



# Reverberant Speech Processing for Human Communication and Automatic Speech Recognition

Tomohiro Nakatani, Armin Sehr, Walter Kellermann nakatani.tomohiro@lab.ntt.co.jp, {sehr,wk}@LNT.de

> NTT Communication Science Laboratories LMS, University of Erlangen-Nuremberg

> > March 26, 2012

Mobile users, distant microphones/loudspeakers























Tasks:

- Rendering Reproduce desired signals at distant ears
- Acquisition Localize sources and capture clean signals from distance

# Challenges:

 Feedback of loudspeaker signals



















Tasks:

- Rendering Reproduce desired signals at distant ears
- Acquisition Localize sources and capture clean signals from distance

# Challenges:

- Feedback of loudspeaker signals
- Noise and interferers
- Reverberation





## Hands-free equipment

for telecommunication and natural human/machine interaction

- for mobile phones / smart phones, mobile computing devices, PDAs
- in car interiors ('command&control', telecommunication, in-car communication, ...)
- for desktop computers, info-/edutainment terminals, interactive TV, game stations, simulators
- for telepresence systems (offices,..., classrooms, ..., auditoria)
- for ambient communication (smart meeting rooms, smart homes, information kiosks, museums and exhibitions, ...)
- ► for voice-driven navigation systems in cars, operating rooms, ...







#### Professional Audio

- equipment for stages and recording studios
- virtual acoustic environments (virtual concert halls, telepresence studios,...)





#### Professional Audio

- equipment for stages and recording studios
- virtual acoustic environments (virtual concert halls, telepresence studios,...)

#### Safety and Surveillance

- acoustic displays in control centers, cockpits
- monitoring in health care environments (advanced 'babyphones')
- acoustic scene analysis (train stations, ...)















#### Tasks:

 Rendering -Reproduce undistorted signals with binaural cues









#### Tasks:

- Rendering -Reproduce undistorted signals with binaural cues
- Acquisition Localize desired source(s) and enhance desired signal(s)







#### Tasks:

- Rendering -Reproduce undistorted signals with binaural cues
- Acquisition Localize desired source(s) and enhance desired signal(s)

## **Challenges:**

 Loudspeaker feedback (howling)







#### Tasks:

- Rendering -Reproduce undistorted signals with binaural cues
- Acquisition Localize desired source(s) and enhance desired signal(s)

#### **Challenges:**

• Noise and interferers







#### Tasks:

- Rendering -Reproduce undistorted signals with binaural cues
- Acquisition Localize desired source(s) and enhance desired signal(s)

**Challenges:** 

#### Reverberation







#### Tasks:

- Rendering -Reproduce undistorted signals with binaural cues
- Acquisition Localize desired source(s) and enhance desired signal(s)

# **Challenges:**

- Loudspeaker feedback (howling)
- Noise and interferers
- Reverberation





- Hearing aids, of course
- ► Headsets, e.g., for
  - mobile phones, mobile computing devices, personal digital assistants
  - hearing protection in noisy environments (construction work, mining,...)
  - active noise cancellation systems



. . .





Voice-controlled home entertainment system (EU project DICIT 2005-2009; see e.g., Marquardt et al., 2009; Youtube)







Voice-controlled home entertainment system (**EU project DICIT** 2005-2009; see e.g., Marquardt et al., 2009; **Youtube**)



featuring

- Multichannel AEC (GFDAF, Buchner/Benesty et al., 2003ff)
- Multibeamforming (Mabande et al., 2009; Kellermann, 1997)
- Source localization (GCF; Brutti et al., 2007)
- Speech/non-speech classification (Omologo, 2009)
- Noise-robust automatic speech recognition (ViaVoice, IBM 2009)







Voice-controlled home entertainment system (**EU project DICIT** 2005-2009; see e.g., Marquardt et al., 2009; **Youtube**)



featuring

- Multichannel AEC (GFDAF, Buchner/Benesty et al., 2003ff)
- Multibeamforming (Mabande et al., 2009; Kellermann, 1997)
- Source localization (GCF; Brutti et al., 2007)
- Speech/non-speech classification (Omologo, 2009)
- Noise-robust automatic speech recognition (ViaVoice, IBM 2009)

Challenge: Reverberation for large source distances in more reverberant rooms





# **Example 2: Meeting recognition system**



#### **Real-time Meeting Browser**

和み日





# **Example 3: Audio postproduction system**









**Part I: Introduction** 







# **Part I: Introduction**

- Fundamentals
- Approaches







## **Part I: Introduction**

### Part II: Multichannel blind inverse filtering

Example applications





## **Part I: Introduction**

- Example applications
  - Professional audio post production
  - Meeting speech recognition with microphone arrays





## **Part I: Introduction**

### Part II: Multichannel blind inverse filtering

- Example applications
- Fundamentals: Dereverberation with inverse filtering







## **Part I: Introduction**

#### Part II: Multichannel blind inverse filtering

- Example applications
- Fundamentals: Dereverberation with inverse filtering
  - What is 'inverse' filtering?
  - Robust 'approximate' inverse filtering





# **Part I: Introduction**

- Example applications
- Fundamentals: Dereverberation with inverse filtering
- Blind inverse filtering





# **Part I: Introduction**

- Example applications
- Fundamentals: Dereverberation with inverse filtering
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model







# **Part I: Introduction**

- Example applications
- Fundamentals: Dereverberation with inverse filtering
- Blind inverse filtering
- Integration with blind source separation





Part I: Introduction

#### Part II: Multichannel blind inverse filtering

## Part III: Robust ASR in reverberant environments

Feature-based approaches







#### **Part I: Introduction**

## Part II: Multichannel blind inverse filtering

## Part III: Robust ASR in reverberant environments

- Feature-based approaches
  - Cepstral mean normalization
  - Model-based feature enhancement





**Part I: Introduction** 

## Part II: Multichannel blind inverse filtering

## Part III: Robust ASR in reverberant environments

- Feature-based approaches
- Model-based approaches




#### **Part I: Introduction**

## Part II: Multichannel blind inverse filtering

- Feature-based approaches
- Model-based approaches
  - Matched training
  - Multi-style training
  - Adaptive training
  - MAP and MLLR adaptation
  - Parametric adaptation tailored to reverberation
  - Frame-wise adaptation





#### **Part I: Introduction**

## Part II: Multichannel blind inverse filtering

- Feature-based approaches
- Model-based approaches
- Decoder-based approaches





#### **Part I: Introduction**

## Part II: Multichannel blind inverse filtering

- Feature-based approaches
- Model-based approaches
- Decoder-based approaches
  - Missing feature techniques
  - Uncertainty decoding







#### **Part I: Introduction**

## Part II: Multichannel blind inverse filtering

- Feature-based approaches
- Model-based approaches
- Decoder-based approaches
- A generic approach: REMOS





**Part I: Introduction** 

Part II: Multichannel blind inverse filtering

Part III: Robust ASR in reverberant environments

Part IV: Summary, Conclusions, and Outlook





# **Fundamental Signal Processing Problems - Formulation**







# **Fundamental Signal Processing Problems - Formulation**









# **Fundamental Signal Processing Problems - Formulation**







**Goal:** Undistorted source signals  

$$\mathbf{y} = \mathbf{W}_{\mathbf{yu}} * \mathbf{u} + \mathbf{W}_{\mathbf{yx}} * \mathbf{x} \stackrel{!}{=} \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$$
  
where  $\mathbf{x} = \mathbf{H}_{\mathbf{xs}} * \mathbf{s} + \mathbf{H}_{\mathbf{xv}} * \mathbf{v} + \mathbf{n}_{\mathbf{x}}$ 







**Goal:** Undistorted source signals  $\mathbf{y} = \mathbf{W}_{\mathbf{vu}} * \mathbf{u} + \mathbf{W}_{\mathbf{vx}} * \mathbf{x} \stackrel{!}{=} \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$ where  $\mathbf{x} = \mathbf{H}_{\mathbf{xs}} * \mathbf{s} + \mathbf{H}_{\mathbf{xv}} * \mathbf{v} + \mathbf{n}_{\mathbf{x}}$ 

## **3 Subproblems:**

 Echo cancellation:  $(\mathbf{W}_{\mathbf{v}\mathbf{u}} + \mathbf{W}_{\mathbf{v}\mathbf{x}} * \mathbf{H}_{\mathbf{x}\mathbf{v}} * \mathbf{W}_{\mathbf{v}\mathbf{u}}) * \mathbf{u} = \mathbf{0}$ 







Goal: Undistorted source signals  $\mathbf{y} = \mathbf{W}_{\mathbf{yu}} * \mathbf{u} + \mathbf{W}_{\mathbf{yx}} * \mathbf{x} \stackrel{!}{=} \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$ where  $\mathbf{x} = \mathbf{H}_{\mathbf{xs}} * \mathbf{s} + \mathbf{H}_{\mathbf{xv}} * \mathbf{v} + \mathbf{n}_{\mathbf{x}}$ 3 Subproblems:

 Source separation and dereverberation:

$$\mathbf{W}_{\mathbf{yx}} * \mathbf{H}_{\mathbf{xs}} * \mathbf{s} = \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$$







Goal: Undistorted source signals  $\mathbf{y} = \mathbf{W}_{\mathbf{yu}} * \mathbf{u} + \mathbf{W}_{\mathbf{yx}} * \mathbf{x} \stackrel{!}{=} \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$ where  $\mathbf{x} = \mathbf{H}_{\mathbf{xs}} * \mathbf{s} + \mathbf{H}_{\mathbf{xv}} * \mathbf{v} + \mathbf{n}_{\mathbf{x}}$ 3 Subproblems:

 Noise and interference suppression:

 $W_{yx}\ast n_x=0$ 





**Goal:** Undistorted source signals  $\mathbf{y} = \mathbf{W}_{\mathbf{yu}} * \mathbf{u} + \mathbf{W}_{\mathbf{yx}} * \mathbf{x} \stackrel{!}{=} \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$ 

# where $\mathbf{x} = \mathbf{H}_{\mathbf{xs}} * \mathbf{s} + \mathbf{H}_{\mathbf{xv}} * \mathbf{v} + \mathbf{n}_{\mathbf{x}}$

## 3 Subproblems:

- Echo cancellation:  $(W_{yu} + W_{yx} * H_{xv} * W_{vu}) * u = 0$
- Source separation and dereverberation:

 $\mathbf{W}_{\mathbf{yx}} * \mathbf{H}_{\mathbf{xs}} * \mathbf{s} = \mathbf{s} * \delta(\mathbf{k} - \mathbf{k}_0)$ 

• Noise and interference suppression:

 $\bm{W_{yx}}\ast\bm{n_x}=\bm{0}$ 

**Components of x**, i.e.,  $H_{xs} * s$ ,  $H_{xv} * v$ ,  $n_x$ , must be separated by W!





# Fundamentals - Room Impulse Response (RIR) properties



## Main characteristic parameters:

*T*<sub>60</sub>: Time for exponential decay of envelope by 60dB **DRR**: Direct-to-Reverberant (Energy) Ratio







#### • Reverberation time T<sub>60</sub>

- $\triangleright$  car  $\approx$  50ms
- $\triangleright$  concert halls  $\approx 1 \dots 2s$





#### • Reverberation time $T_{60}$

- $\triangleright$  car  $\approx$  50ms
- $\triangleright$  concert halls  $\approx 1 \dots 2s$

#### FIR models

- ▷ typically  $L_H \approx T_{60} \cdot f_s/3$  coefficients
- nonminimum-phase
- many zeros close to unit circle





14

- Reverberation time  $T_{60}$ 
  - $\triangleright$  car  $\approx$  50ms
  - $\triangleright$  concert halls  $\approx 1 \dots 2s$

#### FIR models

- ▷ typically  $L_H \approx T_{60} \cdot f_s/3$  coefficients
- nonminimum-phase  $\triangleright$
- many zeros close to unit circle  $\triangleright$









**RIRs for varying source-mic distance** ( $d_1 = 1m$  vs.  $d_2 = 4m$ ,  $T_{60} \approx 900m$ s)







15



**RIRs for varying source-mic distance** ( $d_1 = 1m$  vs.  $d_2 = 4m$ ,  $T_{60} \approx 900m$ s)







**RIRs for varying source-mic distance** ( $d_1 = 1m$  vs.  $d_2 = 4m$ ,  $T_{60} \approx 900m$ s)





**RIRs for varying source-mic distance** ( $d_1 = 1m$  vs.  $d_2 = 4m$ ,  $T_{60} \approx 900m$ s)



#### **RIR, DRR** $\Leftrightarrow$ **Reverberation time** $T_{60}$



Nakatani, Sehr, Kellermann: Reverberant Speech Processing



## Variability with displacements:

#### Mic displacement 4.2cm Shift of RIR by 1 sample: (source distance d=4m): RIR 1 RIR 1 x 10<sup>-5</sup> 0.5 h(n) h(n) -00.1 0.2 0.3 0.4 0.6 0.7 0.8 0.9 0 0.5 t in s x 10<sup>-5</sup> RIR 2 -1.5 L 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 h(n) t in s Difference between RIR1 and a RIR1 shifted by 1 sample \_' 0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1 0.5 t in s Difference between RIR1 and RIR2 x 10<sup>-5</sup> h(n) 0.5 -0 h(n) -0.5-1.5<sup>L</sup> 0 -1 0.1 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 0.2 0.3 0.9 0 1 0.4 0.5 0.6 0.7 0.8 1 t in s t in s System error norm: 0.23dB System error norm: 2.56dB







Clean vs. reverberated with  $T_{60} \approx 900 m$ s,  $d_1 = 4m$  and  $T_{60} \approx 3.1$ s,  $d_2 = 5m$ 







# **Fundamentals - Reverberation in ASR features**

Clean vs. reverberated with  $T_{60} \approx 900 ms$ ,  $d_1 = 4m$  and  $T_{60} \approx 3.1s$ ,  $d_2 = 5m$ 









Pauses of  $c_0$  filled!





