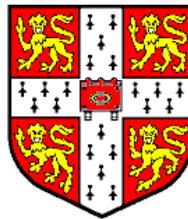


Model-Based Approaches to Speaker and Environment Adaptation

Mark Gales

April 2009



Cambridge University Engineering Department

Overview

- Speaker Adaptation - “Adaptive”
 - linear transform-based adaptation / adaptive training
- Extensions to Linear-Transform Approaches
 - Bayesian adaptive training and inference
 - discriminative mapping functions
 - noisy constrained MLLR
- Noise Robust Speech Recognition - “Predictive”
 - model-based approaches / ML noise estimation
- Extensions to Model-Based Approaches
 - joint uncertainty decoding
 - predictive linear transforms
 - adaptive training / incremental adaptation

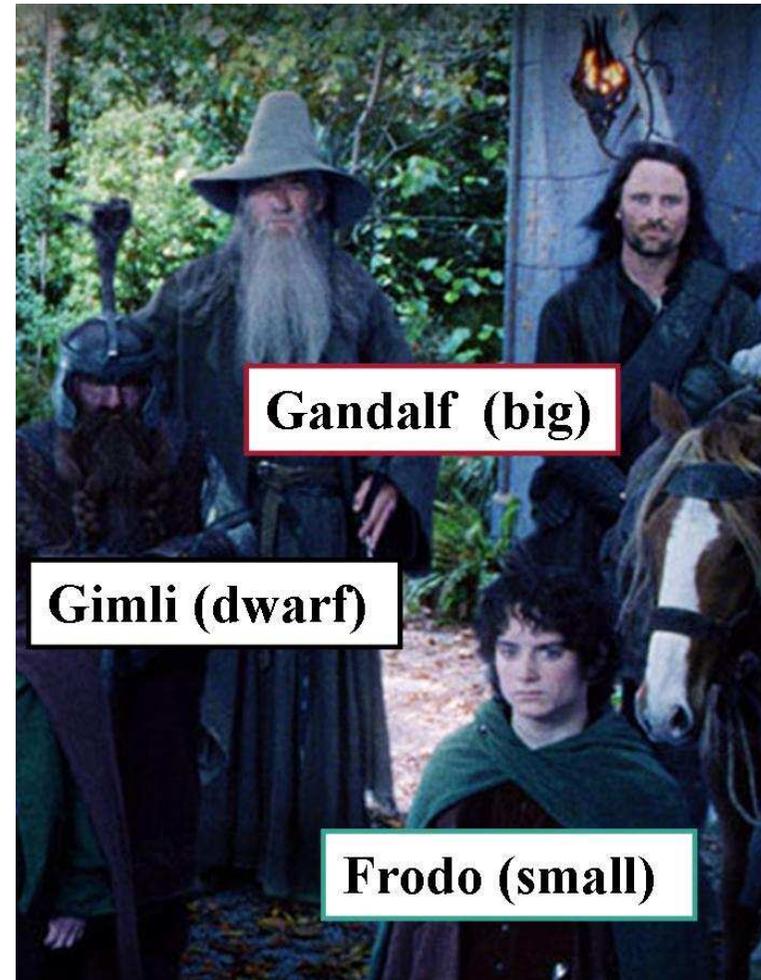


Speaker Adaptation



Speaker Adaptation

- Large differences between speakers
- Linguistic Differences e.g.
 - Accents
tomato in RP/American English
 - Speaker idiosyncrasies
either in English
 - non-native speaker
- Physiological Differences e.g.
 - physical attributes - gender, length of vocal tract
 - transitory effects
cold/stress/public speaking



Adaptation Modes

- Speaker/environment adaptation is an essential part of LVCSR systems
 - obtain the performance of a Speaker/Environment dependent system with orders-of-magnitude less data (30 seconds vs 2000 hours!)
- The **mode** of adaptation depends on the task being investigated
 - **incremental**: results are required causally, the adaptation data is not all available in one block - dictation tasks, car navigation
 - **batch**: all the data is available (or can be used) in one block - BN transcription, CTS transcription

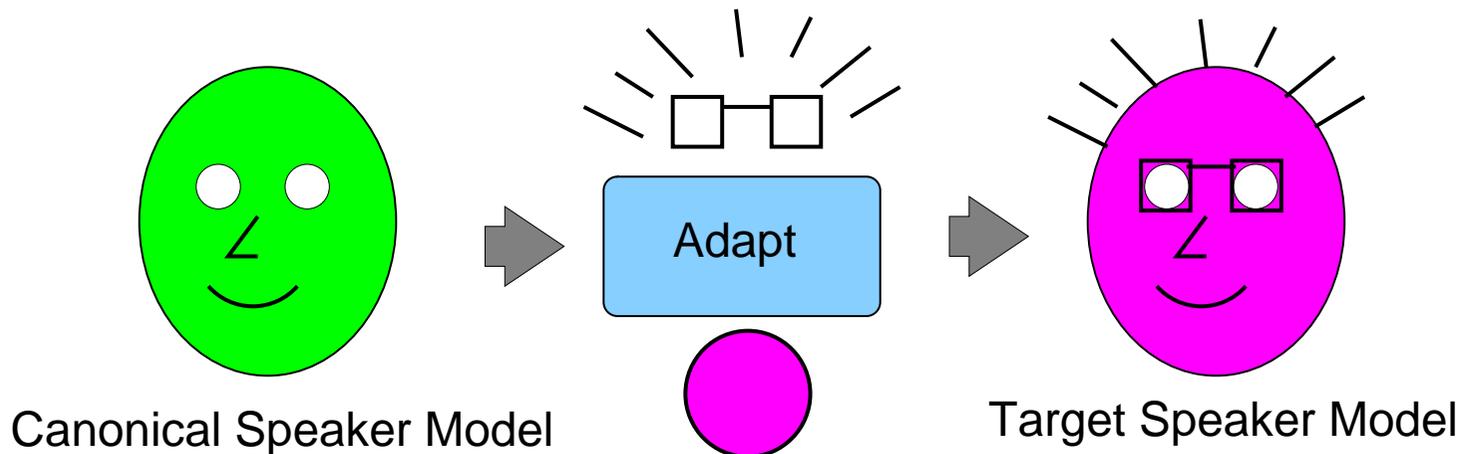
In addition for batch adaptation the adaptation data may be

- **supervised**: the correct transcription of the adaptation data is known (dictation enrolment)
- **unsupervised**: no transcribed adaptation data available, transcription must be hypothesised (BN transcription)



General Adaptation Process

- **Aim:** Modify a “canonical” model to represent a target speaker
 - transformation should require minimal data from the target speaker
 - adapted model should accurately represent target speaker



- Need to determine
 - nature (and complexity) of the speaker transform
 - how to train the “canonical” model that is adapted

