

# Noise Compensation for Subspace Gaussian Mixture Models

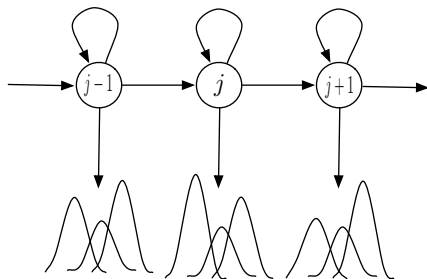
Liang Lu  
University of Edinburgh

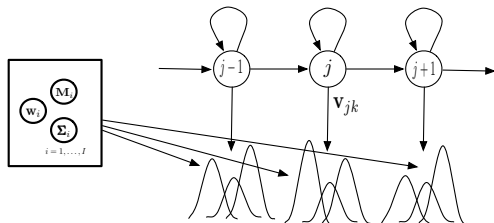
Joint work with KK Chin, A. Ghoshal and S. Renals

- ▶ Motivation
  - ▶ Subspace GMM (SGMM) works well in matched speech condition [Povey et al., 2011]
  - ▶ In mismatched condition (i.e. noise), the gain disappears
- ▶ Goal
  - ▶ Noise compensation for SGMM
- ▶ Method
  - ▶ Model space compensation
  - ▶ Joint uncertainty decoding (JUD) [Liao and Gales, 2005]



# HMM-GMM acoustic model





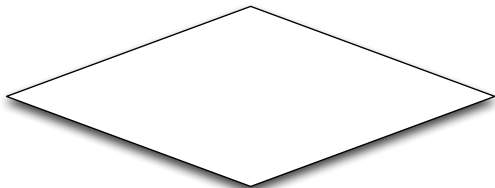
## ▶ Global

- ▶  $\mathbf{M}_i$  is the basis for means
- ▶  $\mathbf{w}_i$  is the basis for weights
- ▶  $\Sigma_i$  is the covariance matrix

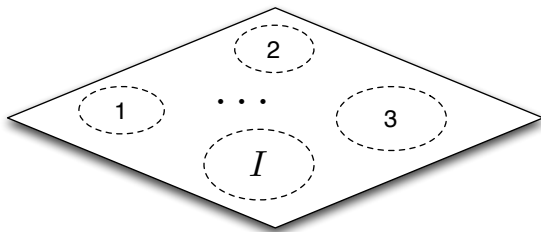
## ▶ State-dependent

- ▶  $\mathbf{v}_{jk}$  is low dimensional vectors (e.g. 40dim)
- ▶ Gaussian mean:  $\boldsymbol{\mu}_{jki} = \mathbf{M}_i \mathbf{v}_{jk}$

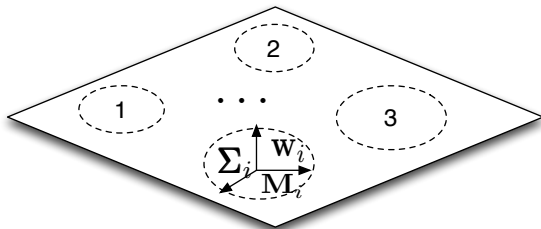
- ▶ More intuitively, suppose we have an acoustic space like this



- ▶ We then partition the whole acoustic space into  $I$  regions.
- ▶ This can be done by learning a GMM using the training data.



- ▶ We then introduce some parameters to structure each region



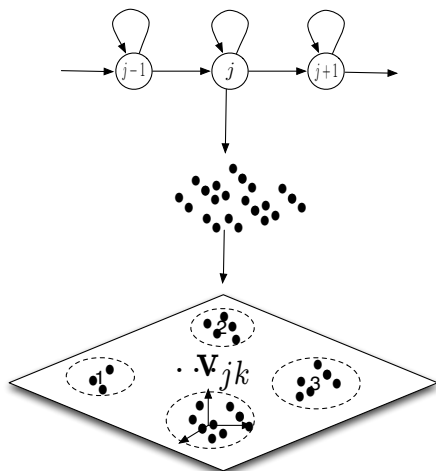
$\Sigma_i$  - model the covariance of this region