# **Dereverberation for Speech Enhancement**

**Basic Idea:** Separate speech production from RIR, equalize the latter









# **Dereverberation for Speech Enhancement**

**Basic Idea:** Separate speech production from RIR, equalize the latter



# 'Blind' problem! (no reference signal for RIR input)





**Basic Idea:** Separate speech production from RIR, equalize the latter



# 'Blind' problem! (no reference signal for RIR input)

## **Distinction:**

# **Partial Deconvolution**

(removes reverberation by RIR inversion, ideally without speech distortion)





Basic Idea: Separate speech production from RIR, equalize the latter



# 'Blind' problem!

(no reference signal for RIR input)

## **Distinction:**

# Partial Deconvolution

(removes reverberation by RIR inversion, ideally without speech distortion)

## $\uparrow$

 Reverberation Suppression (compromise between dereverberation and signal distortion necessary)



## Prior Knowledge on

- A, Speech production models (e.g., source-filter model, HMM) and signal properties (nonwhiteness, nonstationarity, nongaussianity)
- B, Room acoustics parameter (e.g.,  $T_{60}$ )
- C, Location and radiation characteristics of speech source





## Prior Knowledge on

- A, Speech production models (e.g., source-filter model, HMM) and signal properties (nonwhiteness, nonstationarity, nongaussianity)
- B, Room acoustics parameter (e.g.,  $T_{60}$ )
- C, Location and radiation characteristics of speech source

## and some Useful Assumptions

D, Joint moments (e.g., correlation) of signal samples: Small lags characterize speech ⇔ Large lags characterize reverberation





## Prior Knowledge on

- A, Speech production models (e.g., source-filter model, HMM) and signal properties (nonwhiteness, nonstationarity, nongaussianity)
- B, Room acoustics parameter (e.g.,  $T_{60}$ )
- C, Location and radiation characteristics of speech source

## and some Useful Assumptions

- D, Joint moments (e.g., correlation) of signal samples: Small lags characterize speech ⇔ Large lags characterize reverberation
- E, Speech signal statistics change faster than RIRs





# Prior Knowledge on

- A, Speech production models (e.g., source-filter model, HMM) and signal properties (nonwhiteness, nonstationarity, nongaussianity)
- B, Room acoustics parameter (e.g.,  $T_{60}$ )
- C, Location and radiation characteristics of speech source

# and some Useful Assumptions

- D, Joint moments (e.g., correlation) of signal samples: Small lags characterize speech ⇔ Large lags characterize reverberation
- E, Speech signal statistics change faster than RIRs
- F, Multichannel recordings: Speech component is the same ⇔ RIRs are different













#### Single-channel partial deconvolution

Can exploit speech models and properties (A) and correlation and stationarity assumptions (D, E) for identifying RIR estimate







## Single-channel partial deconvolution

- Can exploit speech models and properties (A) and correlation and stationarity assumptions (D, E) for identifying RIR estimate
- Inversion of a single RIR involves [Neely 1979]
  - $\blacktriangleright$  removing the allpass component of the nonminimum-phase RIR  $\rightarrow$  approximated by delay
  - ► inverting zeros close to, or on unit circle → approximation by 'channel shortening'





## Single-channel partial deconvolution

- Can exploit speech models and properties (A) and correlation and stationarity assumptions (D, E) for identifying RIR estimate
- Inversion of a single RIR involves [Neely 1979]
  - removing the allpass component of the nonminimum-phase RIR  $\rightarrow$  approximated by delay
  - ► inverting zeros close to, or on unit circle → approximation by 'channel shortening'
- ► for realization problems see, e.g., [Morjopoulos 1994], [Naylor 2010]










- Can additionally exploit spatial diversity (incl. assumption F) and prior knowledge of source location and radiation characteristic (C) for identifying RIR
- Spatial diversity facilitates RIR identification







- Can additionally exploit spatial diversity (incl. assumption F) and prior knowledge of source location and radiation characteristic (C) for identifying RIR
- Spatial diversity facilitates RIR identification
- Perfect inversion with FIR filters is possible (MINT [Miyoshi 1988])
  - exact knowledge of RIR lengths required
  - no common zeros of RIRs allowed





- Can additionally exploit spatial diversity (incl. assumption F) and prior knowledge of source location and radiation characteristic (C) for identifying RIR
- Spatial diversity facilitates RIR identification
- Perfect inversion with FIR filters is possible (MINT [Miyoshi 1988])
  - exact knowledge of RIR lengths required
  - no common zeros of RIRs allowed
- Indirect approaches often invert in subbands for robustness (e.g. [Naylor 2005])







- Can additionally exploit spatial diversity (incl. assumption F) and prior knowledge of source location and radiation characteristic (C) for identifying RIR
- Spatial diversity facilitates RIR identification
- Perfect inversion with FIR filters is possible (MINT [Miyoshi 1988])
  - exact knowledge of RIR lengths required
  - no common zeros of RIRs allowed
- Indirect approaches often invert in subbands for robustness (e.g. [Naylor] 2005])
- Direct approaches to identify a robust inverse exist (e.g. [Buchner 2004], [Buchner 2010], and below!)





#### **Single-channel Reverberation Suppression**







#### Single-channel Reverberation Suppression

- can exploit speech models and properties (A) and correlation and stationarity assumptions (D, E), e.g., for
  - equalizing the vocal tract IR and suppressing reverberation in the LPC residual (e.g., [Yegnanarayana 2000], [Gaubitch 2006])





#### Single-channel Reverberation Suppression

- can exploit speech models and properties (A) and correlation and stationarity assumptions (D, E), e.g., for
  - equalizing the vocal tract IR and suppressing reverberation in the LPC residual (e.g., [Yegnanarayana 2000], [Gaubitch 2006])
- can exploit prior knowledge on room acoustics (e.g., T<sub>60</sub>) to estimate PSD of reverberation and use spectral subtraction methods as common for additive noise (e.g., [Lebart 2001])











can additionally exploit spatial diversity (incl. assumption F) and prior knowledge on source location and radiation characteristic (C), e.g.,







- can additionally exploit spatial diversity (incl. assumption F) and prior knowledge on source location and radiation characteristic (C), e.g.,
  - beamforming using only prior knowledge of source location and radiation characteristic (C) (e.g., [Griebel 2001])





- can additionally exploit spatial diversity (incl. assumption F) and prior knowledge on source location and radiation characteristic (C), e.g.,
  - beamforming using only prior knowledge of source location and radiation characteristic (C) (e.g., [Griebel 2001])
  - spatial diversity for multichannel spectral subtraction (e.g., [Allen 1977]), or subspace methods (e.g., [Gannot 2003])





- can additionally exploit spatial diversity (incl. assumption F) and prior knowledge on source location and radiation characteristic (C), e.g.,
  - beamforming using only prior knowledge of source location and radiation characteristic (C) (e.g., [Griebel 2001])
  - spatial diversity for multichannel spectral subtraction (e.g., [Allen 1977]), or subspace methods (e.g., [Gannot 2003])
  - spatial diversity complemented by prior knowledge on room acoustics parameter (e.g., [Habets 2005])





#### **Block diagram of ASR system**







#### **Block diagram of ASR system**



#### **Strategies**

A) signal-based approaches







#### **Block diagram of ASR system**



#### **Strategies**

- A) signal-based approaches
- B) feature-based approaches





#### **Block diagram of ASR system**



#### **Strategies**

A) signal-based approaches

C) model-based approaches

B) feature-based approaches





#### **Block diagram of ASR system**



#### **Strategies**

- A) signal-based approaches
- B) feature-based approaches
- C) model-based approaches
- D) decoder-based approaches







# Part II. Multichannel blind inverse filtering







### Two approaches for signal dereverberation





## **Multichannel inverse filtering**



Clinear filtering: 
$$y_t = \sum_{m=1}^{\infty} \sum_{\tau=0}^{\infty} w_{\tau}^{(m)} x_{t-\tau}^{(m)}$$
  
Goal: estimate  $\{w_t^{(m)}\}$  s.t.  $y_t = s_t$ 

*t* : time index *t* : time index {·} : a set of variables for all *t* and *m* 





- Example applications
  - Professional audio post production
  - Meeting recognition with microphone arrays
- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model
- Integration with blind source separation













#### **Dereverberation system for audio post production** [Kinoshita 2008]

 Dereverberation plug-in for Pro Tools: NML RevCon-RR (sold by TAC System, Inc.)







## **Online meeting recognition [Hori 2012]**



#### **Real-time Meeting Browser**





## **Online/offline processing flow of meeting recognition**



## ASR performance w/ and w/o dereverberation



Trained on CSJ (corpus of spontaneous Japanese): headset recording Language model:

Vocabulary size: 156K (LVCSR)



### **Questions to be answered**



- What is inverse filtering ?
- Is the inverse filter robust against interferences ?
- Can we estimate the inverse filter with blind processing ?





What is inverse filtering?

Inversion of room impulse responses (RIRs)

Is the inverse filter robust against interferences ?



but there is a robust 'approximate' inverse filter

Can we estimate the inverse filter with blind processing ?

Yes, we can,

by using cues for distinguishing speech from RIRs





## Part II. Multichannel blind inverse filtering

- Example applications
  - Professional audio post product
  - Meeting recognition with microp

Assume *non-blind processing* for analysis purpose

- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model
- Integration with blind source separation





### Inversion of RIRs = Inversion of matrix transformation







### Matrix/vector representations of RIR convolution/filtering







### **Existence of inverse filter**

$$S_t \longrightarrow \mathbf{X}_t = \mathbf{H}\mathbf{S}_t \xrightarrow{\mathbf{X}_t} y_t = \mathbf{W}^T \mathbf{X}_t \longrightarrow \mathcal{Y}_t$$

• A column vector  ${\bf W}$  is an inverse filter when it satisfies:

$$s_{t} = y_{t} = \mathbf{W}^{T} \mathbf{H} \mathbf{S}_{t} \text{ where } \mathbf{s}_{t} = [s_{t}, s_{t-1}, \cdots, s_{t-K_{0}}]^{T}$$
$$\mathbf{W}^{T} \mathbf{H} = \mathbf{e}^{T} \text{ where } \mathbf{e} = [1, 0, \cdots, 0]^{T}$$

• An inverse filter w exists, when H is invertible, i.e., it is full column rank, and w is obtained as

$$\mathbf{w}^T = \mathbf{e}^T \mathbf{H}^+$$
 where  $\mathbf{H}^+ = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$ 





• H is invertible, or full column rank, if and only if

(#rows of H) >= (#columns of H)

and all columns are linearly independent

 In the case of single source, (#rows of H) >= (#columns of H) is satisfied if and only if

## *M* (#mics) > 1





Multiple-input/output inverse theorem (MINT) [Miyoshi 1988]



– Inverse filter exists when **H** is full column rank



- M (#mics) >N (#sources)
- H(z) does not contain common zeros





## Part II. Multichannel blind inverse filtering

- Example applications
  - Professional audio post production
  - Meeting recognition with microphone arrays
- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model
- Integration with blind source separation




Assumptions for inverse filtering

- Invertible RIRs
- No additive noise
- Time-invariant RIRs

Not realistic !