Form of the Adaptation Transform

- There are a number of standard forms in the literature [1].
- **Maximum A-Posteriori** MAP [2] adaptation: general “robust” estimation
 - in simplest form only adapts “seen” components
- **Speaker Clustering**: Gender-dependent (GD) models are the simplest from:
 - often estimated using MAP adaptation with speaker-independent priorsEigenVoices[3], CAT [4] are more complex forms.
- **Vocal Tract Length Normalisation**: motivated from physiological perspective
- **Linear Transform Adaptation**: dominant form for LVCSR
 - will be the focus of this part of the talk



Form of the Adaptation Transform

- Dominant form for LVCSR are ML-based linear transformations
 - MLLR adaptation of the means [5]

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

- MLLR adaptation of the covariance matrices [6, 7]

$$\boldsymbol{\Sigma}^{(s)} = \mathbf{H}^{(s)} \boldsymbol{\Sigma} \mathbf{H}^{(s)\top}$$

- Constrained MLLR adaptation [7]

$$\boldsymbol{\mu}^{(s)} = \mathbf{A}^{(s)} \boldsymbol{\mu} + \mathbf{b}^{(s)}; \quad \boldsymbol{\Sigma}^{(s)} = \mathbf{A}^{(s)} \boldsymbol{\Sigma} \mathbf{A}^{(s)\top}$$

- Forms may be combined into a hierarchy [8] e.g.

CMLLR \rightarrow MLLRMEAN



ML and MAP Linear Transforms

- Transforms often estimated using ML (with hypothesis \mathcal{H})

$$\mathbf{W}_{\text{ml}}^{(s)} = \arg \max_{\mathbf{W}} \left\{ p(\mathbf{O}^{(s)} | \mathcal{H}; \mathbf{W}) \right\}$$

- where $\mathbf{W}_{\text{ml}}^{(s)} = \begin{bmatrix} \mathbf{A}_{\text{ml}}^{(s)} & \mathbf{b}_{\text{ml}}^{(s)} \end{bmatrix}$
- however not robust to limited training data

- Including transform prior, $p(\mathbf{W})$, to get MAP estimate [9]

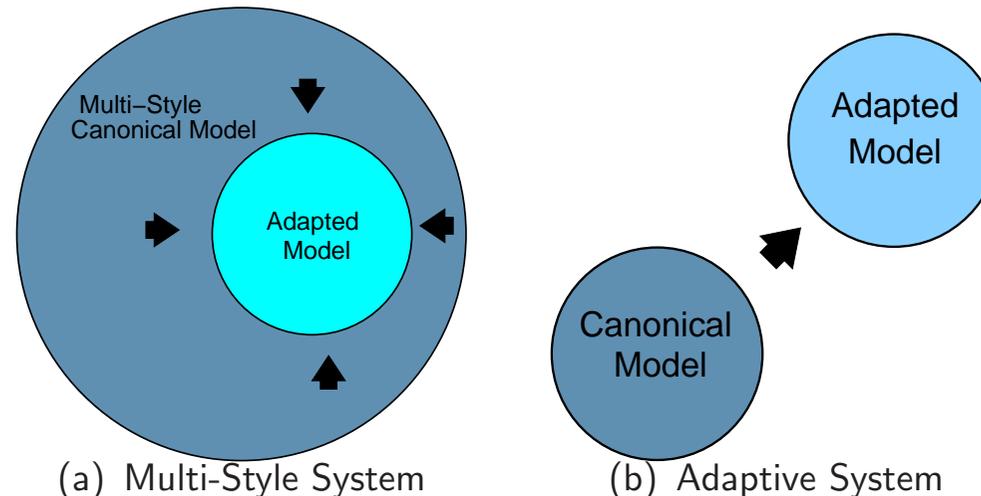
$$\mathbf{W}_{\text{map}}^{(s)} = \arg \max_{\mathbf{W}} \left\{ p(\mathbf{O}^{(s)} | \mathcal{H}; \mathbf{W}) p(\mathbf{W}) \right\}$$

- for MLLR Gaussian is a Gaussian prior for the auxiliary function
- CMLLR prior more challenging ...
- Both approaches rely on expectation-maximisation (EM)



Training a “Good” Canonical Model

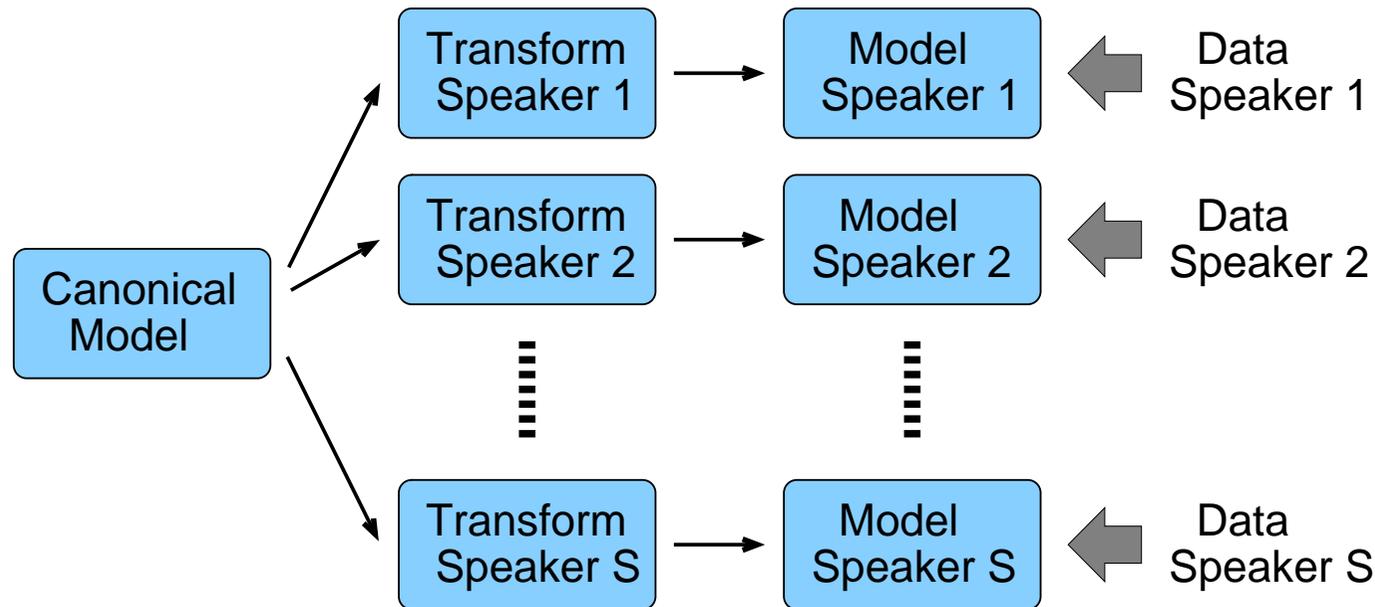
- Standard “multi-style” canonical model
 - treats all the data as a single “homogeneous” block
 - model represents acoustic realisation of phones/words (desired)
 - **and** acoustic environment, speaker, speaking style variations (unwanted)



Two different forms of canonical model:

- **Multi-Style**: adaptation converts a general system to a specific condition;
- **Adaptive**: adaptation converts a “neutral” system to a specific condition [10, 7]