$\mathbf{M}_i$  - span the basis for Gaussian mean

$\mathbf{W}_i$  - span the basis for Gaussian weight

# Subspace Gaussian Mixture Models

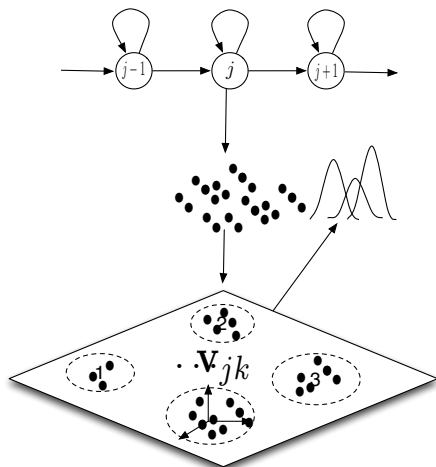
Given a class with some data, such as an HMM state



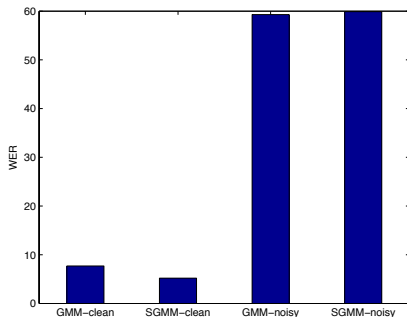


# Subspace Gaussian Mixture Models

Then we learn a GMM for this class



- ▶ Larger modelling power  $\rightarrow$  higher recognition accuracy.
  - ▶ Our systems on Aurora 4, the #Gaussians is 6.4M (SGMM), vs. 50k (GMM).
  - ▶ SGMM vs. GMM  $\rightarrow$  5.2% vs. 7.7% on clean condition
  - ▶ SGMM vs. GMM  $\rightarrow$  59.9% vs. 59.3% on noisy condition
- ▶ Can we do noise compensation for SGMMs ?

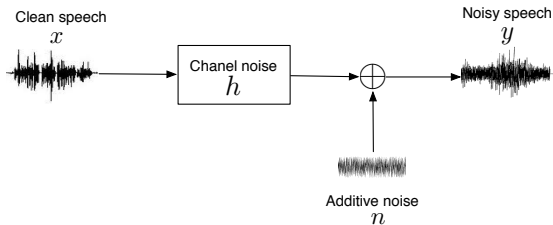


There are numerous work on noise compensation for robust ASR  
[Deng, 2011]

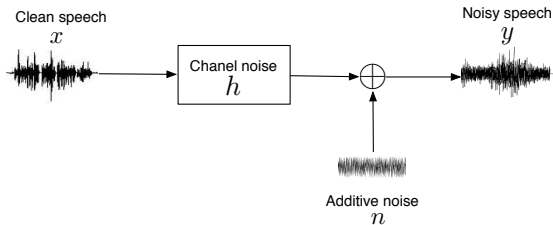
- ▶ Feature domain
  - ▶ Spectral subtraction, cmn/cvn
  - ▶ Cepstral mean square error estimation
  - ▶ Algonquin
  - ▶ Splice
  - ▶ Feature space vector Taylor series (VTS)
- ▶ Model domain
  - ▶ MLLR, noise constraint MLLR
  - ▶ PMC, Data-driven PMC (DPMC), iterative DPMC
  - ▶ VTS, joint uncertainty decoding (JUD)
  - ▶ Linear spline interpolation (LSI)
  - ▶ Unscented transform (UT)
- ▶ Hybrid
  - ▶ Noise adaptive training



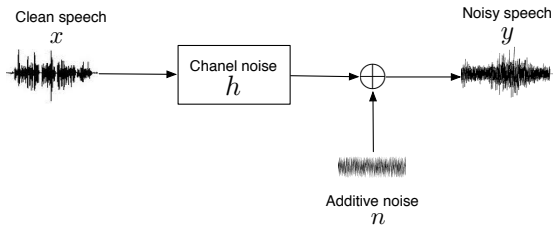
- ▶ Model space compensation for SGMM
- ▶ Not data-driven but using heuristic knowledge
- ▶ Mismatch function  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \alpha)$  [Acero, 1990]
- ▶  $\alpha$  denotes the phase term between noise and speech [Deng et al., 2004].