Inverse filter is too sensitive to modeling errors (noise or RIR change)



Noise-free reverberant case

- Clean speech
- Reverberant speech

   Synthesized using a fixed RIR (RT60=0.5 s)
- Dereverberated speech using an inverse filter for known RIRs (2-channel)

Noisy reverberant case

- Noisy reverberant speech (SNR=30dB)





## Why inverse filter is so sensitive to additive noise







## Standard numerical approach for robustness [Engl 1996]

- Regularization Noisy rev. 🌾 Processed 🌾
  - A general technique for robust matrix inversion
    - Add a very small positive constant  $\delta$  to diagonal of  $\mathbf{H}^T \mathbf{H}$  for calculating the pseudo-inversion of  $\mathbf{H}$

I : Identity matrix

– It can reduce the maximum singular value  $\,\mathcal{X}^{\rm inv}_{max}\,$  of  $\,\widetilde{H}^{\,+}\,$ 

# Noise amplification is greatly mitigated



48



## **Room acoustics motivated approach for robustness**

- Channel shortening Noisy rev. 4 Processed 4
  - Set "direct signal + early reflections" as target signal, and reduce only late reverberation



#### Target to reverberation ratio (TRR) w/ channel shortening is much higher than TRR w/ inverse filtering







• Dereverberation: inversion of RIRs

– Assuming RIRs to be a time-invariant linear system

- Inverse filter exists
  - When we have more microphones than sources
  - But it may be very sensitive to additive noise
- 'Approximate' inverse filter is robust against noise
  - Based on regularization and channel shortening





## Part II. Multichannel blind inverse filtering

- Example applications
  - Professional audio post production
  - Meeting recognition with microphone arrays
- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model
- Integration with blind source separation





## Blind inverse filtering based dereverberation







#### Conventional decorrelation approaches for stationary white signal



- SOS approach assumes *S<sub>t</sub>* to be stationary white Gaussian *Multichannel linear prediction (MCLP)* [Slock 1994], [Abed-Meraim 1997]
- HOS approach assumes  $S_t$  to be an i.i.d. sequence Higher order decorrelation [Sato 1975], [Bellini 1994]





## **Multichannel linear prediction (MCLP)**







## MCLP based decorrelation [Slock 1994], [Abed-Meraim 1997]

• 
$$x_t^{(1)}$$
 is modeled by  
 $x_t^{(1)} = \sum_{\substack{m=1 \ \tau=1}}^M \sum_{\tau=1}^K c_{\tau}^{(m)} x_{t-\tau}^{(m)} + S_t = \mathbf{x}_{t-1}^T \mathbf{c} + S_t \implies S_t = x_t^{(1)} - \mathbf{x}_{t-1}^T \mathbf{c}$   
Predicted signal (= reverberation)  
Prediction error (= direct signal)

where  $\mathbf{c} = [c_1^{(1)}, \dots, c_K^{(1)}, \dots, c_1^{(M)}, \dots, c_K^{(M)}]^T$  is prediction coeffs. - **c** is equivalent to inverse filter **W** 

- **c** can be estimated by minimizing prediction error when sources are stationary and uncorrelated in time  $\hat{\mathbf{c}} = \operatorname{argmin} \sum (x_t^{(1)} - \mathbf{x}_{t-1}^T \mathbf{c})^2$ 
  - Quadratic form: optimized using a closed form solution



### Why dereverberation can be achieved by MCLP



ͲͲ

Let  $Z_t^{(m)} = x_t^{(m)} + \underline{n}_t^{(m)}$  be noisy reverberant observation.

Additive noise (or can be viewed as modeling error) Assume  $x_t^{(m)}$  and  $n_t^{(m)}$  to be uncorrelated, then, the cost function becomes

$$\sum_{t=1}^{T} \left( z_t^{(1)} - \mathbf{z}_{t-1}^T \mathbf{c} \right)^2 = \sum_{t=1}^{T} \left( x_t^{(1)} - \mathbf{x}_{t-1}^T \mathbf{c} \right)^2 + \sum_{t=1}^{T} \left( n_t^{(1)} - \mathbf{n}_{t-1}^T \mathbf{c} \right)^2$$

Cost function for dereverberation

Cost function for noise amplification







#### Problem of decorrelation approach for speech dereverberation







Cues	Speech	RIRs
Inter-channel difference	<b>Common</b> to all the microphone signals	Different for each microphone
Auto-correlation duration	Correlated only within short time interval of the order of <b>30 ms</b>	Correlated within long time interval over 100 ms
Nonstationarity	Stationary only within short time period of the order of <b>30 ms</b>	Stationary over long time period of the order of <b>1000 ms or</b> larger





## Approaches to blind inverse filtering

# • Subspace method (RIR estimation + inversion) -[Furuya 1997], [Gannot 2003], [Gaubitch 2006] Pre-whitening + decorrelation -Second-order statistics (SOS): [Gaubitch 2003], [Furuya 2007], [Triki 2007] -Higher-order statistics (HOS): [Gillespie 2001] Channel shortening -[Gillespie 2003], [Kinoshita 2009] Joint speech and reverberation modeling -[Hopgood 2003], [Buchner (TRINICON) 2010], [Yoshioka 2007], [Nakatani 2008]

#### <u>Cues</u>

Inter-channel difference

Auto-correlation duration

Auto-correlation duration and nonstationarity



## **Pre-whitening + decorrelation**



- A typical method for pre-whitening
  - -Low-dimensional (e.g., 12-dim) single channel linear prediction often used

**Assumption:** pre-whitening can decorrelate only  $\mathbf{s}_t$  in  $\mathbf{x}_t = \mathbf{H}\mathbf{s}_t$ , and we can obtain  $\widetilde{\mathbf{x}}_t = \mathbf{H}\widetilde{\mathbf{s}}_t$  where  $\widetilde{\mathbf{s}}_t$  is an unknown decorrelated speech



61



Introduce constraints so that dereverberation  $\bullet$ reduces only late reverberation



Make derev. *robust* and *do not decorrelate speech* 

- Techniques:
  - Correlation shaping [Gillespie 2003]
  - Multistep MCLP [Kinoshita 2009]





## Multistep MCLP [Gesbert 1997], [Kinoshita 2009]







## Approaches to blind inverse filtering







## Joint speech and reverberation modeling for derev.

#### Model of generative system











#### **Multichannel blind partial deconvolution (MCBPD) by TRINICON**







## Part II. Multichannel blind inverse filtering

- Example applications
  - Professional audio post production
  - Meeting recognition with microphone arrays
- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model
- Integration with blind source separation



68









## **Reformulation of MCLP based on likelihood maximization**

$$L(\mathbf{c}) = \log p_x(\{\mathbf{x}_t\}; \mathbf{c})$$

$$= \sum_{t=1} \log p_x(x_t^{(1)} | \{x_{t'}\}_{t'=1:t-1}; \theta) + const.$$

$$= \sum_{t=1} \log p_s(s_t) + const. \text{ where } s_t = x_t^{(1)} - \mathbf{x}_{t-1}^T \mathbf{c}$$

$$x_t^{(1)} = \mathbf{x}_{t-1}^T \mathbf{c} + s_t$$
Source model

Assume  $p_s(s_t) = N(s_t; 0, 1)$  (stationary white Gaussian), then

$$L(\mathbf{c}) = -(1/2)\sum_{t=1} |x_t^{(1)} - \mathbf{c}^T \mathbf{x}_{t-1}|^2 + const.$$

 $\begin{pmatrix} \max_{\mathbf{c}} L(\mathbf{c}) & \min_{\mathbf{c}} \sum_{t=1}^{\infty} |x_t^{(1)} - \mathbf{c}^T \mathbf{x}_{t-1}|^2 \\ \text{Maximize likelihood} & \text{Minimize prediction error} \end{pmatrix}$ 



## Time-varying Gaussian source model (TVGSM)

30 ms order

1. Each short time segment  $\overline{\mathbf{s}}_{t} = [s_{t} s_{t-1} \cdots s_{t-N+1}]^{T}$  is **stationary multivariate Gaussian**, which can be characterized by

$$p_s(\overline{\mathbf{s}}_t;\mathbf{R}_t) = \mathcal{N}(\overline{\mathbf{s}}_t;\mathbf{0},\mathbf{R}_t)$$

where  $\mathbf{R}_{t} = E\left\{ \overline{\mathbf{s}}_{t} \overline{\mathbf{s}}_{t}^{T} \right\}$  is an autocorrelation matrix

2.  $\mathbf{R}_{t}$  varies over different time segments

$$\theta_s = \{\mathbf{R}_t\}$$
: parameters to be estimated





### MCLP with multivariate source model

$$\mathbf{x}_{t}^{(1)} = \mathbf{x}_{t-1}^{T}\mathbf{c} + s_{t}$$

$$\mathbf{\overline{x}}_{t}^{(1)} = \mathbf{X}_{t-1}\mathbf{c} + \mathbf{\overline{s}}_{t} \text{ or } \begin{bmatrix} \mathbf{x}_{t}^{(1)} \\ \mathbf{x}_{t-1}^{(1)} \\ \vdots \\ \mathbf{x}_{t-1}^{(1)} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_{t-1}^{T} \\ \mathbf{x}_{t-2}^{T} \\ \vdots \\ \mathbf{x}_{t-2}^{T} \end{bmatrix} \mathbf{c} + \begin{bmatrix} s_{t-1} \\ s_{t-1} \\ \vdots \\ \mathbf{x}_{t-1}^{T} \end{bmatrix} \xrightarrow{\mathbf{a}}_{t} \xrightarrow{\mathbf{a}}_{t} \xrightarrow{\mathbf{a}}_{t}$$

$$30 \text{ ms}_{t} \text{ order}$$

$$\mathbf{\overline{x}}_{t-1}^{(1)} \mathbf{\overline{x}}_{t-1} \xrightarrow{\mathbf{\overline{s}}}_{t}$$

• Prediction error  $\overline{\mathbf{S}}_{t}$  is assumed to follow TVGSM



### Likelihood function of MCLP with TVGSM

$$L(\mathbf{c}, \{\mathbf{R}_t\}) = \sum_{t} \log p_s(\overline{\mathbf{s}}_t; \mathbf{c}, \{\mathbf{R}_t\})$$
  
where  $\overline{\mathbf{s}}_t = \overline{\mathbf{x}}_t^{(1)} - \mathbf{X}_{t-1}\mathbf{c}$  and  $p_s(\overline{\mathbf{s}}_t; \mathbf{R}_t) = \mathcal{N}(\overline{\mathbf{s}}_t; 0, \mathbf{R}_t)$ 

$$L(\mathbf{c}, \{\mathbf{R}_t\}) = -\sum_{t} \| \overline{\mathbf{x}}_{t}^{(1)} - \mathbf{X}_{t-1}\mathbf{c} \|_{\mathbf{R}_t} - \frac{\log \| \mathbf{R}_t \|_{\mathbf{R}_t}}{\text{Normalization term}}$$
  
Prediction error  
weighted by  $\mathbf{R}_t^{-1}$   
where  $\| \overline{\mathbf{s}} \|_{\mathbf{R}} = \overline{\mathbf{s}}^T \mathbf{R}^{-1} \overline{\mathbf{s}}$  (quadratic form)





## Iterative optimization procedure







## Importance of time-varying source model





## Blind inverse filtering works in noisy environments



\*TRR: Target-to-reverberation ratio (target = direct signal + early reflections)





## **Computationally efficient implementation**

 Subband decomposition approach [Nakatani 2010], [Yoshioka 2009b]



Computational efficiency largely improves

Real-time factor (RTF)	Time-domain	Subband
using MATLAB	470	
(RT60: 0.5 s, # mics: 2)	170	<b>0.8</b>



## Processing flow with subband decomposition [Nakatani 2010]

- 1.Set analysis parameters
- *D* : prediction delay (*D* should be # of subband samples corresponding to 30 ms, or larger)
- L: length of prediction filter, M: # of mics,
- $m_0$ : index of target channel to be dereverberated
- $\alpha$  : a coeff. for flooring constant (e.g.,  $\alpha$  =  $10^{-4}$  )
- 2.Decompose a multichannel observed signal into a set of subband signals
  - $\chi_{n,f}^{(m)}$ : subband signal (e.g., [Weiss 2000], or STFT can also be used)
    - *m* : channel index, *n* : sample index
    - $\boldsymbol{f}$  : subband index
  - E.g., # of subbands is 512 (including negative frequencies) for 16 kHz sampling
- 3.In each subband f, set initial estimates of source variance  $\sigma_{\mathbf{n},f}$  as

$$\varepsilon_k = \alpha \max_n |x_{n,f}^{(m_0)}|^2$$
$$\sigma_{n,f} = \max\left\{|x_{n,f}^{(m_0)}|^2, \varepsilon_f\right\}$$

where  $\mathcal{E}_{f}$  is a flooring constant for  $\sigma_{n,f}$ 

4.Obtain vector representation of  $\boldsymbol{x}_{n,f}^{(m)}$  in all channels as  $\mathbf{x}_{n,f} = [\mathbf{x}_{n,f}^{(1)^{T}}, \mathbf{x}_{n,f}^{(2)^{T}}, \cdots \mathbf{x}_{n,f}^{(M)^{T}}]^{T}$ where *T* is non-conjugate transposition, and

$$\mathbf{x}_{n,f}^{(m)} = [x_{n,f}^{(m)}, x_{n-1,f}^{(m)}, \cdots x_{n-L+1,f}^{(m)}]^T$$