Adaptive Training



- In adaptive training the training corpus is split into “homogeneous” blocks
 - use adaptation transforms to represent unwanted acoustic factors
 - canonical model **only** represents desired variability
- All forms of linear transform can be used for adaptive training
 - CMLLR adaptive training highly efficient

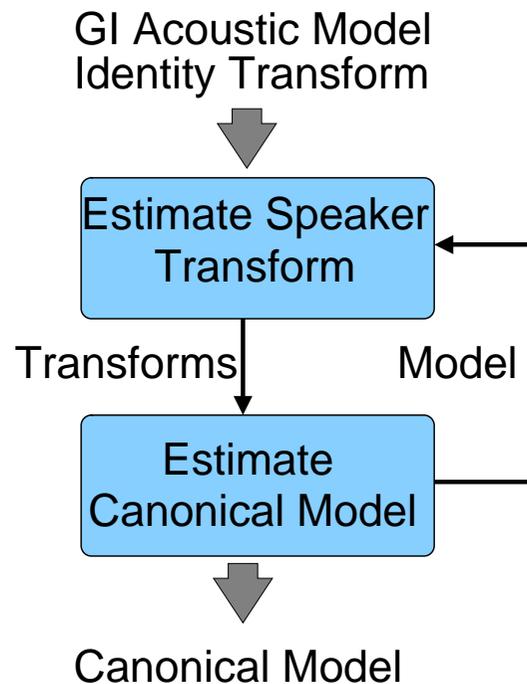


CMLLR Adaptive Training

- The CMLLR likelihood may be expressed as [7]:

$$\mathcal{N}(\mathbf{o}_t; \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top) = \frac{1}{|\mathbf{A}|} \mathcal{N}(\mathbf{A}^{-1}\mathbf{o}_t - \mathbf{A}^{-1}\mathbf{b}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

same as feature normalisation - simply train model in transformed space

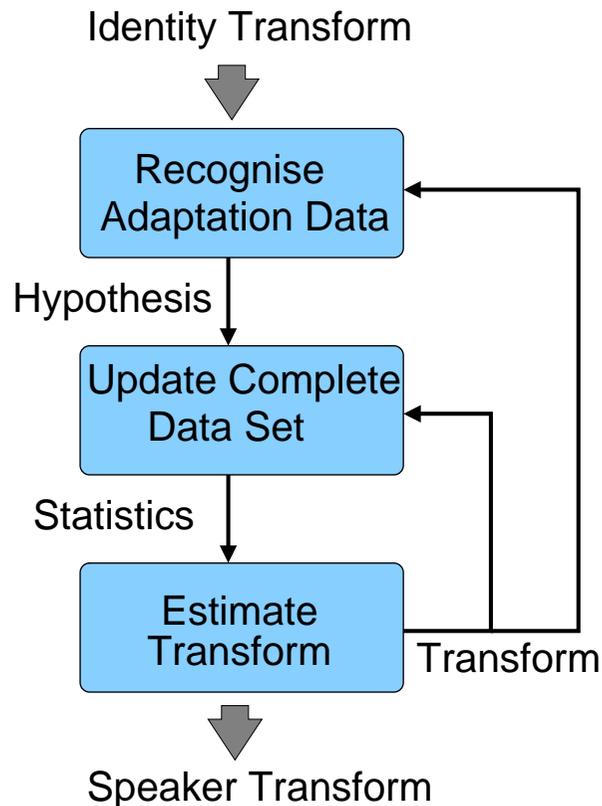


- Interleave Model and transform estimation
- Adaptive canonical model not suited for unadapted initial decode
 - GI model used for initial hypothesis
- MLLR less efficient, but still reasonable



Unsupervised Linear Transformation Estimation

- Estimation of all the transforms is based on EM:
 - requires the **transcription/hypothesis** of the adaptation data
 - iterative process using “current” transform to estimate new transform

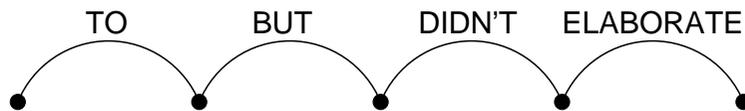


- Two iterative loops for estimation:
 1. estimate hypothesis given transform
 2. update complete-dataset given transform and hypothesisreferred to as **Iterative MLLR** [11]
- For supervised training hypothesis is known
- Confidence-scores can also be used
 - **confidence-based MLLR** [12]

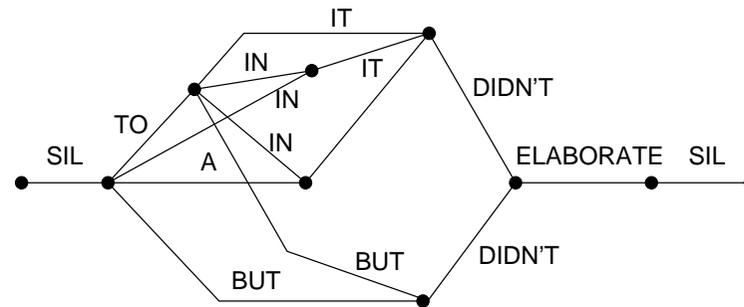


Lattice-Based MLLR

- For unsupervised adaptation hypothesis will be error-full
- Rather than using the 1-best transcription and iterative/confidence MLLR
 - generate a lattice when recognising the adaptation data [12]
 - accumulate statistics over the lattice ([Lattice-MLLR](#))



1-best transcription



Word lattice

- The accumulation of statistics is closely related to obtaining denominator statistics for discriminative training
- No need to re-recognise the data
 - iterate over the transform estimation using the same lattice

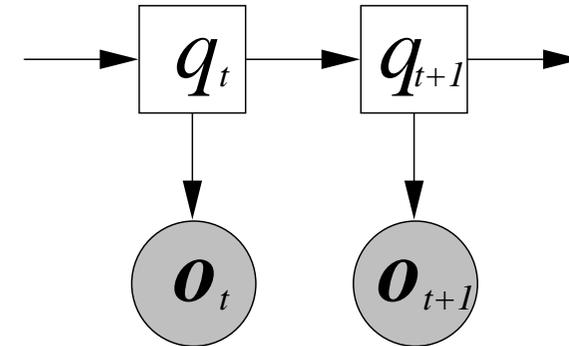
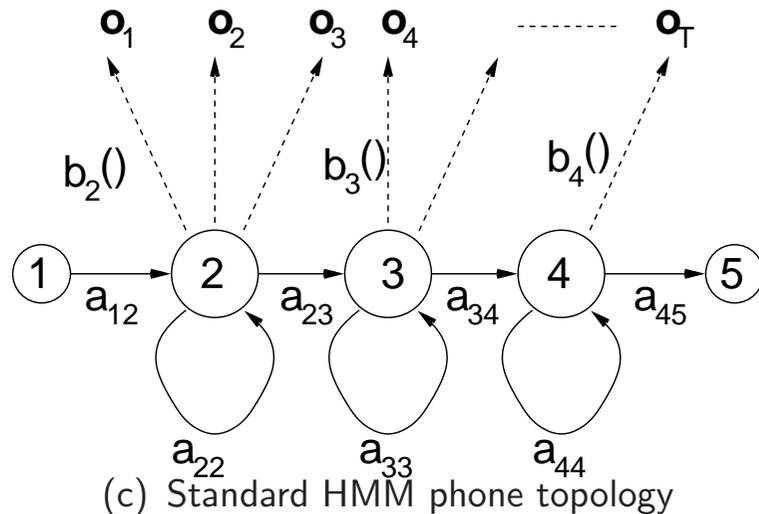


