

# Model-based approaches to adaptive training in reverberant environments

Y.-Q. Wang and M.J.F. Gales

Engineering Department, Cambridge University  
Trumpington St. Cambridge University, CB2 1PZ, U.K.

{yw293, mjfg}@eng.cam.ac.uk

## Abstract

Adaptive training is a powerful approach for building speech recognition systems using non-homogeneous data. This work presents an extension of model-based adaptive training to handle reverberant environments. The recently proposed Reverberant VTS-Joint (RVTSJ) adaptation[1] is used to factor out unwanted additive and reverberant noise variations in multi-conditional training data, yielding a canonical model neutral to noise conditions. An maximum likelihood estimation of the canonical model parameters is described. An initialisation scheme that uses the VTS-based adaptive training to initialise the model parameters is also presented. Experiments are conducted on a reverberant simulated AURORA4 task.

**Index Terms:** reverberant noise robustness, vector Taylor series, adaptive training

## 1. Introduction

Model-based approaches for noise-robust speech recognition have been investigated and extended in a number of ways, e.g. Vector Taylor series (VTS) compensation[2, 3] and joint uncertain decoding (JUD)[2]. However, there has been less work in applying model-based approaches to handle noise in reverberant environments. In [1, 4], two model compensation schemes, reverberant VTS (RVTS) and reverberant VTS-joint (RVTSJ), are proposed, where the underlying acoustic model is adapted to the target environment given an estimated noise model. The RVTSJ compensation enables the joint estimation of additive and reverberant noise, thus gives better performance. RVTS and RVTSJ schemes assume the underlying acoustic model is trained on clean data. Acoustic models trained on multi-conditional data generally give better noise robustness. In [5], a clean corpus was filtered by several room impulse responses (RIRs) to form the multi-conditional data. It was found that the multi-style trained acoustic model (MST) followed by general adaptation, e.g., MLLR, yielded large gains. Recently, [6] shows that using stereo data in multi-style training can further improve the reverberant noise robustness. However, multi-style training forces the noise variations to be modelled by the underlying acoustic model, which could potentially harm the performance. It also degrades the performance when the acoustic model is operated outside the training environments.

Alternatively, adaptive training can be applied to factor out the unwanted noise variation, yielding a canonical model that is neutral to the noise condition. Adaptive training is originally proposed for speaker adaptive training[7], and is extended to handle additive and convolutional noise recently[2, 8]. Moti-

vated by the success of VTS-based adaptive training (VAT) [8], this work investigates the model-based approaches to adaptive training in reverberant environments. RVTSJ is used to compensate acoustic models in both training and testing. The new training algorithm is referred to reverberant adaptive training (RAT). An maximum likelihood (ML) estimation of canonical model parameters in the EM framework is described. An initialisation scheme using VAT canonical model parameters to start RAT training is also presented. This RAT training algorithm was compared with MST and VAT on a reverberant simulated version of AURORA4.

## 2. Reverberant and additive noise

If clean speech is corrupted by additive and short-time channel distortion (a.k.a. convolutional noise) in the Mel-cepstral domain, the mismatch function describing the distortion on the current clean speech parameter  $\mathbf{x}_t$  is given by :

$$\begin{aligned} \mathbf{y}_t^s &= \mathbf{C} \log \left( \exp(\mathbf{C}^{-1}(\mathbf{x}_t^s + \boldsymbol{\mu}_n)) + \exp(\mathbf{C}^{-1}\mathbf{n}_t^s) \right) \\ &= \mathbf{f}(\mathbf{x}_t^s, \boldsymbol{\mu}_n, \mathbf{n}_t^s), \end{aligned} \quad (1)$$