- 5.In each subband *f*, iterate the following until convergence is achieved
  - i. Obtain prediction filter  $\mathbf{c}_f$  as

$$\mathbf{c}_{f} = \left(\sum_{n} \frac{\mathbf{x}_{n-D,f} \mathbf{x}_{n-D,f}^{*T}}{\sigma_{n,f}}\right)^{+} \sum_{n} \frac{\mathbf{x}_{n-D,f} \left(\mathbf{x}_{n,f}^{(m_{0})}\right)^{*}}{\sigma_{n,f}}$$

where+ and\* are Moore-Penrose pseudoinverse and complex conjugate operations. (see [Yoshioka, 2009b] for efficient calculation)

ii. Obtain dereverberated subband signal  $\mathbf{y}_{n,f}$  as

$$\mathbf{y}_{n,f} = x_{n,f}^{(m_0)} - \mathbf{c}_f^{*T} \mathbf{x}_{n-D,f}$$

iii. Update source variance estimates  $\sigma_{n,f}$  as  $\sigma_{n,f} = \max \{ |y_{n,f}|^2, \varepsilon_f \}$ 

6.Compose a dereverberated signal from a set of dereverberated subband signals  $\mathbf{y}_{n,f}$ 



## Part II. Multichannel blind inverse filtering

- Example applications
  - Professional audio post production
  - Meeting recognition with microphone arrays
- Fundamentals: dereverberation with inverse filtering
  - What is inverse filter
  - Robust 'approximate' inverse filter
- Blind inverse filtering
  - Overview of basic approaches
  - Closer look: multichannel linear prediction with time-varying source model

## - Integration with blind source separation



79

### **BSS+dereverberation**



### Approaches:

TRINICON [Buchner 2010]

- MCLP based approach [Yoshioka 2009b, 2011]
- $\bigcirc$


## Generative model for reverberant sound mixture



 $x_{t}^{(m)}$ : reverberant mixture

# Jointly optimized by maximum likelihood estimation approach [Yoshioka 2009b, 2011]





## **Optimization procedure (subband-based implementation)**







## Improvement in signal-to-interference ratio (SIR)



Results averaged over 672 pairs of utterances (TIMIT test set)







# Live demo











### **TRINICON:** general framework for blind MIMO signal processing



Cost function [Buchner 2010]

$$\mathcal{J}(\mathbf{W}_b) = -\sum_{i=0}^{\infty} \beta(i,b) \sum_{j=0}^{N_0} \left\{ \log(\hat{p}_{s,PD}(\mathbf{y}(i,j))) - \log(\hat{p}_{y,PD}(\mathbf{y}(i,j))) \right\}$$

with *PD*-variate pdfs (*P*: source number, *D*: filter length)

- $\hat{p}_{s,PD}(\mathbf{y}(i,j))$  for source (assumed or estimated)
- $\hat{p}_{y,PD}(\mathbf{y}(i,j))$  for output





# Comparison of SOS and HOS by TRINICON [Buchner 2010]



# mics.: 4, # sources: 2,  $T_{60}$ : 700 ms, Source-mic distance: 1.65 m, Recording: 30 sec





## Summary II-2

- Robust blind inverse filtering is possible
  - Using joint speech and reverberation modeling
    - Based only on a few seconds of observation (e.g., 2.5 s)
    - With a relatively small computational cost (e.g., RTF<1)
    - In an online processing manner (e.g., latency=1s)
  - Under low SNR conditions (e.g., 10 dB SNR)
- Future challenges
  - Realtime adaptation of inverse filter [Yoshioka 2009a], [Evers 2011]
  - Single channel inverse filtering [Gillespie 2001]
  - Processing under more adverse noise conditions such as nonstationary diffuse noise
  - Optimal integration of inverse filtering and spectral enhancement based dereverberation



87



#### Part III:

#### Robust Automatic Speech Recognition (ASR) in Reverberant Environments







#### Introduction

- Feature-based Approaches
- Model-based Approaches
- Decoder-based Approaches

#### ▶ REMOS





89



#### Introduction

Feature-based Approaches

Model-based Approaches

Decoder-based Approaches







#### **ASR System**

speech presignal processing feature extraction acoustic language model model recognition recog-



**Block Diagram** 





**Goal:** Dimensionality reduction









**Goal:** Dimensionality reduction









**Goal:** Dimensionality reduction

MFCCs:
--------







Goal: Dimensionality reduction

MFCCs:







**Goal:** Dimensionality reduction







#### **Acoustic Modeling**

#### Hidden Markov Model (HMM) $\lambda$







#### **Acoustic Modeling**

#### Hidden Markov Model (HMM) $\lambda$



Powerful model for:

- temporal variation
- spectral variation







#### **Dispersive Effect of Reverberation**



- Logmelspec features, dB scale
- Dispersive effect of reverberation: features smeared along time axis



#### **Dispersive Effect of Reverberation**



- Time-frequency pattern is changed
- Inter-frame correlation is increased





#### **Dispersive Effect of Reverberation**



Different statistical properties to be captured by acoustic model

Contradiction to conditional independence assumption of HMMs





#### **Explanation of Dispersive Effect**



Time-domain (TD) description of reverberant speech x<sub>t</sub>:

 $x_t = h_t * s_t$ 

RIR typically much longer than analysis window





#### **Explanation of Dispersive Effect**



► Time-domain (TD) description of reverberant speech *x*<sub>t</sub>:

 $x_t = h_t * s_t$ 

RIR typically much longer than analysis window

► Feature-domain (FD) description of *x*<sup>MEL</sup>: melspec convolution

$$\boldsymbol{x}_{n}^{\mathrm{MEL}} = \sum_{\tau=0}^{T_{H}-1} \boldsymbol{h}_{\tau}^{\mathrm{MEL}} \odot \boldsymbol{s}_{n-\tau}^{\mathrm{MEL}}$$

- $m{s}_n^{ ext{MEL}}$ : clo  $m{x}_n^{ ext{MEL}}$ : re  $m{h}_n^{ ext{MEL}}$ : m  $\odot$ : el
  - <sup>2L</sup>: clean-speech feature vector <sup>2L</sup>: reverberant feature vector
  - L: melspec RIR representation element-wise multiplication





#### **Illustration of Melspec Convolution**









#### **Illustration of Melspec Convolution**









#### Illustration of Melspec Convolution











#### **Accuracy of Melspec Convolution**



a) Clean utterance

b) Reverberant utterance

c) Melspec convolution

d) Simple multiplication





#### **Statistical Properties of Reverberant Speech Features**

#### Example: digit "seven"







#### **Statistical Properties of Reverberant Speech Features**

#### **Histograms**







#### **Statistical Properties of Reverberant Speech Features**

#### **Histograms**





Auto-CoVariances (ACVs)



#### Word Accuracy as Function of Reverberation Time



- Task: Read sentences from Wall Street Journal (WSJ 5K task)
- Features: MFCCs
  + Δ + ΔΔ coefficients
- Recognizer: Cross-word triphones, 3 states per triphone, 16 Gaussians per state







[Sehr 2010a]

- Task: Connected digits (TI digits)
- ► Features: MFCCs + ∆ coefficients
- Recognizer: Word-level HMMs, 16 states per digit, 3 Gaussians per state





#### **Strategies**







#### **Strategies**

A) signal-based approaches







#### Strategies

- A) signal-based approaches
- B) feature-based approaches







#### Strategies

A) signal-based approaches

- C) model-based approaches
- B) feature-based approaches




# Strategies for Reverberation-Robust ASR



### Strategies

- A) signal-based approaches
- B) feature-based approaches
- C) model-based approaches
- D) decoder-based approaches





## Strategies for Reverberation-Robust ASR



#### Strategies

- A) signal-based approaches
- B) feature-based approaches
- C) model-based approaches
- D) decoder-based approaches





## Introduction

## Feature-based Approaches

Model-based Approaches

Decoder-based Approaches







### **Three Different Approaches**

- Feature compensation
  - $\Rightarrow$  Example: Cepstral mean normalization (CMN)





### **Three Different Approaches**

- Feature compensation
  - $\Rightarrow$  Example: Cepstral mean normalization (CMN)
- Features insensitive to reverberation
  - $\Rightarrow$  Example: RASTA features





### **Three Different Approaches**

- Feature compensation
  - $\Rightarrow$  Example: Cepstral mean normalization (CMN)
- Features insensitive to reverberation
  - $\Rightarrow$  Example: RASTA features
- Features facilitating the capture of statistical properties
  - $\Rightarrow$  Example: Dynamic features





If impulse response much shorter than STFT analysis window

$$x_t = h_t * s_t$$





If impulse response much shorter than STFT analysis window

$$egin{array}{rcl} X_t &=& n_t * m{s}_t \ X_{n.k}^{ ext{STFT}} |^2 &pprox & |m{H}_k^{ ext{STFT}}|^2 \, |m{S}_{n.k}^{ ext{STFT}}|^2 \end{array}$$





If impulse response much shorter than STFT analysis window

$$egin{array}{rcl} x_t &=& h_t * s_t \ X_{n,k}^{
m STFT}|^2 &pprox & |H_k^{
m STFT}|^2 \; |S_{n,k}^{
m STFT}|^2 \ x_{n,c}^{
m MFCC} &pprox & h_c^{
m MFCC} + s_{n,c}^{
m MFCC} \end{array}$$





If impulse response much shorter than STFT analysis window

$$\begin{array}{lll} x_t &=& h_t \ast s_t \\ |X_{n,k}^{\rm STFT}|^2 &\approx& |H_k^{\rm STFT}|^2 \; |S_{n,k}^{\rm STFT}|^2 \\ x_{n,c}^{\rm MFCC} &\approx& h_c^{\rm MFCC} + s_{n,c}^{\rm MFCC} \end{array}$$

$$x_{n,c}^{\text{CMN}} = x_{n,c}^{\text{MFCC}} - \bar{x}_{c}^{\text{MFCC}}$$





If impulse response much shorter than STFT analysis window

$$egin{array}{rcl} x_t &=& h_t * s_t \ |X_{n,k}^{ ext{STFT}}|^2 &pprox & |H_k^{ ext{STFT}}|^2 \ |S_{n,k}^{ ext{STFT}}|^2 \ x_{n,c}^{ ext{MFCC}} &pprox & h_c^{ ext{MFCC}} + s_{n,c}^{ ext{MFCC}} \end{array}$$

$$\begin{split} x_{n,c}^{\text{CMN}} &= x_{n,c}^{\text{MFCC}} - \bar{x}_{c}^{\text{MFCC}} \\ \bar{x}_{c}^{\text{MFCC}} &= \frac{1}{N} \sum_{n=1}^{N} x_{n,c}^{\text{MFCC}} \approx h_{c}^{\text{MFCC}} + \bar{s}_{c}^{\text{MFCC}} \end{split}$$





If impulse response much shorter than STFT analysis window

$$egin{array}{rcl} x_t &=& h_t * s_t \ |X_{n,k}^{ ext{STFT}}|^2 &pprox & |H_k^{ ext{STFT}}|^2 \; |S_{n,k}^{ ext{STFT}}|^2 \ x_{n,c}^{ ext{MFCC}} &pprox & h_c^{ ext{MFCC}} + s_{n,c}^{ ext{MFCC}} \end{array}$$

$$\begin{split} x_{n,c}^{\text{CMN}} &= x_{n,c}^{\text{MFCC}} - \bar{x}_{c}^{\text{MFCC}} \\ \bar{x}_{c}^{\text{MFCC}} &= \frac{1}{N} \sum_{n=1}^{N} x_{n,c}^{\text{MFCC}} \approx h_{c}^{\text{MFCC}} + \bar{s}_{c}^{\text{MFCC}} \\ x_{n,c}^{\text{CMN}} &\approx h_{c}^{\text{MFCC}} + s_{n,c}^{\text{MFCC}} - (h_{c}^{\text{MFCC}} + \bar{s}_{c}^{\text{MFCC}}) = s_{n,c}^{\text{MFCC}} - \bar{s}_{c}^{\text{MFCC}} = s_{n,c}^{\text{CMN}} \end{split}$$





If impulse response much shorter than STFT analysis window

$$egin{array}{rcl} x_t &=& h_t * s_t \ |X_{n,k}^{ ext{STFT}}|^2 &pprox & |H_k^{ ext{STFT}}|^2 \ |S_{n,k}^{ ext{STFT}}|^2 \ x_{n,c}^{ ext{MFCC}} &pprox & h_c^{ ext{MFCC}} + s_{n,c}^{ ext{MFCC}} \end{array}$$

$$\begin{aligned} x_{n,c}^{\text{CMN}} &= x_{n,c}^{\text{MFCC}} - \bar{x}_{c}^{\text{MFCC}} \\ \bar{x}_{c}^{\text{MFCC}} &= \frac{1}{N} \sum_{n=1}^{N} x_{n,c}^{\text{MFCC}} \approx h_{c}^{\text{MFCC}} + \bar{s}_{c}^{\text{MFCC}} \\ x_{n,c}^{\text{CMN}} &\approx h_{c}^{\text{MFCC}} + s_{n,c}^{\text{MFCC}} - (h_{c}^{\text{MFCC}} + \bar{s}_{c}^{\text{MFCC}}) = s_{n,c}^{\text{MFCC}} - \bar{s}_{c}^{\text{MFCC}} = s_{n,c}^{\text{CMN}} \\ \Rightarrow \text{ convolution compensated} \end{aligned}$$