Extensions to Linear Transform Approaches

- Bayesian Adaptive Training and Inference:
 - HMMs as a dynamic Bayesian network
 - transform parameters embedded in acoustic model
 - integrated (instantaneous) adaptation and recognition
- Discriminative Mapping Transforms:
 - efficient and robust approach to obtaining discriminative linear transforms
- Noisy Constrained MLLR:
 - ML-estimated transform suitable for both noise and speaker adaptation
 - integration into adaptive training framework



Hidden Markov Model - A Dynamic Bayesian Network



- Notation for DBNs:

circles - continuous variables

shaded - observed variables

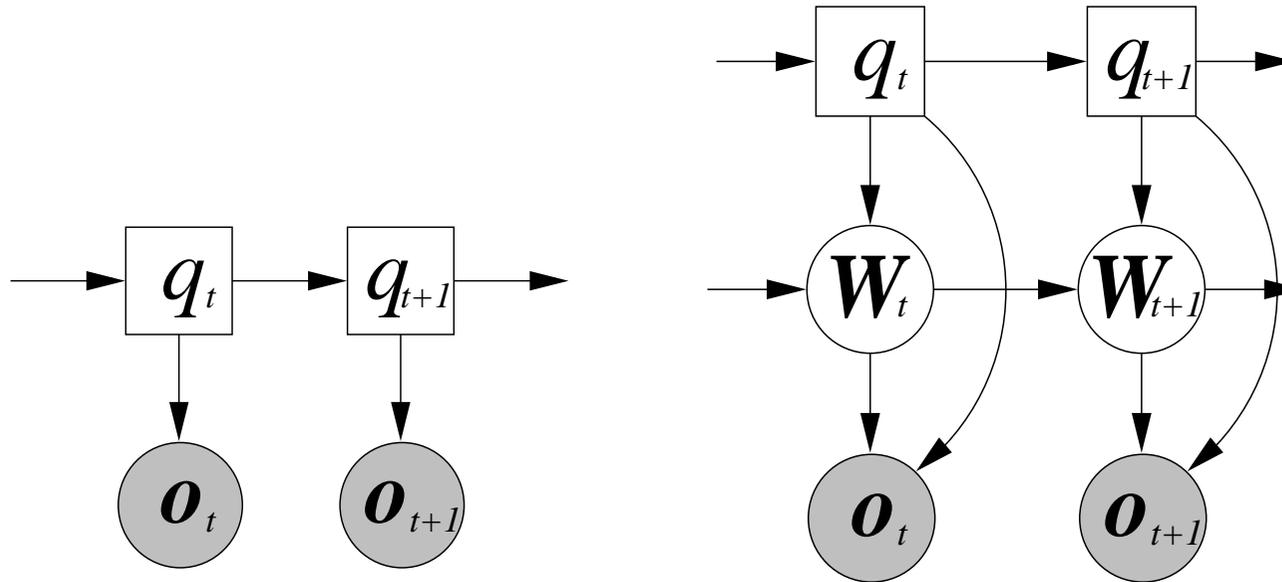
squares - discrete variables

non-shaded - unobserved variables

- Observations conditionally independent of other observations given state.
- States conditionally independent of other states given previous states.
- **Poor model of the speech process - piecewise constant state-space.**



Adaptive Training From Bayesian Perspective



(e) Standard HMM

(f) Adaptive HMM

- Observation additionally dependent on transform \mathbf{W}_t [13]
 - transform same for each homogeneous block ($\mathbf{W}_t = \mathbf{W}_{t+1}$)
 - adaptation integrated into acoustic model - **instantaneous adaptation**
- Need to know the prior transform distribution $p(\mathbf{W})$ (as in MAP scheme)



Inference with Adaptive HMMs

- Acoustic score - marginal likelihood of the whole sequence, $\mathbf{O} = \mathbf{o}_1, \dots, \mathbf{o}_T$
 - still depends on the hypothesis \mathcal{H}
 - point-estimate canonical parameters (standard complexity control schemes)

$$\begin{aligned}
 p(\mathbf{O}|\mathcal{H}) &= \int_{\mathbf{W}} p(\mathbf{O}|\mathcal{H}, \mathbf{W})p(\mathbf{W}) d\mathbf{W} \\
 &= \int_{\mathbf{W}} \sum_{\mathbf{q} \in \mathcal{Q}(\mathcal{H})} P(\mathbf{q}) \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \mathbf{A}\boldsymbol{\mu}^{(q_t)} + \mathbf{b}, \boldsymbol{\Sigma}^{(q_t)})p(\mathbf{W}) d\mathbf{W}
 \end{aligned}$$

- Latent variables makes exact inference impractical
 - need to sum over all possible state-sequences explicitly
 - Viterbi decoding not possible to find best hypothesis
- Need schemes to handle both these problems [13]
 - variational Bayes/MAP N-best supervision/adaptation/rescoring



Utterance Level Bayesian Adaptation

- Initial evaluations on English Conversational Speech recognition task

Bayesian Approx	ML Train	
	SI	SAT
—	32.8	—
ML	35.5	35.2
MAP	32.2	31.8
VB	31.8	31.5

- All experiments use **N-best** supervision
 - ML adaptation much worse than SI - insufficient adaptation data
 - VB yields additional gains over MAP
 - Note: N-best supervision better than 1-best (0.5% for VB-SAT)
- SAT performance better than SI performance
 - gains from adaptive HMM 1.3% absolute over SI baseline
 - integrated adaptation seems to be useful (though implementation an issue)



Discriminative Linear Transforms

- Linear transforms can be trained using **discriminative criteria** [14, 15]
 - estimation using minimum phone error (MPE) training

$$\mathbf{W}_d^{(s)} = \arg \min_{\mathbf{W}} \left\{ \sum_{\mathcal{H}} P(\mathcal{H} | \mathbf{O}^{(s)}; \mathbf{W}) \mathcal{L}(\mathcal{H}, \mathcal{H}^{(s)}) \right\}.$$

- For **unsupervised adaptation** discriminative linear transforms (DLTs) not used
 - estimation highly sensitive to errors in supervision hypothesis
 - more costly to estimate transform than ML training
- Not used for discriminative SAT [16], standard procedure
 1. perform standard ML-training (ML-SI)
 2. perform ML SAT training updating models and transforms (ML-SAT)
 3. estimate **MPE-models** given the **ML-transforms** (MPE-SAT)



