

Speaker and Noise Factorisation on the AURORA4 Task

Yongqiang Wang & Mark Gales

May. 2011



Cambridge University Engineering Department

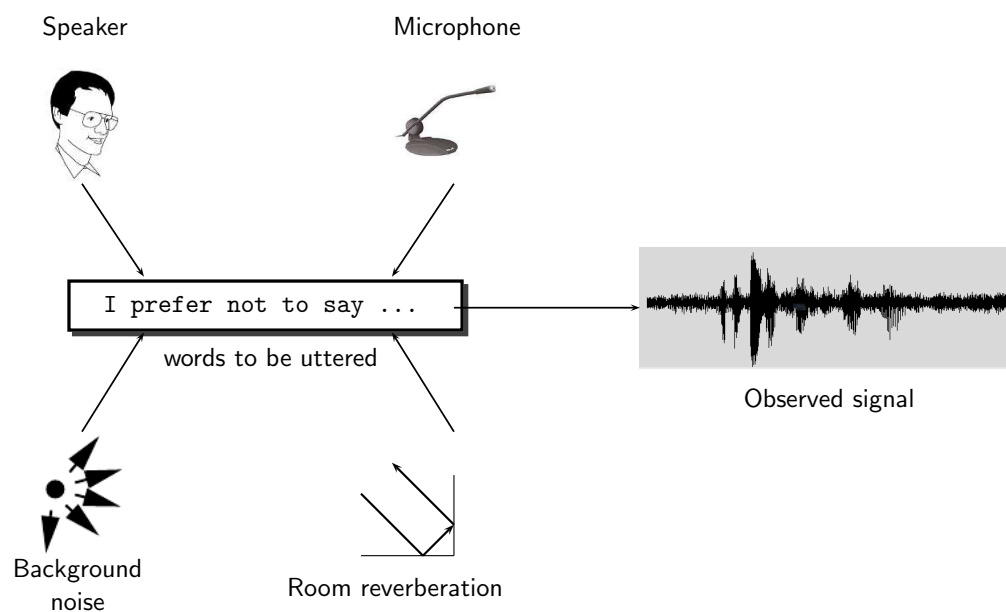
Overview

- Model-based approaches to robust speech recognition
 - Acoustic environment
 - Speaker adaptation schemes
 - Noise robustness schemes
- Handling multiple acoustic factors
 - Acoustic factorisation
 - An example of acoustic factorisation – “Joint” (speaker and noise) adaptation
- Experiments
- Conclusion



Acoustic Environment

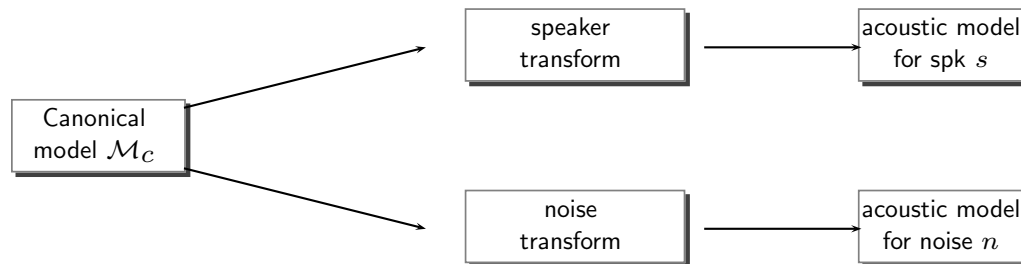
- Speech signals can be influenced by many factors
 - **desired variations**: words to be uttered
 - **unwanted variations**: speaker, noise, channel ...