# **CMN - Illustration 1st-order Highpass Filter**

#### **Clean vs. Highpass Filtered Logmel Features**

No CMN







# **CMN - Illustration 1st-order Highpass Filter**





Nakatani, Sehr, Kellermann: Reverberant Speech Processing



# **CMN - Illustration Reverberation**







## **CMN - Illustration Reverberation**





Nakatani, Sehr, Kellermann: Reverberant Speech Processing



# **CMN** - Discussion

## Approach

Apply CMN to both training and test data





## Approach

- Apply CMN to both training and test data
- $\Rightarrow$  Short impulse responses can be compensated
- + Good for compensating different microphone characteristics or different telephone channels
- $+\,$  Good for compensating coloration due to early reflections





## Approach

- Apply CMN to both training and test data
- $\Rightarrow$  Short impulse responses can be compensated
- + Good for compensating different microphone characteristics or different telephone channels
- + Good for compensating coloration due to early reflections
- But: not suitable for compensating late reverberation





## Approach

- Apply CMN to both training and test data
- $\Rightarrow$  Short impulse responses can be compensated
- + Good for compensating different microphone characteristics or different telephone channels
- + Good for compensating coloration due to early reflections
- But: not suitable for compensating late reverberation

## **Further considerations**

- Reliable only if utterance is long enough (>4 s [Droppo 2008])
- Extensions necessary for different speech activity rates of training and test data [Droppo 2008]





- Speed of spectral changes of speech:
  - $\Rightarrow$  limited by movements of articulators in vocal tract





- Speed of spectral changes of speech:
   ⇒ limited by movements of articulators in vocal tract
- Many non-speech effects:
  - $\Rightarrow$  characterized by short time-invariant impulse responses Examples: microphone characteristics, telephone channels





- Speed of spectral changes of speech:
   ⇒ limited by movements of articulators in vocal tract
- Many non-speech effects:
  - $\Rightarrow$  characterized by short time-invariant impulse responses Examples: microphone characteristics, telephone channels
- Analysis artifacts:
  - $\Rightarrow$  very fast spectral changes





- Speed of spectral changes of speech:
   ⇒ limited by movements of articulators in vocal tract
- Many non-speech effects:
  - $\Rightarrow$  characterized by short time-invariant impulse responses Examples: microphone characteristics, telephone channels
- Analysis artifacts:
  - $\Rightarrow$  very fast spectral changes

## Idea

- Remove very slow and fast spectral changes from features:
   ⇒ bandpass filtering in each channel
- + Insensitivity to slow and fast spectral changes





#### **Calculation of RASTA Features**







#### **RASTA Features**

- Effective for short time-invariant impulse responses (like CMN)
- + Good for compensating different microphone characteristics or different telephone channels
- + Good for compensating coloration due to early reflections





#### **RASTA Features**

- Effective for short time-invariant impulse responses (like CMN)
- + Good for compensating different microphone characteristics or different telephone channels
- + Good for compensating coloration due to early reflections
- Reverberation described by long RIRs

- Therefore: not suitable for compensating late reverberation





#### Idea

- Temporal changes of short-time spectra:
   important for discriminating phonemes
- First and second derivate of static features (∆ and ∆∆ features): ⇒ capture these changes





# Dynamic Features [Furui 1986]

#### Idea

- Temporal changes of short-time spectra:
   important for discriminating phonemes
- ► First and second derivate of static features (∆ and ∆∆ features): ⇒ capture these changes

## △ Feature Calculation

$$\Delta s_n = s_{n+\kappa} - s_{n-\kappa}$$

typical:  $\kappa \in \{1, 2\}$ 





#### Idea

- Temporal changes of short-time spectra:
   important for discriminating phonemes
- ► First and second derivate of static features (∆ and ∆∆ features): ⇒ capture these changes

## △ Feature Calculation

$$\Delta \boldsymbol{s}_{n} = \boldsymbol{s}_{n+\kappa} - \boldsymbol{s}_{n-\kappa}$$
$$\Delta \boldsymbol{s}_{n} = \frac{\sum_{\kappa=1}^{N_{\Delta}} \kappa \cdot \left(\boldsymbol{s}_{n+\kappa} - \boldsymbol{s}_{n-\kappa}\right)}{2 \cdot \sum_{\kappa=1}^{N_{\Delta}} \kappa^{2}}$$

typical:  $\kappa \in \{1, 2\}$  or  $N_\Delta \in \{2, 3, 4\}$ 





### Idea

- Temporal changes of short-time spectra:
   ⇒ important for discriminating phonemes
- ► First and second derivate of static features (∆ and ∆∆ features): ⇒ capture these changes

## △ Feature Calculation

$$\Delta \boldsymbol{s}_{n} = \boldsymbol{s}_{n+\kappa} - \boldsymbol{s}_{n-\kappa}$$
$$\Delta \boldsymbol{s}_{n} = \frac{\sum_{\kappa=1}^{N_{\Delta}} \kappa \cdot \left(\boldsymbol{s}_{n+\kappa} - \boldsymbol{s}_{n-\kappa}\right)}{2 \cdot \sum_{\kappa=1}^{N_{\Delta}} \kappa^{2}}$$

typical:  $\kappa \in \{1,2\}$  or  $N_\Delta \in \{2,3,4\}$ 

 $\Delta\Delta$  features: calculated in a similar way from  $\Delta$  features





- Long-term relations between feature vectors
- Cannot be captured by HMMs





- Long-term relations between feature vectors
- Cannot be captured by HMMs
  - $\Rightarrow$  Mitigation by feature vectors with long temporal reach





- Long-term relations between feature vectors
- Cannot be captured by HMMs
  - $\Rightarrow$  Mitigation by feature vectors with long temporal reach

### **Temporal reach of features**

- Static features: typically 10 ms 40 ms
- ▲ features: typically 20 ms 120 ms
- $\Delta\Delta$  features: typically 30 ms 200 ms





- Long-term relations between feature vectors
- Cannot be captured by HMMs
  - $\Rightarrow$  Mitigation by feature vectors with long temporal reach

### **Temporal reach of features**

- Static features: typically 10 ms 40 ms
- ▲ features: typically 20 ms 120 ms
- ► △△ features: typically 30 ms 200 ms

#### Dynamic features can partly capture long-term relations




# Model-based Feature Enhancement [Krueger 2010]







Linear dynamic model







Linear dynamic model







Switching linear dynamic model







Switching linear dynamic model



$$\begin{aligned} \boldsymbol{s}_n &= \boldsymbol{A}(q_n)\boldsymbol{s}_{n-1} + \boldsymbol{b}(q_n) + \boldsymbol{u}_n \\ \boldsymbol{p}(\boldsymbol{s}_n | \boldsymbol{s}_{n-1}, q_n) &= \mathcal{N}(\boldsymbol{s}_n; \boldsymbol{A}(q_n) \boldsymbol{s}_{n-1} + \boldsymbol{b}(q_n), \boldsymbol{\Sigma}_{\boldsymbol{u}}(q_n)) \end{aligned}$$





Switching linear dynamic model



$$\begin{aligned} \boldsymbol{s}_n &= \boldsymbol{A}(q_n)\boldsymbol{s}_{n-1} + \boldsymbol{b}(q_n) + \boldsymbol{u}_n \\ \boldsymbol{p}(\boldsymbol{s}_n|\boldsymbol{s}_{n-1},q_n) &= \mathcal{N}(\boldsymbol{s}_n;\boldsymbol{A}(q_n)\boldsymbol{s}_{n-1} + \boldsymbol{b}(q_n),\boldsymbol{\Sigma}_{\boldsymbol{u}}(q_n)) \end{aligned}$$

#### Model for non-stationary feature vector sequences of clean speech









115



based on melspec convolution

$$oldsymbol{x}_n = \log\left(\sum_{ au=0}^{ au_H} \exp(oldsymbol{h}_{ au} + oldsymbol{s}_{n- au})
ight) + oldsymbol{v}_n$$

- **v**<sub>n</sub>: captures approximation error
- $h_{0:T_H}$ : based on strictly exponentially decaying RIR model  $\Rightarrow$  Only  $T_{60}$  needs to be estimated





based on melspec convolution

$$\boldsymbol{x}_n = \log \left( \sum_{\tau=0}^{T_H} \exp(\boldsymbol{h}_{\tau} + \boldsymbol{s}_{n-\tau}) \right) + \boldsymbol{v}_n$$
$$= f(\boldsymbol{s}_{n-T_H:n}, \boldsymbol{h}_{0:T_H}) + \boldsymbol{v}_n$$

 $v_n$ :captures approximation error $h_{0:T_H}$ :based on strictly exponentially decaying RIR model $\Rightarrow$  Only  $T_{60}$  needs to be estimated





based on melspec convolution

$$\begin{aligned} \boldsymbol{x}_n &= \log\left(\sum_{\tau=0}^{T_H} \exp(\boldsymbol{h}_{\tau} + \boldsymbol{s}_{n-\tau})\right) + \boldsymbol{v}_n \\ &= f(\boldsymbol{s}_{n-T_{H:n}}, \boldsymbol{h}_{0:T_H}) + \boldsymbol{v}_n \\ \rho(\boldsymbol{v}_n) &= \mathcal{N}(\boldsymbol{v}_n; \boldsymbol{\mu}_{\boldsymbol{v}}, \boldsymbol{\Sigma}_{\boldsymbol{v}}) \end{aligned}$$





based on melspec convolution

$$\begin{aligned} \boldsymbol{x}_{n} &= \log \left( \sum_{\tau=0}^{T_{H}} \exp(\boldsymbol{h}_{\tau} + \boldsymbol{s}_{n-\tau}) \right) + \boldsymbol{v}_{n} \\ &= f(\boldsymbol{s}_{n-T_{H}:n}, \boldsymbol{h}_{0:T_{H}}) + \boldsymbol{v}_{n} \\ p(\boldsymbol{v}_{n}) &= \mathcal{N}(\boldsymbol{v}_{n}; \boldsymbol{\mu}_{\boldsymbol{v}}, \boldsymbol{\Sigma}_{\boldsymbol{v}}) \\ p(\boldsymbol{x}_{n} | \boldsymbol{s}_{n-T_{H}:n}) &= \mathcal{N}(\boldsymbol{x}_{n}; f(\boldsymbol{s}_{n-T_{H}:n}, \boldsymbol{h}_{0:T_{H}}) + \boldsymbol{\mu}_{\boldsymbol{v}}, \boldsymbol{\Sigma}_{\boldsymbol{v}}) \end{aligned}$$

 $v_n$ :captures approximation error $h_{0:T_H}$ :based on strictly exponentially decaying RIR model $\Rightarrow$  Only  $T_{60}$  needs to be estimated









MMSE estimate

 $\hat{s}_n = E\{s_n | x_{1:n}\}$ 





MMSE estimate

$$\hat{\boldsymbol{s}}_{n} = \operatorname{E} \{ \boldsymbol{s}_{n} | \boldsymbol{x}_{1:n} \}$$

$$p(\boldsymbol{s}_{n} | \boldsymbol{x}_{1:n}) = \frac{p(\boldsymbol{x}_{n} | \boldsymbol{s}_{n}, \boldsymbol{x}_{1:n-1}) \ p(\boldsymbol{s}_{n} | \boldsymbol{x}_{1:n-1})}{\int p(\boldsymbol{x}_{n} | \boldsymbol{s}_{n}, \boldsymbol{x}_{1:n-1}) \ p(\boldsymbol{s}_{n} | \boldsymbol{x}_{1:n-1}) \mathrm{d} \boldsymbol{s}_{n}}$$





MMSE estimate

$$\hat{\mathbf{s}}_{n} = E\{\mathbf{s}_{n}|\mathbf{x}_{1:n}\}$$

$$p(\mathbf{s}_{n}|\mathbf{x}_{1:n}) = \frac{p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ p(\mathbf{s}_{n}|\mathbf{x}_{1:n-1})}{\int p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ p(\mathbf{s}_{n}|\mathbf{x}_{1:n-1}) \, \mathrm{d}\mathbf{s}_{n}}$$

$$\approx \frac{p(\mathbf{x}_{n}|\mathbf{s}_{n-T_{H}:n}) \ \sum_{i=1}^{M} p(\mathbf{s}_{n}|\mathbf{s}_{n-1}, q_{n}=i) \ p(q_{n}=i)}{\int p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ \sum_{i=1}^{M} p(\mathbf{s}_{n}|\mathbf{s}_{n-1}, q_{n}=i) \ p(q_{n}=i) \, \mathrm{d}\mathbf{s}_{n}}$$