Discriminative Mapping Functions

- Would like to get aspects of discriminative transform without the problems:
 - train all speaker-specific parameters using ML training
 - train speaker-independent parameters using MPE training
- Applying this to linear transforms yields (as one option) [17]

$$\begin{aligned}\boldsymbol{\mu}^{(s)} &= \mathbf{A}_d \left(\mathbf{A}_{m1}^{(s)} \boldsymbol{\mu} + \mathbf{b}_{m1}^{(s)} \right) + \mathbf{b}_d \\ &= \mathbf{A}_d \boldsymbol{\mu}_{m1}^{(s)} + \mathbf{b}_d\end{aligned}$$

- $\mathbf{W}_{m1}^{(s)} = \begin{bmatrix} \mathbf{A}_{m1}^{(s)} & \mathbf{b}_{m1}^{(s)} \end{bmatrix}$ - speaker-specific ML transform
- $\mathbf{W}_d = \begin{bmatrix} \mathbf{A}_d & \mathbf{b}_d \end{bmatrix}$ - speaker-independent MPE transform

- Yields a composite **discriminative-like** transform

$$\mathbf{A}_d^{(s)} = \mathbf{A}_d \mathbf{A}_{m1}^{(s)}; \quad \mathbf{b}_d^{(s)} = \mathbf{A}_d \mathbf{b}_{m1}^{(s)} + \mathbf{b}_d$$



Training DMTs

- This form of DMT results in the following estimation criterion

$$\mathbf{W}_d = \arg \min_{\mathbf{W}} \left\{ \sum_s \sum_{\mathcal{H}} P(\mathcal{H} | \mathbf{O}^{(s)}; \mathbf{W}, \mathbf{W}_{ml}^{(s)}) \mathcal{L}(\mathcal{H}, \mathcal{H}^{(s)}) \right\}.$$

- posterior $P(\mathcal{H} | \mathbf{O}^{(s)}; \mathbf{W}, \mathbf{W}_{ml}^{(s)})$ based on speaker ML-adapted models
- supervised training of discriminative transform
- Standard DLT update formulae can be used
- Quantity of training data vast compared to available speaker-specific data
 - use large number of base-classes
 - in these experiments 1000 base-classes used
- Can also be used for adaptive training [18]
 - closer to full discriminative adaptive training



Discriminative Adaptive Training with DMTs

- Initial evaluations on English Conversational Speech recognition task

Training Scheme	Transform		WER (%) eval03
	Training	Testing	
SI	—	—	29.2
SI	—	MLLR	27.0
		MLLR+DMT	26.2
DSAT	MLLR	MLLR	26.4
	MLLR	MLLR+DMT	25.6
	DLT	DLT	28.1
	MLLR+DMT	MLLR+DMT	25.3

- All systems trained using MPE (both multi-style and adaptive)
- As expected adaptation helps with the multi-style trained system
 - DMTs help with the multi-style trained system (0.8% absolute)
 - DMTs help with adaptively trained system (1.1% absolute)



Noisy CMLLR

- Linear transforms described are general
 - hierarchies allow very complex forms to be used
 - interesting to examine forms aimed at particular tasks
- Noisy CMLLR is aimed at noise-robust speech [19] recognition

$$p(\mathbf{o}_t; \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}^{(m)}, \mathbf{A}, \mathbf{b}, \boldsymbol{\Sigma}_b) = |\mathbf{A}| \mathcal{N}(\mathbf{A}\mathbf{o}_t + \mathbf{b}; \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}^{(m)} + \boldsymbol{\Sigma}_b)$$

- has the same form as a model-based compensation scheme (JUD)
- Similar to CMLLR, but with an additional bias on the variance
 - CMLLR can be viewed as estimating the “neutral” speech
 - the variance bias, $\boldsymbol{\Sigma}_b$, a level of uncertainty
- Form can be used in an adaptive training/discriminative fashion as well



Noisy CMLLR and Factor Analysis

- The estimation of/adaptive training of NCMLLR related to:
 - shared factor analysis approach for covariance matrix modelling [20]
 - EM-based VTS adaptive training for canonical model estimation [21]
- All treat “clean” speech as a latent variable
 - posterior distribution depends on the form being examined
 - update for canonical models:

$$\hat{\boldsymbol{\mu}}^{(m)} = \frac{\sum_{h=1}^H \sum_{t=1}^T \gamma_t^{(mh)} \mathcal{E} \{ \mathbf{s}_t | \mathbf{o}_t, m \}}{\sum_{h=1}^H \sum_{t=1}^T \gamma_t^{(mh)}}$$

- Discriminative adaptive training also considered [19]

$$\hat{\boldsymbol{\mu}}^{(m)} = \frac{\sum_{h=1}^H \sum_{t=1}^T (\gamma_{\text{num}t}^{(mh)} - \gamma_{\text{den}t}^{(mh)}) \mathcal{E} \{ \mathbf{s}_t | \mathbf{o}_t, m \} + D_m \boldsymbol{\mu}^{(m)} + \tau_p \boldsymbol{\mu}_p^{(m)}}{\sum_{h=1}^H \sum_{t=1}^T (\gamma_{\text{num}t}^{(mh)} - \gamma_{\text{den},t}^{(mh)}) + D_m + \tau_p}$$



Noisy CMLLR vs CMLLR Performance

- Evaluated on engine-on/highway noise condition from the Toshiba data
 - phone-numbers task (unknown digit length sequences)
 - see next section for test data configuration, here run at speaker level
 - simplified training data set-up trained on noise-corrupted WSJ SI-284

System	Adapt (diag)	ENON		HWY	
		ML	MPE	ML	MPE
Multi-style	—	1.2	0.8	6.7	5.0
	CMLLR	0.3	0.3	2.4	2.0
	NCMLLR	0.5	0.6	2.1	1.9
Adaptive Training	CMLLR	0.3	0.2	2.1	1.5
	NCMLLR	0.3	0.2	1.8	1.2

- Adaptive training again shows gains over multi-style training
 - NCMLLR out-performs CMLLR at low SNR conditions
 - MPE gains larger when using adaptive training



Speaker Adaptation Summary

- Speaker adaptation an important part of speech recognition systems
- Linear transform-based adaptation still dominant form for LVCSR adaptation
 - extensively used in CU-HTK and other evaluation systems
 - needs to be able to handle errors in the hypotheses
 - need to be able to discriminatively estimated transforms
- Adaptive training a theoretically very interesting extension
 - use adaptation transforms during training
 - allows a “neutral” speaker model to be generated
- Gains for speaker adaptive training still disappointing ...
- Though simple (just a linear transform) still issues to be addressed
 - e.g. integrating adaptation into acoustic model efficiently ...