where the superscript (and also the subscript hereafter)  $s$  denotes static parameter, and  $\mathbf{y}_t$ ,  $\mathbf{h}$  and  $\mathbf{n}_t$  are the noisy speech, convolutional and additive noise, respectively,  $\mathbf{C}$  is the (truncated) DCT matrix. It is usually assumed that  $\mathbf{n}$  is Gaussian distributed with mean  $\boldsymbol{\mu}_n$  and diagonal matrix  $\boldsymbol{\Sigma}_n$ , and  $\mathbf{h}$  is a unknown constant. Note that the above mismatch function relies on the assumption that the effective length of the impulse response of the channel is shorter than the size of the analysis window (typically 25ms). However, in a reverberant environment, due to the late-reflection caused by multiple acoustic paths from the speaker to the microphone, the reverberant time  $T_{60}$  which is the time needed for reflections sound to decay 60dB, is usually ranging from 200ms to 800ms or even longer. This is significantly larger than the size of analysis window. The long reverberant time causes the clean speech not only distorted by the additive noise, but also blurred by several previous frames. In [1], a mismatch function describing the joint effect of additive and reverberant noise is derived:

$$\begin{aligned} \mathbf{z}_t^s &= \mathbf{C} \log \left( \sum_{\delta=0}^n \exp(\mathbf{C}^{-1}(\mathbf{x}_{t-\delta}^s + \boldsymbol{\mu}_{1\delta})) + \exp(\mathbf{C}^{-1}\mathbf{n}_t^s) \right) \\ &= \mathbf{g}(\mathbf{x}_t^s, \dots, \mathbf{x}_{t-n}^s, \boldsymbol{\mu}_1, \mathbf{n}_t^s) \end{aligned} \quad (2)$$

where  $\mathbf{z}_t$  is the additive and reverberant noise corrupted speech,  $\boldsymbol{\mu}_1 = [\boldsymbol{\mu}_{10}^T, \dots, \boldsymbol{\mu}_{1n}^T]^T$  is referred to as reverberant noise.

### 2.1. Model compensation

Given the additive and convolutional noise mismatch function in Eq. (1), VTS can be used to approximate it for every Gaus-

---

This work was partially supported by Google research award and DARPA under the GALE and RATS program. The paper does not necessarily reflect the position or the policy of US Government and no official endorsement should be inferred.

sian component  $m$  in the following form:

$$\mathbf{y}_t^s|m \approx \mathbf{f}(\boldsymbol{\mu}_{\text{sx}}^{(m)}, \boldsymbol{\mu}_n, \boldsymbol{\mu}_{\text{sn}}) + \mathbf{J}_x^{(m)}(\mathbf{x}_t^s - \boldsymbol{\mu}_{\text{sx}}^{(m)}) + \mathbf{J}_n^{(m)}(\mathbf{n}_t^s - \boldsymbol{\mu}_n),$$

where  $\boldsymbol{\mu}_n$  and  $\boldsymbol{\mu}_x^{(m)}$  are the means of additive noise and  $m$ -th Gaussian respectively, and

$$\mathbf{J}_x^{(m)} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\boldsymbol{\mu}_{\text{sx}}^{(m)}, \boldsymbol{\mu}_{\text{sn}}, \boldsymbol{\mu}_n}, \quad \mathbf{J}_n^{(m)} = \mathbf{I} - \mathbf{J}_x^{(m)} \quad (3)$$

This yields the following VTS-based model compensation form for the static parameter:

$$\boldsymbol{\mu}_{\text{sy}}^{(m)} = \mathbf{f}(\boldsymbol{\mu}_{\text{sx}}^{(m)}, \boldsymbol{\mu}_n, \boldsymbol{\mu}_{\text{sn}}), \quad \boldsymbol{\Sigma}_{\text{sy}}^{(m)} = \text{diag}(\mathbf{J}_x^{(m)} \boldsymbol{\Sigma}_{\text{sx}}^{(m)T} \mathbf{J}_x^{(m)}) \quad (4)$$

The dynamic parameter compensation can be derived using the continuous time approximation, i.e.,

$$\boldsymbol{\mu}_{\Delta y}^{(m)} = \mathbf{J}_x^{(m)} \boldsymbol{\mu}_{\Delta x}^{(m)}, \quad \boldsymbol{\Sigma}_{\Delta y}^{(m)} = \text{diag}(\mathbf{J}_x^{(m)} \boldsymbol{\Sigma}_{\Delta x}^{(m)T} \mathbf{J}_x^{(m)}) \quad (5)$$

where the subscript  $\Delta$  denotes the delta parameter. The delta-delta parameter is compensated in a similar way. For notation convenience, only the delta parameter is considered in the following discussion.