MMSE estimate

$$\hat{\mathbf{s}}_{n} = E\{\mathbf{s}_{n}|\mathbf{x}_{1:n}\}$$

$$p(\mathbf{s}_{n}|\mathbf{x}_{1:n}) = \frac{p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ p(\mathbf{s}_{n}|\mathbf{x}_{1:n-1})}{\int p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ p(\mathbf{s}_{n}|\mathbf{x}_{1:n-1}) \,\mathrm{d}\mathbf{s}_{n}}$$

$$\approx \frac{p(\mathbf{x}_{n}|\mathbf{s}_{n-T_{H}:n}) \ \sum_{i=1}^{M} p(\mathbf{s}_{n}|\mathbf{s}_{n-1}, q_{n}=i) \ p(q_{n}=i)}{\int p(\mathbf{x}_{n}|\mathbf{s}_{n}, \mathbf{x}_{1:n-1}) \ \sum_{i=1}^{M} p(\mathbf{s}_{n}|\mathbf{s}_{n-1}, q_{n}=i) \ p(q_{n}=i) \,\mathrm{d}\mathbf{s}_{n}}$$

 $\Rightarrow$  Inference performed by bank of iterated extended Kalman filters





### Discussion

- Approach tailored to reverberant feature vector sequences
- long-term relations explicitely captured by observation model





### Discussion

- Approach tailored to reverberant feature vector sequences
- Iong-term relations explicitely captured by observation model
- + Promising results reported on AURORA 5 task (connected digits)
- + Moderate computational complexity
- + Latency of only a few frames





### Discussion

- Approach tailored to reverberant feature vector sequences
- long-term relations explicitely captured by observation model
- + Promising results reported on AURORA 5 task (connected digits)
- + Moderate computational complexity
- + Latency of only a few frames

Suitable for online recognition in reverberant environments





[Petrick 2008]Harmonicity-based feature analysis[Thomas 2008]Frequency-domain linear prediction[Wölfel 2009]Particle filter-based feature enhancement[Kumar 2010]Cepstral inverse filtering





### Introduction

Feature-based Approaches

Model-based Approaches

Decoder-based Approaches







Mismatch between clean HMM and reverberant data







Feature-based: "dereverberate" data







Model-based: "reverberate" acoustic model







- Model-based: "reverberate" acoustic model
- Adjust acoustic model to statistical properties of reverberant data







- Model-based: "reverberate" acoustic model
- Adjust acoustic model to statistical properties of reverberant data







Record training data in target environment





- Record training data in target environment
  - + Training data perfectly capture statistical properties
  - Extremely high effort





- Record training data in target environment
  - + Training data perfectly capture statistical properties
  - Extremely high effort
- Generate training data by convolution with RIR [Giuliani 1999, Stahl 2001, Matassoni 2002]





- Record training data in target environment
  - + Training data perfectly capture statistical properties
  - Extremely high effort
- Generate training data by convolution with RIR [Giuliani 1999, Stahl 2001, Matassoni 2002]
  - + Significantly reduced effort
  - + Only slight degradation in recognition performance [Stahl 2001]





- Record training data in target environment
  - + Training data perfectly capture statistical properties
  - Extremely high effort
- Generate training data by convolution with RIR [Giuliani 1999, Stahl 2001, Matassoni 2002]
  - + Significantly reduced effort
  - + Only slight degradation in recognition performance [Stahl 2001]

### **Multi-Style Training**

Use training data from many different rooms





- Record training data in target environment
  - + Training data perfectly capture statistical properties
  - Extremely high effort
- Generate training data by convolution with RIR [Giuliani 1999, Stahl 2001, Matassoni 2002]
  - + Significantly reduced effort
  - + Only slight degradation in recognition performance [Stahl 2001]

### **Multi-Style Training**

- Use training data from many different rooms
  - + Robust HMMs
  - + Very flexible
  - Discrimination capability reduced compared to matched training











### Matched Training: Modeling Accuracy

#### Histograms







# Matched Training: Modeling Accuracy

#### Histograms





Auto-CoVariances (ACVs)



# **Multi-Style Training**












































- Capture only linguistic variabilities by acoustic model
- Remove acoustic variabilities by appropriate transforms





- Capture only linguistic variabilities by acoustic model
- Remove acoustic variabilities by appropriate transforms

## Approach

Multi-style training with dereverberated data





- Capture only linguistic variabilities by acoustic model
- Remove acoustic variabilities by appropriate transforms

## Approach

- Multi-style training with dereverberated data
- Similar to noise-adaptive training [Deng 2000] or model-independent adaptive training [Gales 2001]





- Capture only linguistic variabilities by acoustic model
- Remove acoustic variabilities by appropriate transforms

## Approach

- Multi-style training with dereverberated data
- Similar to noise-adaptive training [Deng 2000] or model-independent adaptive training [Gales 2001]
  - + long-term relations partly removed by dereverberation
  - + room dependency reduced
  - $\Rightarrow$  discrimination capability increased compared to multi-style training





- Capture only linguistic variabilities by acoustic model
- Remove acoustic variabilities by appropriate transforms

## Approach

- Multi-style training with dereverberated data
- Similar to noise-adaptive training [Deng 2000] or model-independent adaptive training [Gales 2001]
  - + long-term relations partly removed by dereverberation
  - + room dependency reduced
  - $\Rightarrow$  discrimination capability increased compared to multi-style training
- Successfully applied, e.g., in [Kinoshita 2006]

















































### Approaches

- Maximum A Posteriori adaptation (MAP) [Gauvain 1994]
- Maximum Likelihood Linear Regression (MLLR) [Legetter 1995, Gales 1998]





### Approaches

- Maximum A Posteriori adaptation (MAP) [Gauvain 1994]
- Maximum Likelihood Linear Regression (MLLR) [Legetter 1995, Gales 1998]
- Successfully used for speaker and noise adaptation
- Can also be used for reducing mismatch due to reverberation





### MLLR

Adaptation of the HMM mean vectors

$$\mu_{X} = D\mu_{S} + d$$





### MLLR

Adaptation of the HMM mean vectors and covariance matrices

$$\mu_X = D\mu_S + d$$
  
 $\Sigma_{XX} = E \Sigma_{SS} E^T$ 





### MLLR

Adaptation of the HMM mean vectors and covariance matrices

$$\mu_X = D\mu_S + d$$
  
 $\Sigma_{XX} = E \Sigma_{SS} E^T$ 

Transformation parameters **D**, **d**, **E** estimated by EM algorithm





### MLLR

Adaptation of the HMM mean vectors and covariance matrices

$$\mu_X = D\mu_S + d$$
  
 $\Sigma_{XX} = E \Sigma_{SS} E^T$ 

Transformation parameters **D**, **d**, **E** estimated by EM algorithm

- Supervised MLLR: known transcription
- Unsupervised MLLR: during recognition





## MLLR

Adaptation of the HMM mean vectors and covariance matrices

$$\mu_X = D\mu_S + d$$
  
 $\Sigma_{XX} = E \Sigma_{SS} E^T$ 

Transformation parameters D, d, E estimated by EM algorithm

- Supervised MLLR: known transcription
- Unsupervised MLLR: during recognition

## CMLLR (Constrained MLLR)

Same transformation matrix for mean vector and covariance matrix

$$\mu_{X} = D\mu_{S} + d$$
  
$$\Sigma_{XX} = D\Sigma_{SS} D^{T}$$

+ fewer adaptation parameters





#### Illustration: Example Matched Training on Reverberated Data







- Very accurate description of statistical properties by reverberant training/adaptation data
- Loss of accuracy: only when turning data into model





- Very accurate description of statistical properties by reverberant training/adaptation data
- Loss of accuracy: only when turning data into model
- Reverberant training: requires large amount of reverberant data





- Very accurate description of statistical properties by reverberant training/adaptation data
- Loss of accuracy: only when turning data into model
- Reverberant training: requires large amount of reverberant data
- Data-driven adaptation: moderate amount of reverberant data (but more than model-based adaptation)



130



- Very accurate description of statistical properties by reverberant training/adaptation data
- Loss of accuracy: only when turning data into model
- Reverberant training: requires large amount of reverberant data
- Data-driven adaptation: moderate amount of reverberant data (but more than model-based adaptation)

### **Main Limitation**

Conventional HMMs cannot accurately capture long-term relations





## **Parametric Model-Based Approaches**

# Illustration description of the clean-speech acoustic environment training data (e.g., set of RIRs) adapted HMM clean-speech HMM reverberation representation C















- proposed in [Raut 2006, Hirsch 2008, Sehr 2009]
- based on melspec convolution







- proposed in [Raut 2006, Hirsch 2008, Sehr 2009]
- based on melspec convolution
  - + long-term relations considered for HMM parameter estimation
  - + no adaptation utterances necessary







- proposed in [Raut 2006, Hirsch 2008, Sehr 2009]
- based on melspec convolution
  - + long-term relations considered for HMM parameter estimation
  - + no adaptation utterances necessary
  - reduced accuracy due to approximation errors
  - additional loss of accuracy when mapping combination to HMM







#### Discussion

- proposed in [Raut 2006, Hirsch 2008, Sehr 2009]
- based on melspec convolution
  - + long-term relations considered for HMM parameter estimation
  - + no adaptation utterances necessary
  - reduced accuracy due to approximation errors
  - additional loss of accuracy when mapping combination to HMM

### **Main Limitation**

Conventional HMMs cannot accurately capture long-term relations




#### Mean Adaptation Approach [Raut 2006, Hirsch 2008, Sehr 2009]







## Mean Adaptation Approach [Raut 2006, Hirsch 2008, Sehr 2009]



**Adaptation Equation** 

$$\mu_{X^{\text{MEL}}}(l,j) = \sum_{\boldsymbol{\rho}} \beta(l,j,j-\boldsymbol{\rho}) \ \mu_{S^{\text{MEL}}}(l,j-\boldsymbol{\rho})$$

 $\begin{array}{ll} \beta(I,j,i) & \text{state-level reverberation representation:} \\ & \text{describes energy dispersion from state } i \text{ to } j \text{ in channel } I \\ i,j & \text{state indices} \\ I & \text{mel channel index} \end{array}$ 





# **Estimation of Reverberation Representation [Hirsch 2008]**







# Estimation of Reverberation Representation [Hirsch 2008]



$$h_t^2 = rac{6 \log(10)}{T_{60M}} \cdot \exp\left(-rac{6 \log(10)}{T_{60M}} \cdot t
ight) \ , \quad ext{for} \quad t \geq 0$$







# **Estimation of Reverberation Representation [Hirsch 2008]**







### **Emission pdf of Conventional HMMs**

 $p(\boldsymbol{x}_n|j)$ 

 $\Rightarrow$  conditional independence assumption





### **Emission pdf of Conventional HMMs**

 $p(\boldsymbol{x}_n|j)$ 

 $\Rightarrow$  conditional independence assumption

## **Conditional Emission pdf**

capturing long-term relationships by

 $p(\pmb{x}_n|j,\pmb{x}_{1:n-1})$ 





## **Emission pdf of Conventional HMMs**

 $p(\boldsymbol{x}_n|j)$ 

 $\Rightarrow$  conditional independence assumption

# Conditional Emission pdf

capturing long-term relationships by

 $p(\boldsymbol{x}_n|j,\boldsymbol{x}_{1:n-1})$ 

Approximation by: Context-aware Methods:

- Frame-wise HMM adaptation
- REMOS





# **Conventional Adaptation versus Context-Aware Methods**



(a) Conventional HMM adaptation





# **Conventional Adaptation versus Context-Aware Methods**







# **Conventional Adaptation versus Context-Aware Methods**







$$oldsymbol{x}_n ~pprox ~\log(\exp(oldsymbol{h}_0+oldsymbol{s}_n)+\exp(oldsymbol{r}_n))$$

$$\boldsymbol{\mu}_{\boldsymbol{x}_n}(j) = \log(\exp(\boldsymbol{h}_0 + \boldsymbol{\mu}_{\boldsymbol{s}}(j)) + \exp(\boldsymbol{r}_n))$$

 $\boldsymbol{r}_n$  late reverberation *j* state index





$$oldsymbol{x}_n ~\approx~ \log(\exp(oldsymbol{h}_0+oldsymbol{s}_n)+\exp(oldsymbol{r}_n))$$

$$\boldsymbol{\mu}_{\boldsymbol{x}_n}(j) = \log(\exp(\boldsymbol{h}_0 + \boldsymbol{\mu}_{\boldsymbol{s}}(j)) + \exp(\boldsymbol{r}_n))$$

 $\boldsymbol{r}_n$  late reverberation *j* state index

### Autoregressive Modeling [Takiguchi 2006]

$$r_n = a + x_{n-1}$$

a prediction coefficient





$$oldsymbol{x}_n ~pprox ~\log(\exp(oldsymbol{h}_0+oldsymbol{s}_n)+\exp(oldsymbol{r}_n))$$

$$\boldsymbol{\mu}_{\boldsymbol{x}_n}(j) = \log(\exp(\boldsymbol{h}_0 + \boldsymbol{\mu}_{\boldsymbol{s}}(j)) + \exp(\boldsymbol{r}_n))$$