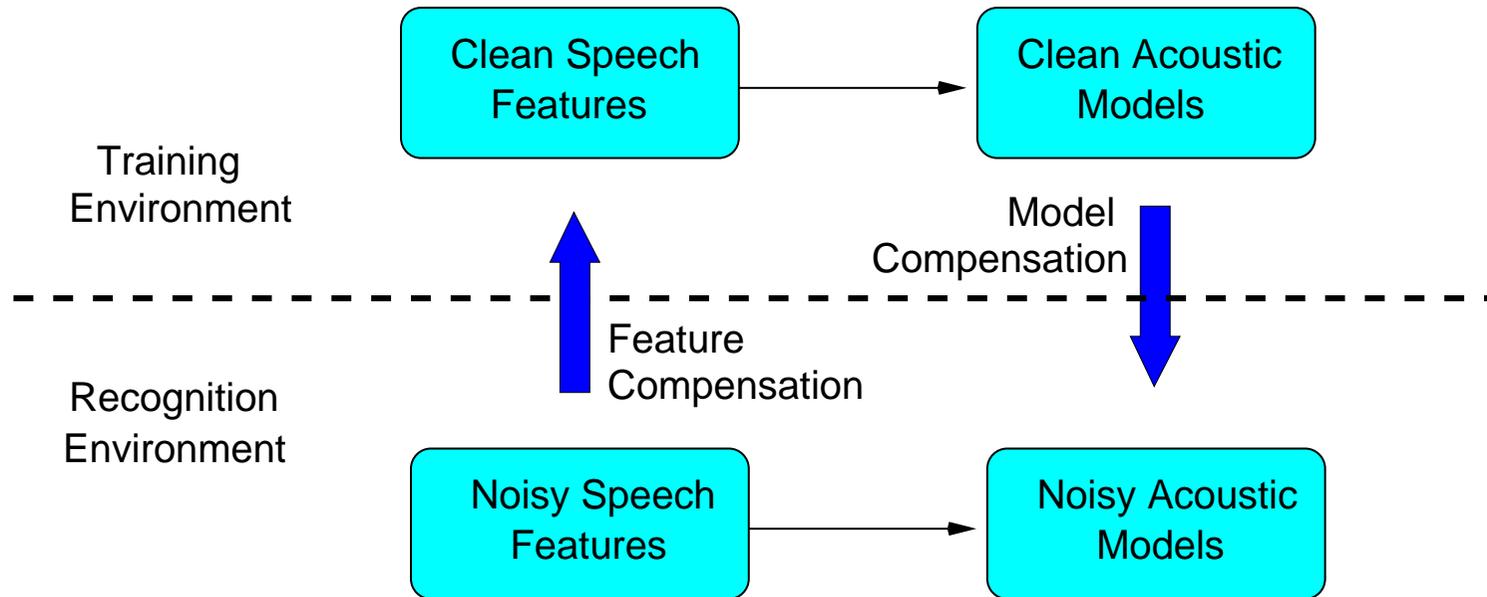
Noise-Robust ASR



Noise Robust ASR - In-Car Navigation



Noise Compensation Approaches



- Two main approaches:
 - **feature** compensation: “clean” the noisy features
 - **model** compensation: “corrupt” the clean models
- This work concentrates on **model compensation** approaches
 - VTS and JUD examples, **predictive** model compensation schemes



Mismatch Functions

- Speech data is normally parameterised in the Cepstral domain, thus

$$\mathbf{y}_t^s = \frac{1}{2} \mathbf{C} \log \left(\exp(2\mathbf{C}^{-1}\mathbf{x}_t^s + 2\mathbf{C}^{-1}\mathbf{h}^s) + \exp(2\mathbf{C}^{-1}\mathbf{n}_t^s) \right) = \mathbf{x}_t^s + \mathbf{h}^s + f(\mathbf{x}_t^s, \mathbf{n}_t^s, \mathbf{h}^s)$$

\mathbf{C} is the DCT, **magnitude**-based Cepstra

- non-linear relationship between the clean speech, noise and corrupted speech
 - not possible to get simple expression for all parameterisations
- This has assumed sufficient smoothing to remove all “cross” terms
 - some sites use **interaction likelihoods** or **phase-sensitive** functions [22, 23]
 - given $\mathbf{x}_t^s, \mathbf{h}^s$ and \mathbf{n}_t^s there is a distribution

$$\mathbf{y}_t^s \sim \mathcal{N}(\mathbf{x}_t^s + \mathbf{h}_t^s + f(\mathbf{x}_t^s, \mathbf{n}_t^s, \mathbf{h}^s), \Phi)$$



Mismatch function optimisation

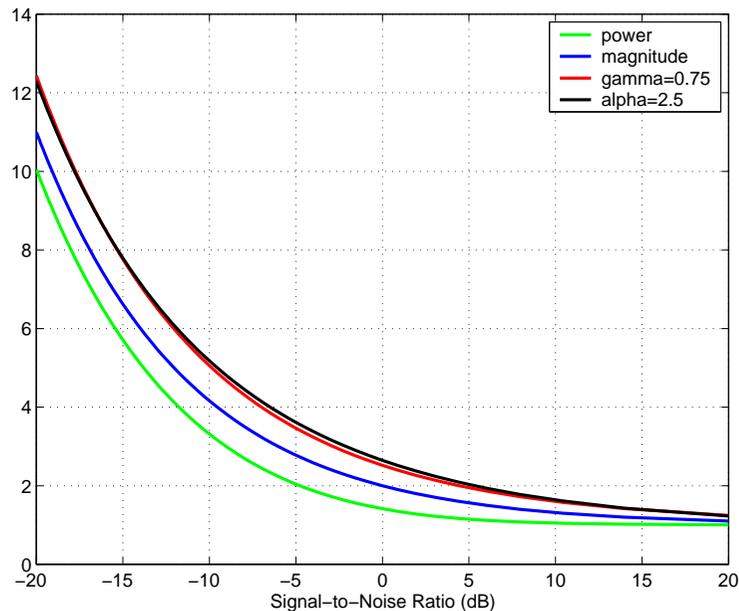
- The mismatch function is only an approximation - variations possible

- γ -optimised - tunable parameter γ , ignoring \mathbf{h}^s

$$\mathbf{y}_t^s = \mathbf{x}_t^s + \frac{1}{\gamma} \mathbf{C} \log \left(1 + \exp \left(\gamma \mathbf{C}^{-1} (\mathbf{n}_t^s - \mathbf{x}_t^s) \right) \right)$$