For the reverberant and additive noise mismatch function in Eq. (2), the current observation  $\mathbf{z}_t$  is a function of the preceding clean speech frames  $\mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-n}$  (ignoring the dynamic parameters). These clean speech frames should be inferred from the noise corrupted observations. In practice, this is computationally intractable. An approximation form was proposed in [4], where  $\mathbf{z}_t$  is assumed to depend on an extended vector  $\bar{\mathbf{x}}_t$ , which is generated by the current Gaussian component. Figure 1 illustrates the dynamic Bayesian network of this approximated model, where  $q_t$  and  $w_t$  are the indicator of current state and Gaussian component respectively.

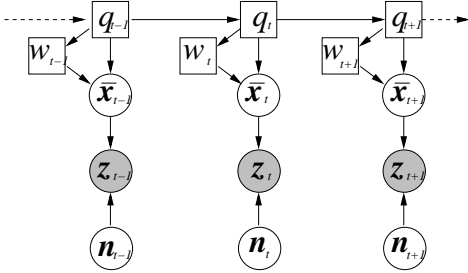


Figure 1: Approximate reverberant dynamic Bayesian network. For clarity, the dynamic parameters are ignored.

The form of  $\bar{\mathbf{x}}_t$  is chosen such that when there is no reverberant noise, the compensated model will back off to the standard VTS compensation, and the Gaussian distribution is used to model  $\bar{\mathbf{x}}_t$  conditioning on the current component  $m$ , i.e.,

$$\bar{\mathbf{x}}_t = \begin{bmatrix} \mathbf{x}_t \\ \Delta \mathbf{x}_t \\ \Delta^2 \mathbf{x}_t \\ \Delta^3 \mathbf{x}_t \\ \vdots \end{bmatrix} = \mathbf{W} \begin{bmatrix} \mathbf{x}_{t+w} \\ \dots \\ \mathbf{x}_{t-n-w} \end{bmatrix}; \quad (6)$$

and  $\bar{\mathbf{x}}_t|m \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_x^{(m)}, \bar{\boldsymbol{\Sigma}}_x^{(m)})$ , where  $w$  is the window size to calculate the dynamic parameter,  $\mathbf{W}$  is a square and invertible matrix, which maps a sequence of statics to static plus the first, second and higher order dynamics. Given the model statistics  $\bar{\boldsymbol{\mu}}_x^{(m)}$ , it is easy to derive the statistics of spliced vector  $\mathbf{x}_e^s = (\mathbf{x}_t^s, \dots, \mathbf{x}_{t-n}^s)^T$ . For example, the static and delta

mean vectors of  $\mathbf{x}_e$  can be obtained by

$$\boldsymbol{\mu}_{\text{sx}_e}^{(m)} = \mathcal{E}\{\mathbf{x}_e^s|m\} = \mathbf{P}_s \bar{\boldsymbol{\mu}}_x^{(m)}, \quad \boldsymbol{\mu}_{\Delta \mathbf{x}_e}^{(m)} = \mathcal{E}\{\Delta \mathbf{x}_e|m\} = \mathbf{P}_\Delta \bar{\boldsymbol{\mu}}_x^{(m)} \quad (7)$$

where  $\mathbf{P}_s$  and  $\mathbf{P}_\Delta$  are the matrices that map  $\bar{\mathbf{x}}_t$  to  $\mathbf{x}_e^s$  and  $\Delta \mathbf{x}_e$ .