 $\boldsymbol{r}_n$  late reverberation *j* state index

### Autoregressive Modeling [Takiguchi 2006]

$$\boldsymbol{r}_n = \boldsymbol{a} + \boldsymbol{x}_{n-1}$$

# Moving-Average Modeling [Sehr 2011]

$$m{r}_n = \log \left( \sum_{ au=1}^{T_H} \exp(\mu_{m{h}_{ au}} + m{s}_{n- au}) 
ight)$$





- + Overcomes conditional independence assumption
- + Accurate modeling of long-term relations





- + Overcomes conditional independence assumption
- + Accurate modeling of long-term relations
- Increased computational complexity
- Increased effort for integration into ASR systems





- + Overcomes conditional independence assumption
- + Accurate modeling of long-term relations
- Increased computational complexity
- Increased effort for integration into ASR systems
- Full potential not yet demonstrated





- + Overcomes conditional independence assumption
- + Accurate modeling of long-term relations
- Increased computational complexity
- Increased effort for integration into ASR systems
- Full potential not yet demonstrated

Promising direction for future research





[Couvreur 2001]	Reverberant training of several HMMs + model selection
[Sehr 2010b]	Training of reverberant HMMs on stereo data
[Gales 2011]	Extension of MLLR and VTS to reverberant environments





## Introduction

Feature-based Approaches

Model-based Approaches

Decoder-based Approaches







Modify the decoding algorithm to increase reverberation robustness







Modify the decoding algorithm to increase reverberation robustness

### **Two Approaches**

#### Missing feature techniques

- $\Rightarrow$  Distinguish between reliable and unreliable observations
- $\Rightarrow$  Estimate or discard the unreliable parts



141



Modify the decoding algorithm to increase reverberation robustness

## Two Approaches

## Missing feature techniques

- $\Rightarrow$  Distinguish between reliable and unreliable observations
- $\Rightarrow$  Estimate or discard the unreliable parts

## Uncertainty decoding

- $\Rightarrow$  Combined with signal or feature enhancement techniques
- $\Rightarrow$  Exploit reliability information about enhanced data





Modify the decoding algorithm to increase reverberation robustness

### **Two Approaches**

#### Missing feature techniques

- $\Rightarrow$  Distinguish between reliable and unreliable observations
- $\Rightarrow$  Estimate or discard the unreliable parts

### Uncertainty decoding

- $\Rightarrow$  Combined with signal or feature enhancement techniques
- $\Rightarrow$  Exploit reliability information about enhanced data

Decoder-based approaches bridge the gap between feature-based and model-based approaches



141



# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition





# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition

# Main Steps

- Mask estimation: Mark observations as either reliable or missing
- Handle missing data appropriately





# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition

## Main Steps

- Mask estimation: Mark observations as either reliable or missing
- Handle missing data appropriately

## How to handle missing data?





# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition

## Main Steps

- Mask estimation: Mark observations as either reliable or missing
- Handle missing data appropriately

### How to handle missing data?

 Marginalization: Eliminate unreliable data by integration over corresponding dimensions





# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition

# Main Steps

- Mask estimation: Mark observations as either reliable or missing
- Handle missing data appropriately

## How to handle missing data?

- Marginalization: Eliminate unreliable data by integration over corresponding dimensions
- Bounded marginalization: Exploit known bounds of the missing data for integration





# Key Ideas

- Partition the observations into reliable and missing components
- Use only the reliable components for recognition

# Main Steps

- Mask estimation: Mark observations as either reliable or missing
- Handle missing data appropriately

## How to handle missing data?

- Marginalization: Eliminate unreliable data by integration over corresponding dimensions
- Bounded marginalization: Exploit known bounds of the missing data for integration
- Data imputation: Determine state-dependent estimates for the unreliable data, given the reliable data





## [Palomäki 2004]

- Modulation filtering for the mask estimation
- Bounded marginalization for handling missing data





## [Palomäki 2004]

- Modulation filtering for the mask estimation
- Bounded marginalization for handling missing data

### [Gemmeke 2011]

- Oracle masks based on clean and reverberant features
- Semi-Oracle masks based on clean features and estimated RIRs
- Gaussian-dependent bounded imputation







### **Conventional Feature Enhancement Methods**







# **Uncertainty Decoding**

### **Conventional Feature Enhancement Methods**



▶ Use only point estimate  $\hat{s}_n$  of clean features





# **Uncertainty Decoding**

## **Conventional Feature Enhancement Methods**



- Use only point estimate  $\hat{s}_n$  of clean features
- Contribution of each Gaussian component *m*:

$$p(\hat{m{s}}_n|m) = \mathcal{N}(\hat{m{s}}_n; m{\mu}_{m{s}}^{(m)}, \Sigma_{m{s}}^{(m)})$$





[Droppo 2002, Deng 2005, Liao 2008, Haeb-Umbach 2011]

#### Feature Enhancement Combined with Uncertainty Decoding






#### Feature Enhancement Combined with Uncertainty Decoding



Signal/feature enhancement inevitably introduces distortions





#### Feature Enhancement Combined with Uncertainty Decoding



- Signal/feature enhancement inevitably introduces distortions
- ► Use reliability information in addition to point estimate  $\Rightarrow$  Use  $p(\hat{s}_n | s_n)$  instead of  $\hat{s}_n$





**Mismatch Model** 

$$\begin{aligned} \hat{\boldsymbol{s}}_n &= \boldsymbol{s}_n + \boldsymbol{b}_n \\ p(\boldsymbol{b}_n) &= \mathcal{N}(\boldsymbol{b}_n; \boldsymbol{0}, \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \\ p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n, m) &\approx p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n) = p(\boldsymbol{b}_n) \end{aligned}$$





**Mismatch Model** 

$$\hat{oldsymbol{s}}_n = oldsymbol{s}_n + oldsymbol{b}_n$$
 $p(oldsymbol{b}_n) = \mathcal{N}(oldsymbol{b}_n; oldsymbol{0}, \Sigma_{oldsymbol{b}_n})$ 
 $p(\hat{oldsymbol{s}}_n | oldsymbol{s}_n, m) \approx p(\hat{oldsymbol{s}}_n | oldsymbol{s}_n) = p(oldsymbol{b}_n)$ 

### **Contribution of Gaussian Component** *m*

$$p(\hat{\mathbf{s}}_n|m) = \int p(\hat{\mathbf{s}}_n, \mathbf{s}_n|m) \, \mathrm{d}\mathbf{s}_n = \int p(\hat{\mathbf{s}}_n|\mathbf{s}_n, m) \, p(\mathbf{s}_n|m) \, \mathrm{d}\mathbf{s}_n$$





**Mismatch Model** 

$$\hat{oldsymbol{s}}_n = oldsymbol{s}_n + oldsymbol{b}_n$$
 $p(oldsymbol{b}_n) = \mathcal{N}(oldsymbol{b}_n; oldsymbol{0}, \Sigma_{oldsymbol{b}_n})$ 
 $p(\hat{oldsymbol{s}}_n | oldsymbol{s}_n, m) \approx p(\hat{oldsymbol{s}}_n | oldsymbol{s}_n) = p(oldsymbol{b}_n)$ 

### **Contribution of Gaussian Component** *m*

$$\begin{split} p(\hat{\boldsymbol{s}}_n|m) &= \int p(\hat{\boldsymbol{s}}_n, \boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n = \int p(\hat{\boldsymbol{s}}_n|\boldsymbol{s}_n, m) \, p(\boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n \\ &\approx \quad \mathcal{N}(\hat{\boldsymbol{s}}_n; \boldsymbol{\mu}_{\boldsymbol{s}}^{(m)}, \boldsymbol{\Sigma}_{\boldsymbol{s}}^{(m)} + \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \end{split}$$





**Mismatch Model** 

$$\begin{aligned} \hat{\boldsymbol{s}}_n &= \boldsymbol{s}_n + \boldsymbol{b}_n \\ p(\boldsymbol{b}_n) &= \mathcal{N}(\boldsymbol{b}_n; \boldsymbol{0}, \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \\ p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n, m) &\approx p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n) = p(\boldsymbol{b}_n) \end{aligned}$$

### **Contribution of Gaussian Component** m

1

$$\begin{split} p(\hat{\boldsymbol{s}}_n|m) &= \int p(\hat{\boldsymbol{s}}_n, \boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n = \int p(\hat{\boldsymbol{s}}_n|\boldsymbol{s}_n, m) \, p(\boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n \\ &\approx \quad \mathcal{N}(\hat{\boldsymbol{s}}_n; \boldsymbol{\mu}_{\boldsymbol{s}}^{(m)}, \boldsymbol{\Sigma}_{\boldsymbol{s}}^{(m)} + \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \end{split}$$

▶ Unreliable features  $\Rightarrow$  large  $\Sigma_{b_n}$   $\Rightarrow$  little effect on Viterbi score





**Mismatch Model** 

$$\begin{aligned} \hat{\boldsymbol{s}}_n &= \boldsymbol{s}_n + \boldsymbol{b}_n \\ p(\boldsymbol{b}_n) &= \mathcal{N}(\boldsymbol{b}_n; \boldsymbol{0}, \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \\ p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n, m) &\approx p(\hat{\boldsymbol{s}}_n | \boldsymbol{s}_n) = p(\boldsymbol{b}_n) \end{aligned}$$

### **Contribution of Gaussian Component** *m*

$$\begin{split} p(\hat{\boldsymbol{s}}_n|m) &= \int p(\hat{\boldsymbol{s}}_n, \boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n = \int p(\hat{\boldsymbol{s}}_n|\boldsymbol{s}_n, m) \, p(\boldsymbol{s}_n|m) \, \mathrm{d}\boldsymbol{s}_n \\ &\approx \quad \mathcal{N}(\hat{\boldsymbol{s}}_n; \boldsymbol{\mu}_{\boldsymbol{s}}^{(m)}, \boldsymbol{\Sigma}_{\boldsymbol{s}}^{(m)} + \boldsymbol{\Sigma}_{\boldsymbol{b}_n}) \end{split}$$

- ▶ Unreliable features  $\Rightarrow$  large  $\Sigma_{b_n}$   $\Rightarrow$  little effect on Viterbi score
- Main challenge: Estimation of time-variant feature cov. Σ<sub>b<sub>n</sub></sub>



[Delcroix 2009, Delcroix 2011a]

### Key Idea

- Strong reverberation
  - $\Rightarrow$  Large effect of speech enhancement
  - $\Rightarrow$  Large mismatch between clean and enhanced features





[Delcroix 2009, Delcroix 2011a]

### Key Idea

- Strong reverberation
  - $\Rightarrow$  Large effect of speech enhancement
  - $\Rightarrow$  Large mismatch between clean and enhanced features
- Effect of speech enhancement captured by  $\hat{\boldsymbol{b}}_n = \boldsymbol{x}_n \hat{\boldsymbol{s}}_n$ 
  - $\Rightarrow$  Mismatch covariance assumed proportional to difference between observed and enhanced features





[Delcroix 2009, Delcroix 2011a]

Key Idea

- Strong reverberation
  - $\Rightarrow$  Large effect of speech enhancement
  - $\Rightarrow$  Large mismatch between clean and enhanced features
- Figure 6 Effect of speech enhancement captured by  $\hat{\boldsymbol{b}}_n = \boldsymbol{x}_n \hat{\boldsymbol{s}}_n$ 
  - ⇒ Mismatch covariance assumed proportional to difference between observed and enhanced features

Model elements of time-variant diagonal mismatch cov. matrix  $\Sigma_{b_n}$  as

$$(\boldsymbol{\Sigma}_{\boldsymbol{b}_n})_{ii} = \alpha_i \, \hat{b}_{n,i}^2$$





[Delcroix 2009, Delcroix 2011a]

Key Idea

- Strong reverberation
  - $\Rightarrow$  Large effect of speech enhancement
  - $\Rightarrow$  Large mismatch between clean and enhanced features
- Figure 6 Effect of speech enhancement captured by  $\hat{\boldsymbol{b}}_n = \boldsymbol{x}_n \hat{\boldsymbol{s}}_n$ 
  - ⇒ Mismatch covariance assumed proportional to difference between observed and enhanced features

Model elements of time-variant diagonal mismatch cov. matrix  $\Sigma_{b_n}$  as

$$(\boldsymbol{\Sigma}_{\boldsymbol{b}_n})_{ii} = \alpha_i \hat{\boldsymbol{b}}_{n,i}^2$$

•  $\alpha$  is estimated by EM algorithm using adaptation data





#### [Delcroix 2009, Delcroix 2011a]









[Delcroix 2009, Delcroix 2011a]

#### Discussion

– Accounting for time-variant covariance matrix  $\Sigma_{\boldsymbol{b}_n}$  increases computational complexity





[Delcroix 2009, Delcroix 2011a]