- Phase-sensitive - tunable parameter α , in theory $-1 \leq \alpha \leq 1$

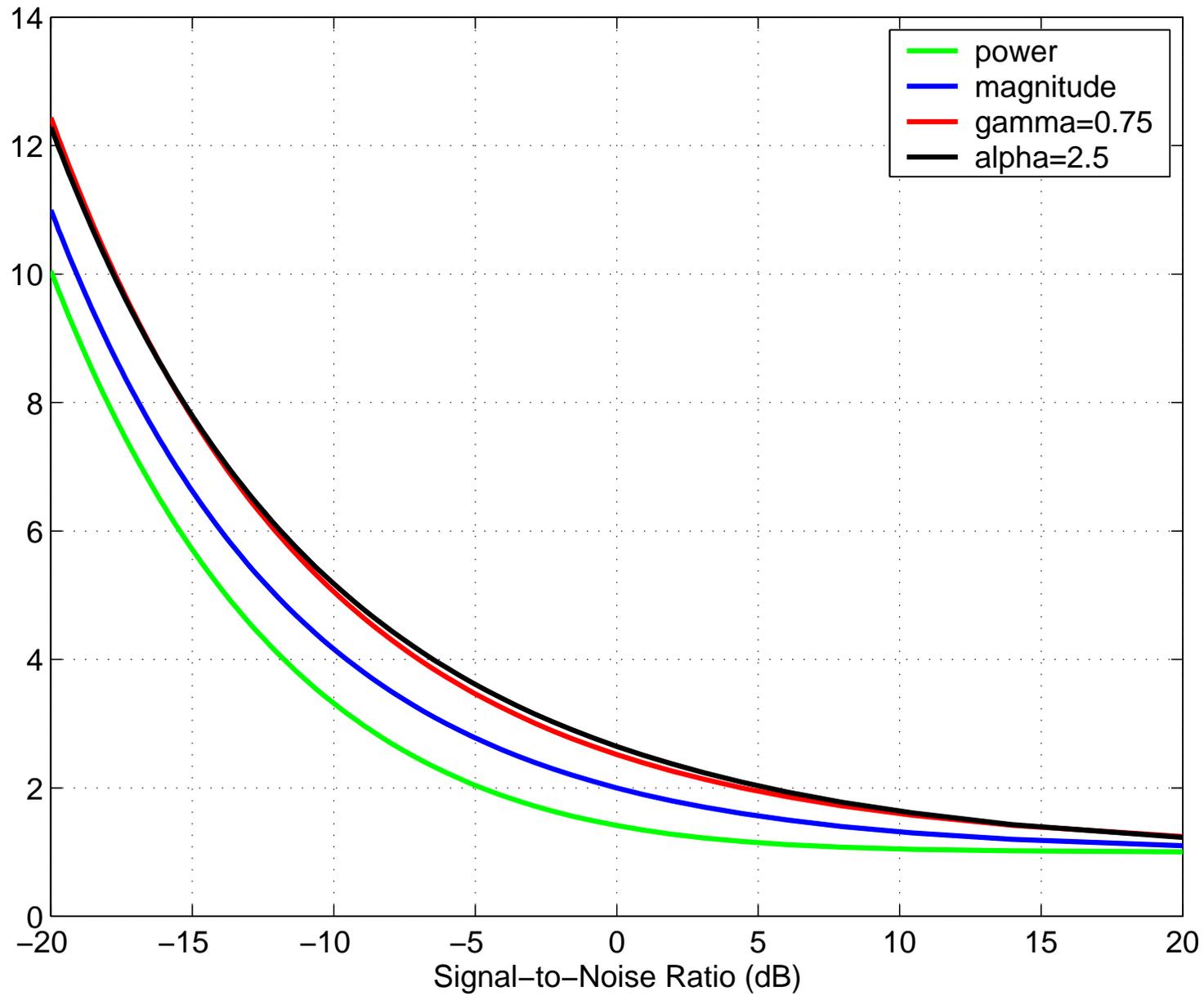
$$\mathbf{y}_t^s = \mathbf{x}_t^s + \frac{1}{2} \mathbf{C} \log \left(1 + \exp \left(2 \mathbf{C}^{-1} (\mathbf{n}_t^s - \mathbf{x}_t^s) \right) + 2\alpha \exp \left(\mathbf{C}^{-1} (\mathbf{n}_t^s - \mathbf{x}_t^s) \right) \right)$$



Ratio of corrupted speech magnitude to clean speech magnitude

- magnitude ($\alpha = 1, \gamma = 1$)
- power ($\alpha = 0, \gamma = 2$)
- $\alpha = 2.5$ (AURORA tuned [24])
- $\gamma = 0.75$ (AURORA tuned [25])
- $\gamma = 1.0$ used in this work





Delta and Delta-Delta Parameters

- Aim to ‘reduce’ HMM conditional independence assumptions
 - standard to add **delta** and **delta-delta** [26] parameters

$$\mathbf{y}_t = \begin{bmatrix} \mathbf{y}_t^s \\ \Delta \mathbf{y}_t^s \\ \Delta^2 \mathbf{y}_t^s \end{bmatrix}; \quad \Delta \mathbf{y}_t^s = \frac{\sum_{i=1}^n w_i (\mathbf{y}_{t+i}^s - \mathbf{y}_{t-i}^s)}{2 \sum_{i=1}^n w_i^2}$$

- Two versions used to represent the impact of noise on these [27, 28]

$$\Delta \mathbf{y}_t^s \approx \frac{\partial \mathbf{y}_t^s}{\partial t} \quad \text{OR} \quad \Delta \mathbf{y}_t^s = \mathbf{D} \begin{bmatrix} \mathbf{y}_{t-1}^s \\ \mathbf{y}_t^s \\ \mathbf{y}_{t+1}^s \end{bmatrix}$$

- the second is more accurate, but more statistics required to be stored
- need to compensate all model parameters for best performance
- For enhancement can simply base deltas on static “clean” features



Model-Based Compensation

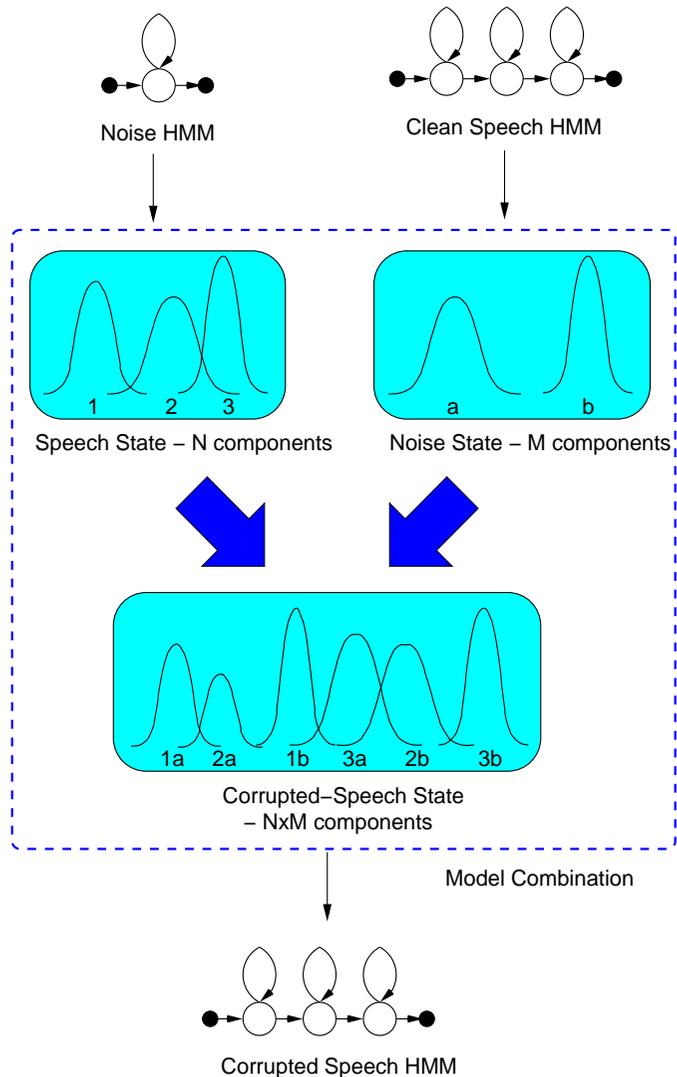
- Could retrain system using noise-corrupted training data
 - need to have all training data available and corrupt it with noise
 - slow - **single-pass retraining** [29] a faster approximation
- Model-based compensation approximates SPR [29]

$$\boldsymbol{\mu}_y^{(m)} = \mathcal{E}\{\mathbf{y}|m\}; \quad \boldsymbol{\Sigma}_y^{(m)} = \text{diag}\left(\mathcal{E}\{\mathbf{y}\mathbf{y}^T|m\} - \boldsymbol{\mu}_y^{(m)}\boldsymbol{\mu}_y^{(m)T}\right)$$