Using statistics of the extended vector, VTS is extended in [1] to handle the mismatch function in Eq. (2). The expansion is performed at  $(\boldsymbol{\mu}_{\text{sx}_e}^{(m)}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_n)$ :

$$\mathbf{z}_t^s|m \approx \mathbf{g}(\boldsymbol{\mu}_{\text{sx}_e}^{(m)}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_n) + [\mathbf{J}_{\text{xe}}^{(m)}, \mathbf{J}_{\text{ne}}^{(m)}] \begin{bmatrix} \mathbf{x}_e^s - \boldsymbol{\mu}_{\text{sx}_e}^{(m)} \\ \mathbf{n}_t^s - \boldsymbol{\mu}_n \end{bmatrix} \quad (8)$$

where

$$\mathbf{J}_{\text{xe}}^{(m)} = [\mathbf{J}_{\text{x0}}^{(m)}, \dots, \mathbf{J}_{\text{xn}}^{(m)}]; \quad \mathbf{J}_{\text{ne}}^{(m)} = \mathbf{I} - \sum_{\delta=0}^n \mathbf{J}_{\text{x}\delta}^{(m)} \quad (9)$$

and  $\mathbf{J}_{\text{x}\delta}^{(m)} = \frac{\partial \mathbf{g}}{\partial \mathbf{x}_e^s} \Big|_{\boldsymbol{\mu}_{\text{sx}_e}^{(m)}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_n}$ . Along with the continuous time approximation assumption, this yields the mean compensation form in the following:

$$\boldsymbol{\mu}_{\text{sz}}^{(m)} = \mathbf{g}(\boldsymbol{\mu}_{\text{sx}_e}^{(m)}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_n), \quad \boldsymbol{\mu}_{\Delta \mathbf{z}}^{(m)} = \sum_{\delta=0}^n \mathbf{J}_{\text{x}\delta}^{(m)} \boldsymbol{\mu}_{\Delta \mathbf{x}\delta}^{(m)} \quad (10)$$

This is referred to as RVT SJ. It is possible to compensate the variance as well. However, in the initial investigation in [1], it was found that variance compensation is not quite effective, thus the standard VTS variance compensation was used, i.e.,  $\boldsymbol{\Sigma}_z^{(m)} = \boldsymbol{\Sigma}_y^{(m)}$ . This is also adopted in this work.

## 2.2. Noise estimation

Given the acoustic model parameter  $\mathcal{M} = \{\bar{\boldsymbol{\mu}}_x^{(m)}, \boldsymbol{\Sigma}_x^{(m)}\}$ , the noise model parameter  $\Phi = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_n)$  is estimated using EM. The following auxiliary function is maximised, i.e.,

$$\hat{\Phi} = \arg \max \sum_{t,m} \gamma_t^{(m)} \log p(\mathbf{z}_t; \boldsymbol{\mu}_z^{(m)}, \boldsymbol{\Sigma}_z^{(m)}) \quad (11)$$

where  $\gamma_t^{(m)}$  is the posterior of component  $m$  at time  $t$ , given the current hypothesis and current noise estimates  $\Phi$ . VTS is again used to expand  $\boldsymbol{\mu}_z^{(m)}$  using the current noise estimates  $\Phi$ , yielding the following update formula using the second-order method:

$$\begin{bmatrix} \hat{\boldsymbol{\mu}}_1 \\ \hat{\boldsymbol{\mu}}_n \end{bmatrix} = \left( \sum_{t,m} \gamma_t^{(m)} \mathbf{J}^{(m)T} \boldsymbol{\Sigma}_{\text{sz}}^{(m)-1} \mathbf{J}^{(m)} + \alpha \mathbf{I} \right)^{-1} \times \left( \sum_{t,m} \gamma_t^{(m)} \mathbf{J}^{(m)T} \boldsymbol{\Sigma}_{\text{sz}}^{(m)-1} \left( \boldsymbol{\mu}_{\text{sz}}^{(m)} - \mathbf{J}^{(m)} \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_n \end{bmatrix} \right) \right) \quad (12)$$

where

$$\mathbf{J}^{(m)} = [\mathbf{J}_{10}^{(m)}, \dots, \mathbf{J}_{1n}^{(m)}, \mathbf{J}_{\text{ne}}^{(m)}], \quad \mathbf{J}_{1\delta}^{(m)} = \frac{\partial \mathbf{g}}{\partial \boldsymbol{\mu}_{1\delta}} \Big|_{\boldsymbol{\mu}_{\text{sx}_e}^{(m)}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_n}$$

and  $\alpha$  is used to improve the stability of noise estimation. Note that the above updating formula only considers the static parameters in the auxiliary function. It is possible to include all the compensated parameters. A good initialisation of noise parameters is important. In [1], the standard VTS-based noise estimation and a guess of  $T_{60}$  value were used to initialise the reverberant and additive noise. This is also adopted in this work. During the noise estimation, it is also important to ensure that every update increases the auxiliary. More details about the noise estimation is given in [1].