- This work will consider model-based approach to robust speech recognition

Model-based Framework to Robust Speech Recognition

- *Canonical model* is built to model the desired variations
- *Transforms* are used to adapt canonical model to different acoustic conditions
 - Different transforms are developed to handle specific acoustic factors
 - * Speaker adaptation
 - * Noise compensation



- Transforms can be combined to handle multiple acoustic factors

Speaker Adaptation

- Linear transforms are widely used to adapt the HMM parameters

- Mean transform

$$\boldsymbol{\mu}^{(s)} = \mathbf{A}^{(s)} \boldsymbol{\mu} + \mathbf{b}^{(s)}$$

- Linear transform $\mathbf{W}^{(s)} = [\mathbf{A}^{(s)}, \mathbf{b}^{(s)}]$ represents speaker s 's characteristic

- Large number of parameters to be estimated

- ↳ not possible to estimate transforms robustly from single utterances

- ↳ not suitable for very rapid adaptation

- Originally designed for speaker adaptation

- Also used as general linear transform for environmental adaptation



Noise Compensation Schemes

- A mismatch function is defined for the impact of environment
 - In the cepstral domain, clean speech \mathbf{x} and corrupted speech \mathbf{y} are related:

$$\begin{aligned} \mathbf{y} &= \mathbf{f}(\mathbf{x}, \mathbf{h}, \mathbf{n}) \\ &= \mathbf{C} \log \left(\exp \left(\mathbf{C}^{-1}(\mathbf{x} + \mathbf{h}) \right) + \exp \left(\mathbf{C}^{-1}\mathbf{n} \right) \right) \end{aligned} \quad \left| \begin{array}{l} \mathbf{h} \text{ channel distortion} \\ \mathbf{n} \text{ noise cepstral} \end{array} \right.$$

- Model parameters are compensated using nonlinear transform:
 - VTS approximation

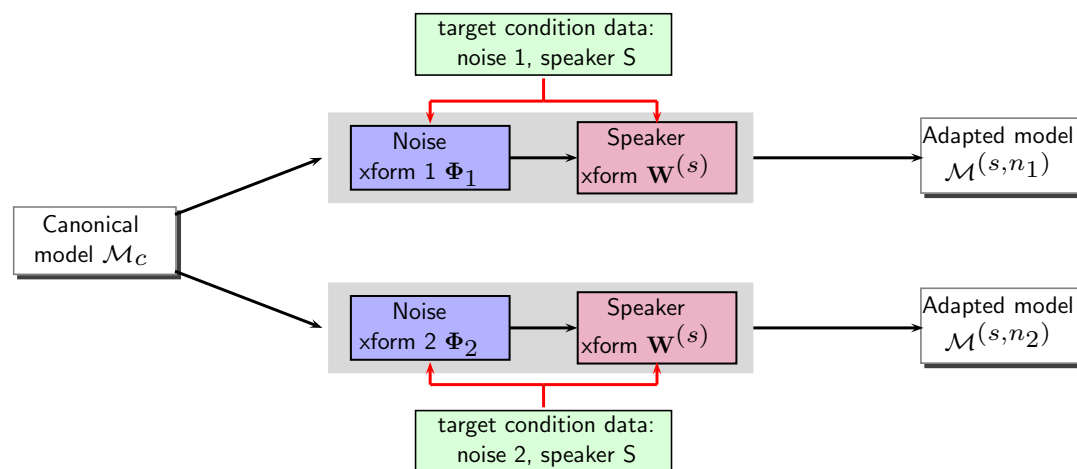
$$\begin{aligned} \boldsymbol{\mu}_y &\approx \mathbf{f}(\boldsymbol{\mu}_x, \boldsymbol{\mu}_h, \boldsymbol{\mu}_n) \\ \boldsymbol{\Sigma}_y &\approx \text{diag} \left(\mathbf{J}_x \boldsymbol{\Sigma}_x \mathbf{J}_x^T + \mathbf{J}_n \boldsymbol{\Sigma}_n \mathbf{J}_n^T \right) \end{aligned} \quad \left| \begin{array}{l} \boldsymbol{\Phi} = (\boldsymbol{\mu}_h, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) \\ \text{noise transform parameter} \end{array} \right.$$

- Noise transform $\boldsymbol{\Phi}$ has a small number of parameters
 - ↳ noise transform parameter can be estimated at per utterance basis
 - ↳ suitable for very rapid adaptation