#### Discussion

- Accounting for time-variant covariance matrix  $\Sigma_{b_n}$  increases computational complexity
- + Can be combined with static variance compensation and mean adaptation by MLLR
- $+ \,$  Independent of enhancement algorithm  $\Rightarrow$  Highly flexible
- + Has been used successfully also for non-stationary interferences [Delcroix 2011b]





[Delcroix 2009, Delcroix 2011a]

#### Discussion

- Accounting for time-variant covariance matrix  $\Sigma_{\boldsymbol{b}_n}$  increases computational complexity
- + Can be combined with static variance compensation and mean adaptation by MLLR
- + Independent of enhancement algorithm  $\Rightarrow$  Highly flexible
- + Has been used successfully also for non-stationary interferences [Delcroix 2011b]

Promising approach for interconnection of signal/feature-based methods and ASR systems





### Introduction

Feature-based Approaches

Model-based Approaches

Decoder-based Approaches









### **Online Model Combination**



- $\begin{array}{rl} \text{CSM:} & \text{clean-speech model} \\ & \Rightarrow \text{HMM network} \end{array}$
- RVM: reverberation model

combination of CSM and RVM:

 $\Rightarrow$  context-aware acoustic model





### **Online Model Combination**



- $\begin{array}{rl} \text{CSM:} & \text{clean-speech model} \\ & \Rightarrow \text{HMM network} \end{array}$
- RVM: reverberation model

combination of CSM and RVM:

 $\Rightarrow$  context-aware acoustic model

#### Advantages

CSM and RVM are trained independently

- changing environment: adjust only RVM
- changing task: adjust only CSM



high degree of flexibility

[Sehr 2010c]





## **REMOS Decoding [Sehr 2010c]**



#### Extended Viterbi Algorithm:

finds most likely path through CSM





## **REMOS Decoding [Sehr 2010c]**



#### Extended Viterbi Algorithm:

finds most likely path through CSM

#### Inner Optimization:

- accounts for RVM and previous observations
- determines most likely contributions of CSM and RVM to current observation





Online combination of model outputs from clean-speech HMM and reverberation model capturing long-term relations:





Online combination of model outputs from clean-speech HMM and reverberation model capturing long-term relations:

### **Combination Operator**

$$\mathbf{x}_n = \mathbf{f}(\mathbf{s}_n, \mathbf{s}_{n-T_H:n-1}, \mathbf{h}_n, \mathbf{a}_n)$$

$$= \log(\exp(\mathbf{h}_n + \mathbf{s}_n) + \exp(\mathbf{r}_n + \mathbf{a}_n))$$

- *r<sub>n</sub>*: logmelspec late reverberation estimate*a<sub>n</sub>*: captures approxima
  - tion error of  $r_n$
- h<sub>n</sub>: logmelspec representation of direct sound component of RIR





Online combination of model outputs from clean-speech HMM and reverberation model capturing long-term relations:

### **Combination Operator**

$$\boldsymbol{x}_n = \boldsymbol{f}(\boldsymbol{s}_n, \boldsymbol{s}_{n-T_H:n-1}, \boldsymbol{h}_n, \boldsymbol{a}_n)$$

$$= \log(\exp(\mathbf{h}_n + \mathbf{s}_n) + \exp(\mathbf{r}_n + \mathbf{a}_n))$$

### Late Reverberation Estimate

$$\boldsymbol{r}_n = \log \left( \sum_{\tau=1}^{T_H} \exp(\boldsymbol{\mu}_{\tau}^H + \boldsymbol{s}_{n-\tau}) \right)$$

*r*<sub>n</sub>: logmelspec late reverberation estimate

a<sub>n</sub>: captures approximation error of r<sub>n</sub>

*h*<sub>n</sub>: logmelspec representation of direct sound component of RIR

 $\mu^{H}_{1:\mathcal{T}_{H}} \colon \underset{\text{mean vectors of log-melspec representation for late reverber-ation}}{\text{mean vectors of log-melses}}$ 



#### **Illustration of Generative Model**







Conditional emission pdf is decomposed into reverberation model and clean HMM:

$$p(\boldsymbol{x}_n|j, \boldsymbol{x}_{1:n-1}) = \int p(\boldsymbol{x}_n|\boldsymbol{s}_n, \boldsymbol{x}_{1:n-1}) p(\boldsymbol{s}_n|j) \, \mathrm{d}\boldsymbol{s}_n,$$





Conditional emission pdf is decomposed into reverberation model and clean HMM:

$$p(\boldsymbol{x}_n|j, \boldsymbol{x}_{1:n-1}) = \int p(\boldsymbol{x}_n|\boldsymbol{s}_n, \boldsymbol{x}_{1:n-1}) p(\boldsymbol{s}_n|j) \, \mathrm{d}\boldsymbol{s}_n,$$

#### **Reverberation Model:**

$$p(\boldsymbol{x}_n | \boldsymbol{s}_n, \boldsymbol{x}_{1:n-1}) = \iint p(\boldsymbol{h}_n) p(\boldsymbol{a}_n) \, \delta(\boldsymbol{x}_n - \boldsymbol{f}(\boldsymbol{s}_n, \boldsymbol{s}_{n-T_H:n-1}, \boldsymbol{h}_n, \boldsymbol{a}_n)) \, \mathrm{d}\boldsymbol{h}_n \, \mathrm{d}\boldsymbol{a}_n$$







### Approximation of Conditional Emission pdf:

by maximum values of integrand

 $p(\boldsymbol{x}_n|j, \boldsymbol{x}_{1:n-1}) \approx p(\hat{\boldsymbol{h}}_n) p(\hat{\boldsymbol{a}}_n) p(\hat{\boldsymbol{s}}_n|j)$ 





#### **Approximation of Conditional Emission pdf:** by maximum values of integrand

$$p(\boldsymbol{x}_n|j, \boldsymbol{x}_{1:n-1}) \approx p(\hat{\boldsymbol{h}}_n) p(\hat{\boldsymbol{a}}_n) p(\hat{\boldsymbol{s}}_n|j)$$

maximum values  $\hat{h}_n$ ,  $\hat{a}_n$ ,  $\hat{s}_n$  determined by inner optimization

$$(\hat{\boldsymbol{h}}_n, \hat{\boldsymbol{a}}_n, \hat{\boldsymbol{s}}_n) = \operatorname*{argmax}_{(\boldsymbol{h}_n, \boldsymbol{a}_n, \boldsymbol{s}_n)} p(\boldsymbol{h}_n) p(\boldsymbol{a}_n) p(\boldsymbol{s}_n | j)$$

subject to 
$$\boldsymbol{x}_n = \boldsymbol{f}(\boldsymbol{s}_n, \boldsymbol{s}_{n-T_H:n-1}, \boldsymbol{h}_n, \boldsymbol{a}_n)$$





























# **Modeling Accuracy of REMOS**

#### Example: digit "seven"







## **Modeling Accuracy of REMOS**

#### **Histograms**







# **Modeling Accuracy of REMOS**

#### Histograms





Auto-CoVariances (ACVs)



## **Recognition Results [Sehr 2010c]**



#### Setup

- Task: Connected digits (TI digits)
- Features: Logmelspec coefficients
- Recognizer: Word-level HMMs, 16 states/digit, 1 Gaussian/state

### Rooms:

	$T_{60}$	DRR
A:	300 ms	4.0dB
B:	700 ms	-4.0dB
C:	900 ms	-4.0dB


#### Discussion

- + Approach tailored to reverberant feature vector sequences
- + Long-term relations explicitely captured by reverberation model
- + Reverberation exploited for discrimination
- + Very promising results in logmelspec domain





#### Discussion

- + Approach tailored to reverberant feature vector sequences
- + Long-term relations explicitely captured by reverberation model
- + Reverberation exploited for discrimination
- + Very promising results in logmelspec domain
- Inner optimization increases decoding complexity
- Implementation requires changes in decoding routines







#### Discussion

- + Approach tailored to reverberant feature vector sequences
- + Long-term relations explicitely captured by reverberation model
- + Reverberation exploited for discrimination
- + Very promising results in logmelspec domain
- Inner optimization increases decoding complexity
- Implementation requires changes in decoding routines

#### Promising direction for future research







# **Dereverberation for Signal Enhancement**

- Close to 12 dB DRR gain with  $T_{60} \approx 0.7s$  (offline) with
  - 4 mics, d=1.65m, no noise (TRINICON, 2 sources)
  - ► 8 mics, d=2m, SNR=10 dB (MCLP)







# **Dereverberation for Signal Enhancement**

### State-of-the-art

- Close to 12 dB DRR gain with  $T_{60} \approx 0.7s$  (offline) with
  - 4 mics, d=1.65m, no noise (TRINICON, 2 sources)
  - ► 8 mics, d=2m, SNR=10 dB (MCLP)

## Challenges

- Larger distances, more reverberant rooms
- Robustness to speech-like interferers, nonstationary/diffuse noise, transient echo cancellation residuals
- Robust tracking of time-varying acoustics
- Low-latency ( $\ll$  1*s*) and efficient real-time implementations
- Joint optimization with spectral subtraction techniques





# **IV.** Summary, Conclusions, and Outlook (cont'd)

**Dereverberation as preprocessing for ASR** 







# **Dereverberation as preprocessing for ASR**

**Example:** 20 k WSJ convolved with RIRs ( $T_{60} = 0.78s$ , d = 2m), NTT ASR

WER[%]	Preproc.	Acoustic model
85.5	none	clean speech
43.4	none	multi-condition training
26.1	1-ch derev	multi-condition training
14.2	2-ch derev	clean w/ unsuperv.
		speaker adaptation by MLLR





# **Dereverberation as preprocessing for ASR**

**Example:** 20 k WSJ convolved with RIRs ( $T_{60} = 0.78s$ , d = 2m), NTT ASR

WER[%]	Preproc.	Acoustic model
85.5	none	clean speech
43.4	none	multi-condition training
26.1	1-ch derev	multi-condition training
14.2	2-ch derev	clean w/ unsuperv.
		speaker adaptation by MLLR

**Challenges** for approaching close-talk performance

- Transition from reverberated signals to real recordings
- Self-adaptation to changing acoustics and front-ends, including
  - variable number and changing, unconstrained positions of talkers
  - different nodes in distributed microphone arrays
- Joint optimization with ASR methods to handle reverberation and noise











- Feature-based techniques
  - account for the inter-frame relations caused by dispersion
  - efficiently exploit predictability of reverberation





- Feature-based techniques
  - account for the inter-frame relations caused by dispersion
  - efficiently exploit predictability of reverberation
- Model-based techniques
  - could not yet show their full potential, as
  - framewise adaptation and optimization is computationally complex





- Feature-based techniques
  - account for the inter-frame relations caused by dispersion
  - efficiently exploit predictability of reverberation
- Model-based techniques
  - could not yet show their full potential, as
  - framewise adaptation and optimization is computationally complex
- Decoder-based techniques
  - compromise between the above regarding complexity





#### State-of-the-art

- Feature-based techniques
  - account for the inter-frame relations caused by dispersion
  - efficiently exploit predictability of reverberation
- Model-based techniques
  - could not yet show their full potential, as
  - framewise adaptation and optimization is computationally complex
- Decoder-based techniques
  - compromise between the above regarding complexity

# Outlook

- Integration into state-of-the art ASR systems
  - expected soon for signal enhancement- and feature-based methods
  - model-based methods must become more efficient for widespread use









165



Blind deconvolution of the acoustic paths seems to come closer







- Blind deconvolution of the acoustic paths seems to come closer
- Less ambitious algorithms are also effective and their progress follows the typical DSP objectives
  - increase algorithmic performance and robustness
  - reduce computational load
  - integrate with other functionalities





- Blind deconvolution of the acoustic paths seems to come closer
- Less ambitious algorithms are also effective and their progress follows the typical DSP objectives
  - increase algorithmic performance and robustness
  - reduce computational load
  - integrate with other functionalities

As a follow-up to the CHIME Challenge 2011

⇒ Next Challenge for Reverberation-robust Speech Processing is underway!





### We are especially grateful to

- Dr. Keisuke Kinoshita, Dr. Marc Delcroix, Dr. Shoko Araki, Dr. Mehrez Souden, and Dr. Takaaki Hori (NTT)
- Dr. Herbert Buchner, Edwin Mabande and Lutz Marquardt (formerly LMS)
- Roland Maas and Christian Hofmann (LMS)

for their contributions to the course material





### We are especially grateful to

- Dr. Keisuke Kinoshita, Dr. Marc Delcroix, Dr. Shoko Araki, Dr. Mehrez Souden, and Dr. Takaaki Hori (NTT)
- Dr. Herbert Buchner, Edwin Mabande and Lutz Marquardt (formerly LMS)
- Roland Maas and Christian Hofmann (LMS)

for their contributions to the course material and wish to acknowledge the support of parts of the LMS by

Deutsche Forschungsgemeinschaft (DFG) under contract number KE 890/4-1





# ご清聴ありがとうございました