## 3. Reverberant adaptive training

The above noise estimation assumes that the acoustic model is trained from clean data. The adaptive training framework can be applied in which both the acoustic model and the noise model

are trained in a full ML framework on multi-conditional data. This is a powerful technique to factor out the unwanted acoustic factors, such as speaker differences and noise distortions, yielding a canonical model  $\mathcal{M}_c$  that models only the relevant phoneme variations. The adaptive training framework is extended in this work to handle to reverberant and additive noise distortions. This is referred to as reverberant adaptive training (RAT).

In adaptive training, both the canonical model  $\mathcal{M}_c$  and a set of noise models  $\Phi$  are iteratively estimated using EM. First, give the current canonical model, the noise models  $\Phi$  are estimated for each utterance<sup>1</sup>, then the canonical model  $\mathcal{M}_c$  is updated given the current noise models. Multiple iterations may be performed to interleave optimisation in the EM framework. With RAT, the following auxiliary is used:

$$\mathcal{Q}(\mathcal{M}_c, \{\Phi^{(u)}\}) = \sum_{u,t,m} \gamma_t^{(mu)} \log p(\mathbf{z}_t^{(u)}; \hat{\boldsymbol{\mu}}_z^{(mu)}, \hat{\boldsymbol{\Sigma}}_z^{(mu)}) \quad (13)$$

where  $u$  is the index of utterance. For example,  $\gamma_t^{(mu)}$  is the posterior of component  $m$  at time  $t$  for the  $u$ -th utterance.

Given the canonical model, estimating the reverberant and additive noise parameter is described in section 2.2. After the noise parameter updated, the canonical model parameter is re-trained. As in RVTSJ adaptation, only the mean is compensated for reverberation, in RAT, only the extended mean  $\bar{\boldsymbol{\mu}}_x^{(m)}$  will be updated. The auxiliary function, where only terms dependent on the  $\bar{\boldsymbol{\mu}}_x^{(m)}$  are shown, is:

$$\mathcal{Q}(\hat{\mathcal{M}}_c; \mathcal{M}_c) = -\frac{1}{2} \sum_u \sum_m \gamma^{(mu)} \hat{\boldsymbol{\mu}}_z^{(mu)\top} \boldsymbol{\Sigma}_z^{(mu)-1} \left( \hat{\boldsymbol{\mu}}_z^{(mu)} - 2\boldsymbol{\Gamma}_z^{(mu)} \right) \quad (14)$$

where  $\gamma^{(mu)} = \sum_t \gamma_t^{(mu)}$  and  $\boldsymbol{\Gamma}_z^{(mu)} = \frac{1}{\gamma^{(mu)}} \sum_t \gamma_t^{(mu)} \mathbf{z}_t$ . Similar as in the noise estimation, the VTS is applied to expand  $\hat{\boldsymbol{\mu}}_{sz}^{(m)}$  using the current canonical model estimates, i.e.,

$$\hat{\boldsymbol{\mu}}_z^{(m)} \approx \begin{bmatrix} \boldsymbol{\mu}_{sz}^{(m)} \\ \boldsymbol{\mu}_{\Delta z}^{(m)} \end{bmatrix} + \begin{bmatrix} \mathbf{J}_{x_e}^{(mu)} & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_{x_e}^{(mu)} \end{bmatrix} \begin{bmatrix} \mathbf{P}_s \\ \mathbf{P}_\Delta \end{bmatrix} (\hat{\boldsymbol{\mu}}_x^{(m)} - \bar{\boldsymbol{\mu}}_x^{(m)}) \quad (15)$$