Speaker and Noise Compensation: Batch Mode

- In real environment, there exists multiple acoustic attributes, e.g., speaker and noise
 - Transforms need to be combined to adapt to the target condition
- Batch mode combination
 - An example: VTS combined with MLLR (“VTS+MLLR”)

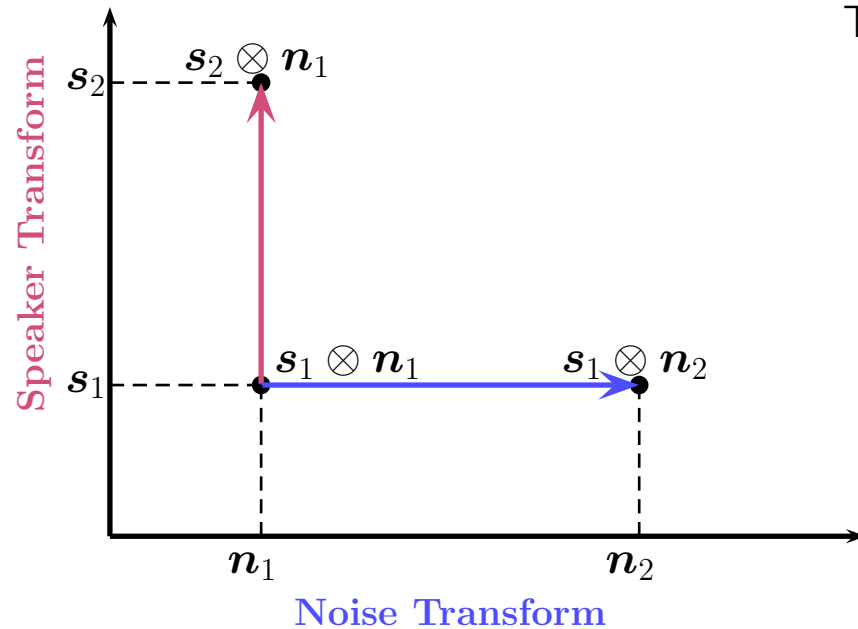


$$\boldsymbol{\mu}_y^{(s)} = \mathbf{A}^{(s)} \mathbf{f}(\boldsymbol{\mu}_x, \boldsymbol{\mu}_h, \boldsymbol{\mu}_n) + \mathbf{b}^{(s)}$$

$$\boldsymbol{\Sigma}_y^{(s)} = \boldsymbol{\Sigma}_y$$

- Transforms are estimated for each speaker and noise conditions
 - * VTS are estimated per utterance
 - * Linear transform $\mathbf{W}^{(s)}$ must be estimated using multiple utterances

Acoustic Factorisation



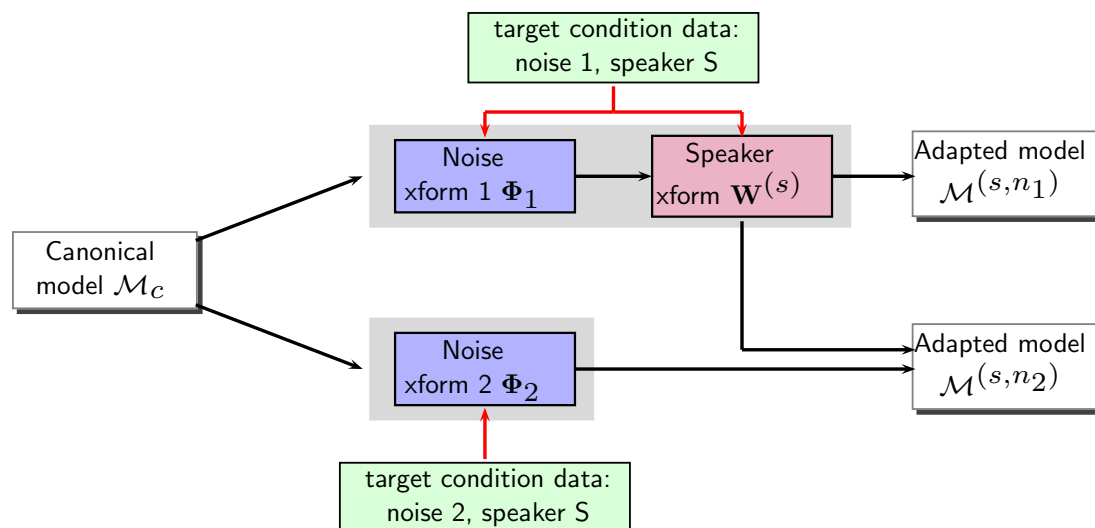
Transforms are "orthogonal":