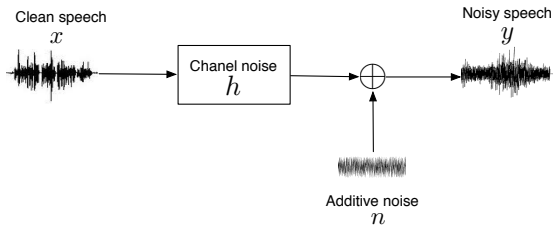
- ▶ Model space compensation for SGMM
- ▶ Not data-driven but using heuristic knowledge
- ▶ Mismatch function  $y = f(x, h, n, \alpha)$  [Acero, 1990]
- ▶  $\alpha$  denotes the phase term between noise and speech [Deng et al., 2004].



- ▶ Model space compensation for SGMM
- ▶ Not data-driven but using heuristic knowledge
- ▶ Mismatch function  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \alpha)$  [Acero, 1990]
- ▶  $\alpha$  denotes the phase term between noise and speech [Deng et al., 2004].



- ▶ Model space compensation for SGMM
- ▶ Not data-driven but using heuristic knowledge
- ▶ Mismatch function  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  [Acero, 1990]
- ▶  $\boldsymbol{\alpha}$  denotes the phase term between noise and speech [Deng et al., 2004].



The mismatch function is

$$\begin{aligned} \mathbf{y} &= f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \alpha) \\ &= \mathbf{x} + \mathbf{h} + \mathbf{C} \log \left[ \mathbf{1} + \exp(\mathbf{C}^{-1}(\mathbf{n} - \mathbf{x} - \mathbf{h})) \right. \\ &\quad \left. + \underbrace{2\alpha \bullet \exp(\mathbf{C}^{-1}(\mathbf{n} - \mathbf{x} - \mathbf{h})/2)}_{\text{phase term}} \right]. \end{aligned} \tag{1}$$

where  $\mathbf{C}$  be the DCT matrix.



- ▶ **Aim:** estimate  $\mu_y$  and  $\Sigma_y$  for each Gaussian component.
- ▶ **Difficulty:**  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  is highly nonlinear, no analytic solution!
- ▶ **Solution:** Vector Taylor series (VTS) approximation [Moreno et al., 1996]
- ▶ **Cost:** Real time factor  $> 100$ , memory  $> 10G$  for (medium size) SGMM with 6.4M Gaussian
- ▶ **Inelegant:** Direct apply VTS will destroy the compact of structure of SGMMs



- ▶ **Aim:** estimate  $\mu_y$  and  $\Sigma_y$  for each Gaussian component.
- ▶ **Difficulty:**  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  is highly nonlinear, no analytic solution!
- ▶ **Solution:** Vector Taylor series (VTS) approximation [Moreno et al., 1996]
- ▶ **Cost:** Real time factor  $> 100$ , memory  $> 10G$  for (medium size) SGMM with 6.4M Gaussian
- ▶ **Inelegant:** Direct apply VTS will destroy the compact of structure of SGMMs



- ▶ **Aim:** estimate  $\mu_y$  and  $\Sigma_y$  for each Gaussian component.
- ▶ **Difficulty:**  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  is highly nonlinear, no analytic solution!
- ▶ **Solution:** Vector Taylor series (VTS) approximation [Moreno et al., 1996]
- ▶ **Cost:** Real time factor  $> 100$ , memory  $> 10G$  for (medium size) SGMM with 6.4M Gaussian
- ▶ **Inelegant:** Direct apply VTS will destroy the compact of structure of SGMMs



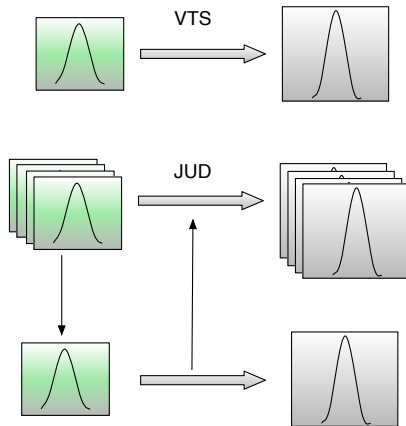
- ▶ **Aim:** estimate  $\mu_y$  and  $\Sigma_y$  for each Gaussian component.
- ▶ **Difficulty:**  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  is highly nonlinear, no analytic solution!
- ▶ **Solution:** Vector Taylor series (VTS) approximation [Moreno et al., 1996]
- ▶ **Cost:** Real time factor  $> 100$ , memory  $> 10G$  for (medium size) SGMM with 6.4M Gaussian
- ▶ **Inelegant:** Direct apply VTS will destroy the compact of structure of SGMMs



