] JIANG L A, ALBOTA M A, HAUPT R W, et al. Laser vibrometry from a moving ground vehicle [J]. Applied Optics, 2011, 50(15):2263-2273.
- 作者简介



**芮小博**,2021年于天津大学获得博士学 位,现为天津大学副研究员、硕士生导师,研 究方向为声学检测技术。

E-mail:ruixiaobo@tju.edu.cn

**Rui Xiaobo** received his Ph. D. degree from Tianjin University in 2021. He is currently an

Associate Researcher and master supervisor with Tianjin University. His research interest is acoustic detection technology.



**孔欣玥**,2022 年于天津大学获得学士学 位,现为天津大学在读研究生,研究方向为 语音信号处理。

E-mail:kxy7372517@163.com

**Kong Xinyue** received her B. Sc. degree from Tianjin University in 2022. She is

currently a M. Sc. candidate at Tianjin University. Her research interest is speech signal processing.



**伍洲**(通信作者),2019年于中国科学 院大学获得博士学位,现为中国科学院空天 信息创新研究院研究员、博士生导师,研究 方向为光学精密测量技术。

E-mail:wz@aircas.ac.cn

Wu Zhou(Corresponding author) received

her Ph. D. degree from Chinese Academy of Sciences in 2019. She is currently a researcher and Ph. D. supervisor at the Aerospace Information Research Institute, Chinese Academy of Sciences. Her research focus is on optical precision measurement technology.