Differentiating the auxiliary function and equating to zero gives the following update:

$$\begin{aligned} \hat{\bar{\boldsymbol{\mu}}}_x^{(m)} &= \bar{\boldsymbol{\mu}}_x^{(m)} + \zeta \left( \sum_u \gamma^{(mu)} \mathbf{K}^{(mu)\top} \boldsymbol{\Sigma}_z^{(mu)-1} \mathbf{K}^{(mu)\top} + \beta \mathbf{I} \right)^{-1} \\ &\times \left( \sum_u \gamma^{(mu)} \mathbf{K}^{(mu)\top} \boldsymbol{\Sigma}_z^{(mu)-1} (\boldsymbol{\Gamma}_z^{(mu)} - \boldsymbol{\mu}_z^{(mu)}) \right) \quad (16) \end{aligned}$$

where  $\mathbf{K}^{(mu)} = \begin{bmatrix} \mathbf{J}_{x_e}^{(mu)} \mathbf{P}_s \\ \mathbf{J}_{x_e}^{(mu)} \mathbf{P}_\Delta \end{bmatrix}$ ,  $\beta$  is a parameter used to stabilise the canonical model parameter update, and  $\zeta$  is the step size. For every update, the step size  $\zeta$  is set as 1 initially. Due to the approximation made in Eq. (15), it is necessary to check the auxiliary function after every update to ensure the auxiliary is increasing. If not, a simple back-off procedure, similar to the one used in [2], is used to reduce the step size until the auxiliary increases. Since the auxiliary function of canonical model involves all the utterance, to make this back-off procedure possible, it is necessary to save  $\boldsymbol{\Gamma}_z^{(mu)}$  for each component  $m$  and each utterance  $u$ , provided  $\gamma^{(mu)}$  is not zero. This is impractical for media/large vocabulary task. An approxima-

tion of the auxiliary function in Eq. (14) is used: for each component update, the summation over  $u$  is done on a subset  $\mathcal{U}_m = \{u | \gamma^{(mu)} \geq \theta_m\}$ , where  $\theta_m$  is choose such that the top  $N$  utterances are in this subset.  $N = 156$  is used in this work and it was found this yields good increase in likelihood.

Since the auxiliary function is highly non-linear, it is crucial to have a good initialisation of the canonical model parameters. In this work, the following strategy is used: the standard model parameters,  $\boldsymbol{\mu}_x^{(m)}, \boldsymbol{\Sigma}_x^{(m)}$  are set as the parameter obtained by the VTS-based adaptive training (VAT), whilst for  $\bar{\boldsymbol{\mu}}_x^{(m)}$ , it is assumed that  $\boldsymbol{\mu}_{x\delta}^{(m)}, \delta = 0 \dots n$  is a smooth trajectory starting from  $\boldsymbol{\mu}_{x0}^{(m)} = \boldsymbol{\mu}_x^{(m)}$ ; hence the reconstruction error is minimised. This amounts to the following optimisation problem :

$$\begin{aligned} \min \quad & \sum_{\delta=1}^n w_\delta \|\mathbf{Q}_\delta \bar{\boldsymbol{\mu}}_x^{(m)} - \boldsymbol{\mu}_x^{(m)}\|^2 \\ \text{s.t.} \quad & \mathbf{Q}_0 \bar{\boldsymbol{\mu}}_x^{(m)} = \boldsymbol{\mu}_x^{(m)} \end{aligned} \quad (17)$$

where  $\mathbf{Q}_\delta$  is the matrix that maps  $\bar{\boldsymbol{\mu}}_x^{(m)}$  to  $\boldsymbol{\mu}_{x\delta}^{(m)}$ ,  $w_\delta = 10^{-\delta \frac{\Delta}{T_{60}}}$ ,  $\Delta$  is the shift of the analysis window ( 10ms ), and  $T_{60}$  the median of reverberation time in the multi-style training data ( 400ms in this work ).