- ▶ **Aim:** estimate  $\mu_y$  and  $\Sigma_y$  for each Gaussian component.
- ▶ **Difficulty:**  $\mathbf{y} = f(\mathbf{x}, \mathbf{h}, \mathbf{n}, \boldsymbol{\alpha})$  is highly nonlinear, no analytic solution!
- ▶ **Solution:** Vector Taylor series (VTS) approximation [Moreno et al., 1996]
- ▶ **Cost:** Real time factor  $> 100$ , memory  $> 10G$  for (medium size) SGMM with 6.4M Gaussian
- ▶ **Inelegant:** Direct apply VTS will destroy the compact of structure of SGMMs

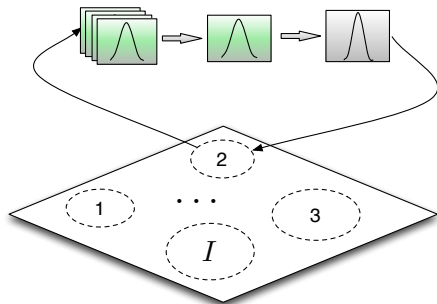


- **Solution:** Joint uncertainty decoding (JUD)



VTS vs. JUD

- ▶ Applying JUD to SGMM



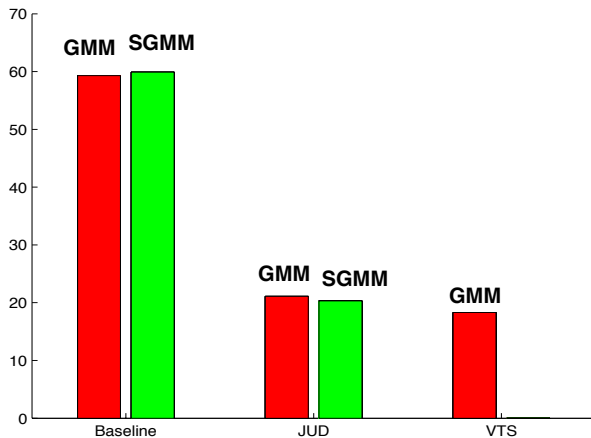
- ▶ **Cost:** Real time factor  $\sim 10$  for SGMM with 6.4M Gaussians

- ▶ Database
  - ▶ Aurora 4 dataset
  - ▶ Clean speech and noisy speech with SNR [5db - 15db]
  - ▶ Close-talking microphone and desk-mounted microphone
  - ▶ ~ 15 hour training data
  - ▶ 330 testing utterances
- ▶ System configuration
  - ▶ 39dim MFCC
  - ▶ #triphone states: 3.1k (GMM) vs. 3.9k (SGMM)
  - ▶ #Gaussians: 50k (GMM) vs. 6.4M (SGMM)
  - ▶ #regression classes: 112 (GMM) vs. 400 (SGMM)

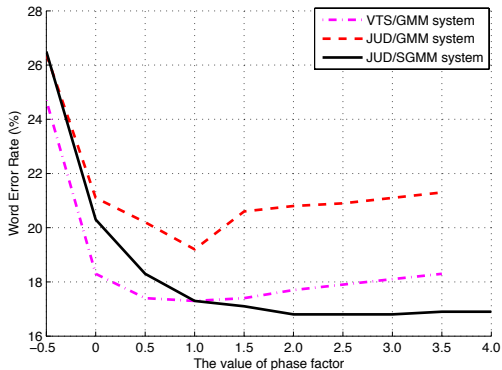




# Noise compensation experiments



Results by tuning the value of phase factors.