- From (s_1, n_1) to condition (s_1, n_2) :
 - Speaker transform is locked as s_1
 - Only need to update noise transform

- Flexible transforms:
 - Speaker transform can be used in a range of noise conditions
- Transforms are used in a factorised fashion, but are estimated jointly

Joint Speaker and Noise Compensation

- “Joint” transform: an example of acoustic factorisation



$$\boldsymbol{\mu}_y^{(s)} = \mathbf{f}(\mathbf{A}^{(s)} \boldsymbol{\mu}_x + \mathbf{b}^{(s)}, \boldsymbol{\mu}_h, \boldsymbol{\mu}_n)$$

$$\boldsymbol{\Sigma}_y^{(s)} = \text{diag}(\mathbf{J}_x \boldsymbol{\Sigma}_x \mathbf{J}_x^T + \mathbf{J}_n \boldsymbol{\Sigma}_n \mathbf{J}_n^T)$$

- Configuration investigated:
 - Estimate speaker transform in a particular noise condition (\mathbf{n}_1)
 - Use obtained speaker transform in all noise conditions (\mathbf{n}_1 and \mathbf{n}_2)
 - ↳ Only noise transform need to be estimated for \mathbf{n}_2 , very rapid adaptation

Experiments

- AURORA4 task
 - Derived from WSJ0 5k-word closed vocabulary dictation task
 - 7138 utterances from 83 speakers in training
 - 330 utterances from 8 speakers in testing
 - * 14 test sets are defined

	Headset	Far-field micro.
No noise added	set A (test01)	set C (test08)
Noise added	set B (test02-07)	set D (test09-14)

- Acoustic model training:
 - Front-end: MFCC+C0+1st, 2nd derivatives
 - xwrđ triphone, 3140 clustered states, 16-component system



Batch Mode Experiments

Schemes	set A	set B	set C	set D	Avg.
VTS	6.9	15.1	11.8	23.3	17.8
+MLLR	5.0	12.1	9.0	19.8	14.7
Joint	5.0	12.1	8.6	19.7	14.6

- speaker and noise transforms are estimated for each of the 14 test sets
 - Noise transforms varies from utterance to utterance
- Speaker and noise adaptation outperforms noise adaptation (“VTS”) only
- Similar performances achieved by “Joint” and “VTS+MLLR” transform
- Batch mode adaptation requires a block of utterances to estimate transform
 - ↳ Not very flexible to be used



Factorisation Experiments

Spk. Est.	Schemes	A	B	C	D	Avg.
—	VTS	6.9	15.1	11.8	23.3	17.8
01	+MLLR	5.0	20.2	16.5	28.0	22.2
	Joint	5.0	14.1	10.4	22.3	16.7
04	+MLLR	10.2	19.7	19.7	28.0	22.5
	Joint	7.0	12.5	11.0	20.4	15.4

- Speaker transform and noise transforms are estimated in a factorised mode
 - Speaker transform estimated from test01/04 and fixed for all other sets
 - Noise transform estimated per utterances
- Speaker transform estimated from test04 improves average performance
- Performance of factorised adaptation (15.4%) is close to batch mode (14.6%)



Conclusion

- Handling multiple acoustic factors is important in realistic acoustic environment
- A powerful and flexible approach: acoustic factorisation
- An example of acoustic factorisation is presented: “Joint”
 - Allows very rapid speaker and noise adaptation
 - Speaker transforms are used across multiple noise conditions
 - The power and flexibility of model-based framework is demonstrated
- * “Joint” outperformed advanced front-end with MLLR for both batch and factorisation mode