## 4. Experiments

A reverberant version of the AURORA4 task was used for evaluation. The original AURORA4 task is derived from Wall Street Journal (WSJ0) 5k-word dictation task. The WSJ0 training set, consisted of 7138 utterances from 83 speakers, recorded by close-talking microphones, were used as the clean training set in the experiments. To create a multi-conditional training set with reverberation and background additive noise, the clean training set was passed through the simulation tool in [9]. Two RIRs, recorded in an office environment (“office1”) and a living room environment (“liv1”) were used to filter the clean training set, with the reverberant time ranging from 200ms to 600ms. 6 types of background noises which were used in AURORA4 task were also added, with the SNR ranging from 10dB to 20dB, matching the configuration in AURORA4 task. For the test sets, the 330 utterance from 8 speakers in the AURORA4 set A were filtered by two RIRs, “office1” and “office2”, where the latter was recorded in another office environment and not observed during training. For each RIR, there were two background noise conditions, “clean” and “restaurant”, where for the latter, the noise from test04 in AURORA4 task were extracted and added to the reverberant signal at the SNR ranging from 5dB to 15dB. Note the creation of these sets were different from the method in [1], where the reverberant noise was added after background noise distortion.

The HTK frontend was used to derive a 39-dimensional feature vector, consisting of 12 MFCCs, extracted from magnitude spectrum, appended with zeroth cepstrum, delta and delta-delta coefficients. A Cross-word triphone model with 3140 tied states and 16 components per state was built. This model topology was used for all the acoustic models throughout the experiments. For the extended model statistics  $\bar{\boldsymbol{\mu}}_x^{(m)}$ , the feature vector was appended with high-order DCT elements of an appropriate window width.  $n = 10, w = 4$  were used as the length of history frames and the window length used for calculating the dynamic parameters, respectively. The standard bi-gram LM for the AURORA4 task was used in decoding. All the adaptation in this work were performed in a unsupervised mode and the noise models were all estimated at the utterance level.

Experiments were first run using the clean-trained acoustic

<sup>1</sup>It is assumed in this work, each utterance has a unique noise condition, thus a homogeneous block.

noise condition		adaptation		
rev.	add.	—	VTS	RVTSJ
—	clean	7.1	6.9	7.3
	rest.	60.0	19.5	20.0
office1	clean	70.3	43.7	25.4
	rest.	97.3	51.6	43.1
office2	clean	60.1	30.9	16.5
	rest.	97.6	48.8	47.1

Table 1: Performance (in WER %) of the clean-trained acoustic model operated in a noisy and/or reverberant environment.

model. VTS-based noise model parameter was first estimated using multiple EM iterations. The acoustic model was then compensated and used to generate supervision hypothesis. This supervision hypothesis was used for updating VTS noise model and estimating the RVTSJ noise model as well. For comparison, experiments were also run on the 01(clean condition) and 04(restaurant noise, no reverberation) sets in the AURORA4 task. Results are shown in Table 1. As expected, the clean-trained acoustic model is quite fragile to the environment: its performance were greatly impacted by the additive and/or the reverberant noise while reverberation seems to be a more detrimental factor than additive noise. Recognition in a noisy and reverberant environment is the most challenging task, as both distortions cause large mismatch between the training and testing data. Performing VTS compensation significantly reduced the mismatch caused by noise and reverberation. RVTSJ adaptation on non-reverberant data gave similar but slightly worse performance (less than 0.5% absolute degradation). This is felt to be a limitation of current noise estimation method when the reverberant noise ( $\mu_{1\delta}, \delta = 1 \dots n$ ) approaches to  $-\infty$ . However, when the reverberation is presented in the data, RVTSJ gave large gains over VTS. This demonstrated that RVTSJ is modelling the impact of reverberation, which is not modelled well by VTS.

In the second set of experiments, multi-conditional training data was used to build acoustic models. Firstly, as in [6], stereo data were used to build a MST system. Starting from this MST model, VAT system was build, which in turn serves as an initialisation for the RAT to begin. The parameter  $\beta$  was reduced from 8 to 1 for 4 iterations of RAT canonical model re-estimation. This was followed by 2 iterations of noise model update. This process is repeated one more time to yield the final RAT model. An initial decoding using MST system without adaptation was run. Results are shown in line 2, Table 1. Compared with performance in Table 1, in the environment observed during training, “office1”, MST model works quite well, producing better performance than adapting clean-trained acoustic models to the target environment. However, when it was operated in “office2”, an environment not observed during training, performances were degraded. This demonstrates multi-style training introduces bias toward the training environments while not generalised well to other unseen conditions. The MST system was also adapted by an CMLLR transform for each speaker. As a general adaptation scheme, CMLLR is powerful to reduce the mismatch at testing. As shown in line 3, Table 2, this yields large error reduction in all the environments. This was consistent with the finding in [6, 5]. The hypothesis obtained by MST+CMLLR will be used as supervisions for the following VTS/RVTSJ adaptation experiments. Compared with the CMLLR adapted MST system, VAT system with VTS adaptation was worse performed when there is only reverberation distortion, but gave gains when the additive noise is also presented in