- ▶ JUD/SGMM system achieved **16.8%** WER on Aurora 4 database

- ▶ The phase term is very effective for noise compensation
- ▶ Similar improvements were also observed in other studies, e.g. [Li et al., 2009]
- ▶ The reasons maybe it can compensate for the linearization bias and performs domain compensation [Li et al., 2009]
- ▶ Our insight is it may helps to avoid the over estimation of the noise model



- ▶ The phase term is very effective for noise compensation
- ▶ Similar improvements were also observed in other studies, e.g. [Li et al., 2009]
- ▶ The reasons maybe it can compensate for the linearization bias and performs domain compensation [Li et al., 2009]
- ▶ Our insight is it may helps to avoid the over estimation of the noise model



- ▶ The phase term is very effective for noise compensation
- ▶ Similar improvements were also observed in other studies, e.g. [Li et al., 2009]
- ▶ The reasons maybe it can compensate for the linearization bias and performs domain compensation [Li et al., 2009]
- ▶ Our insight is it may helps to avoid the over estimation of the noise model



- ▶ The phase term is very effective for noise compensation
- ▶ Similar improvements were also observed in other studies, e.g. [Li et al., 2009]
- ▶ The reasons maybe it can compensate for the linearization bias and performs domain compensation [Li et al., 2009]
- ▶ Our insight is it may helps to avoid the over estimation of the noise model



- ▶ SGMM is a promising alternative for acoustic modelling
- ▶ Noise compensation using JUD works well for SGMMs
- ▶ The phase term is particular effective for the noise compensation
- ▶ Future works will be on noise adaptive training, compensation in log-spectral domain.





- ▶ With JUD, the marginal likelihood can be obtained as

$$p(\mathbf{y} | m) \approx |\mathbf{A}^{(r)}| \mathcal{N} \left( \mathbf{A}^{(r)} \mathbf{y} + \mathbf{b}^{(r)}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m + \boldsymbol{\Sigma}_b^{(r)} \right). \quad (2)$$

- ▶ The transformation is done in the feature space, applied to each frame
- ▶ Computation is saved since that the  $\#frame \ll \#Gaussians$
- ▶ The transformation should be diagonalized in GMM systems, but not in SGMM system since we used full covariance matrix



Table: GMM systems with  $\alpha = 0$ .

Methods	Clean	Avg
Clean model	7.7	<b>59.3</b>
MTR model	12.7	<b>26.9</b>
VTS	7.3	<b>18.3</b>
JUD	7.0	<b>21.1</b>

Table: SGMM systems with  $\alpha = 0$ .

Methods	Clean	Avg
Clean model	5.2	<b>59.9</b>
MTR model	6.8	<b>22.2</b>
JUD	5.3	<b>20.3</b>



Acero, A. (1990).

*Acoustic and Environmental Robustness in Automatic Speech Recognition.*

PhD thesis, Carnegie Mellon University.



Deng, L., Droppo, J., and Acero, A. (2004).

Enhancement of log mel power spectra of speech using a phase-sensitive model of the acoustic environment and sequential estimation of the corrupting noise.

*IEEE Transactions on Speech and Audio Processing*,  
12(2):133–143.



Droppo, J., Acero, A., and Deng, L. (2002).

Uncertainty decoding with SPLICE for noise robust speech recognition.

In *Proc. ICASSP*. IEEE.



Gales, M. (1995).

*Model-based techniques for noise robust speech recognition.*

PhD thesis, Cambridge University.





Hu, Y. and Huo, Q. (2006).

An HMM compensation approach using unscented transformation for noisy speech recognition.

*Chinese Spoken Language Processing*, pages 346–357.



Li, J., Deng, L., Yu, D., Gong, Y., and Acero, A. (2009).

A unified framework of HMM adaptation with joint compensation of additive and convolutive distortions.

*Computer Speech & Language*, 23(3):389–405.



Liao, H. and Gales, M. (2005).

Joint uncertainty decoding for noise robust speech recognition.

In *Proc. INTERSPEECH*. Citeseer.




Moreno, P., Raj, B., and Stern, R. (1996).

A vector Taylor series approach for environment-independent speech recognition.

In *Proc. ICASSP*, volume 2, pages 733–736. IEEE.





Povey, D., Burget, L., Agarwal, M., Akyazi, P., Kai, F.,  Ghoshal, A., Glembek, O., Goel, N., Karafiát, M., Rastrow, A., Rose, R., Schwarz, P., and Thomas, S. (2011).

The subspace Gaussian mixture model—A structured model for speech recognition.

*Computer Speech & Language*, 25(2):404–439.