Systems	Adaptation	office1		office2		Avg.
		clean	rest.	clean	rest.	
MST	—	19.9	38.3	37.4	63.6	39.8
	CMLLR	14.1	30.1	14.7	48.4	26.8
VAT	VTS	15.1	29.0	18.7	44.5	26.8
	RVTSJ	14.9	29.6	18.3	43.6	26.6
RAT	RVTSJ	13.7	28.5	15.0	42.2	24.9

Table 2: Performance of multi-style trained and adaptively trained acoustic models in reverberant environments.

the environment. This is due to that VTS is designed to compensate the impact of additive noise, while CMLLR can be used for general adaption. Based on the VAT canonical model, the extended model statistics  $\{\bar{\mu}_x^{(m)}\}$  were initialised by solving the optimisation in Eq. (17). Given the extended model statistics, RVTSJ adaptation of VAT was performed, which gave small gains in average (0.2%). RVTSJ adaptation of the RAT system further improves the performance. Compared with VAT system, RAT yields 0.5% to 1.3% absolute gains for the office1 environment, and 2.3% to 3.7% for the office2 environment. This demonstrates that RVTSJ adaptation models the impact of both the reverberant and additive noise, while RAT produces a canonical model neutral to these distortions to some extent. On average, RAT system gave the best performance, yielding 1.9% absolute gains over CMLLR adapted MST system.

## 5. Conclusions

This work investigates model based approaches to adaptive training in reverberant environments. A new adaptive training algorithm, reverberant adaptive training (RAT), is proposed, where the RVTSJ adaptation is used in both testing and training stage. An ML estimation of canonical model parameters in the EM framework is presented. Experiments conducted on a reverberant simulated AURORA4 task demonstrates RAT produces a canonical model neutral (to some extent) to the reverberant and additive noise variations and provides better results than adapting the multi-style trained acoustic model using CMLLR.

## 6. References

- [1] Y.-Q. Wang and M. J. F. Gales, “Improving reverberant VTS for hands-free robust speech recognition,” in *Proc. ASRU-2011*.
- [2] H. Liao and M. J. F. Gales, “Joint uncertainty decoding for robust large vocabulary speech recognition,” Tech. Rep. CUED/F-INFENG/TR552, University of Cambridge, 2006.
- [3] J. Li, D. Yu, Y. Gong, and A. Acero, “High-performance HMM adaptation with joint compensation of additive and convolutive distortions via vector Taylor series,” in *Proc. ASRU-2007*.
- [4] M. J. F. Gales and Y.-Q. Wang, “Model-based approaches to handling additive noise in reverberant environments,” in *Workshop on Hands-free Speech Communication and Microphone Arrays*, 2011.
- [5] M. Matassoni et al., “Hidden Markov model training with contaminated speech material for distant-talking speech recognition,” *Computer Speech & Language*, vol. 16, no. 2, pp. 205–223, 2002.
- [6] A. Sehr et al., “Multi-style training of HMMs with stereo data for reverberation-robust speech recognition,” in *Workshop on Hands-free Speech Communication and Microphone Arrays*, 2011.
- [7] T. Anastasakos, J. McDonough, R. Schwartz, and J. Makhoul, “A compact model for speaker adaptive training,” in *Proc. ICSLP-96*, pp. 1137–1140.
- [8] O. Kalinli, M. L. Seltzer, and A. Acero, “Noise adaptive training using a vector Taylor series approach for noise robust automatic speech recognition,” in *Proc. ICASSP-2009*.
- [9] H. G. Hirsch and H. Finster, “The simulation of realistic acoustic input scenarios for speech recognition systems,” in *Proc. ICSLP*, 2005.