- Due to non-linearities no closed form solution - approximations required [29]
 - **Monte-Carlo**-style: generate “speech” and “noise” observations and combine
 - **Log-Add**: only transform the mean
 - **Log-Normal**: sum of two log-normal variables approximately log-normal
 - **Vector Taylor series**: first or higher order expansions used [30]
- Referred to here as **predictive schemes** - model parameters implicitly found
 - contrast to **adaptive speaker** transforms - explicit parameter estimation



Model-Based Compensation Procedure



- Each speech/noise pair considered
 - yields final component
- Also multiple-states possible
 - 3-D Viterbi decoding [31]
- Iterative schemes also possible:
 - iterative PMC [29]
 - Algonquin [22]
- Commonly used configuration:
 - single state
 - single component

Vector Taylor Series

- **Vector Taylor Series (VTS)** one popular approximation [32, 30]
 - Taylor series expansion about “current” parameter values
 - for these expression ignore impact of convolutional distortion
 - mismatch function approximated using first order series

$$\mathbf{y}_t^s \approx \boldsymbol{\mu}_x^s + f(\boldsymbol{\mu}_x^s, \boldsymbol{\mu}_n^s) + \nabla_x f(\mathbf{x}, \mathbf{n})|_{\boldsymbol{\mu}_x^s, \boldsymbol{\mu}_n^s} (\mathbf{x}_t^s - \boldsymbol{\mu}_x^s) + \nabla_n f(\mathbf{x}, \mathbf{n})|_{\boldsymbol{\mu}_x^s, \boldsymbol{\mu}_n^s} (\mathbf{n}_t^s - \boldsymbol{\mu}_n^s)$$

where $f(\mathbf{x}, \mathbf{n})$ is the mismatch function from previous slide (ignoring \mathbf{h}^s)

- Gives simple approach to estimating noise parameters

$$\boldsymbol{\mu}_y^{(m)s} = \mathcal{E}\{\mathbf{y}_t^s | m\} \approx \boldsymbol{\mu}_x^{(m)s} + f(\boldsymbol{\mu}_x^{(m)s}, \boldsymbol{\mu}_n^s)$$

$$\boldsymbol{\Sigma}_y^{(m)s} \approx \mathbf{A} \boldsymbol{\Sigma}_x^{(m)s} \mathbf{A}^\top + (\mathbf{I} - \mathbf{A}) \boldsymbol{\Sigma}_n^{(m)s} (\mathbf{I} - \mathbf{A})^\top; \quad \mathbf{A} = \frac{\partial \mathbf{y}^s}{\partial \mathbf{x}^s}$$



Noise Parameter Estimation

- In practice the noise model parameters, μ_n, μ_h, Σ_n , are not known
 - need to be estimated from test data
 - simplest approach - use VAD and start/end frames to estimate noise
- Also possible to use ML estimation [32, 33, 24]

$$\left\{ \hat{\mu}_n, \hat{\mu}_h, \hat{\Sigma}_n \right\} = \operatorname{argmax}_{\mu_n, \mu_h, \Sigma_n} \left\{ p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mu_n, \mu_h, \Sigma_n; \lambda_x) \right\}$$

- VTS approximation yields simple approach to find μ_n, μ_h
 - first/second-order approaches to find Σ_n
 - simple statistics for auxiliary function
- Parameters estimated in the same fashion as unsupervised adaptation
 - need to have hypothesis \mathcal{H}



Extensions to Model-Based Approaches

- Joint Uncertainty Decoding:
 - derived from joint clean/corrupted speech modelling
 - attempts to speed up model compensation process
- Predictive Linear Transforms:
 - efficiently handles changes in the feature-vector correlations
- Adaptive Training using VTS/JUD:
 - training systems with a wide-range of back-ground noise conditions
- Incremental Adaptation:
 - using model-based related schemes in a causal fashion



Minimum Mean-Square Error Estimates

- Estimate the clean speech $\hat{\mathbf{x}}_t$ given the corrupted speech \mathbf{y}_t
 - to handle non-linearity partition space using an R -component GMM, then

$$\hat{\mathbf{x}}_t = \mathcal{E}\{\mathbf{x}_t|\mathbf{y}_t\} = \sum_{r=1}^R P(r|\mathbf{y}_t) \mathcal{E}\{\mathbf{x}_t|\mathbf{y}_t, r\}$$

- Model the joint-distribution for each component, then [34]

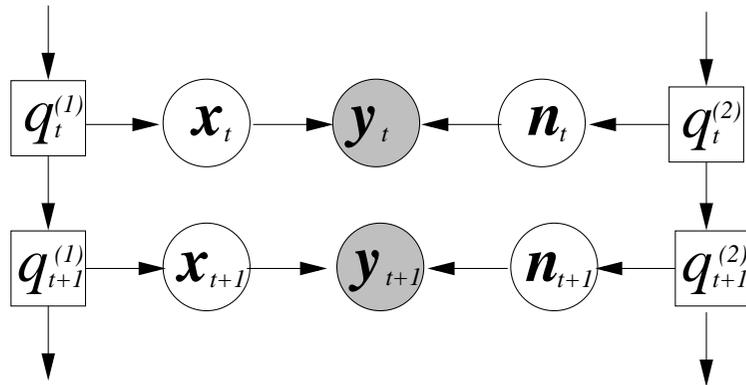
$$\begin{bmatrix} \mathbf{y}_t \\ \mathbf{x}_t \end{bmatrix} \Big| r \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_y^{(r)} \\ \boldsymbol{\mu}_x^{(r)} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{yy}^{(r)} & \boldsymbol{\Sigma}_{yx}^{(r)} \\ \boldsymbol{\Sigma}_{xy}^{(r)} & \boldsymbol{\Sigma}_{xx}^{(r)} \end{bmatrix} \right)$$

$$\mathcal{E}\{\mathbf{x}_t|\mathbf{y}_t, r\} = \boldsymbol{\mu}_x^{(r)} + \boldsymbol{\Sigma}_{xy}^{(r)} \boldsymbol{\Sigma}_{yy}^{(r)-1} (\mathbf{y}_t - \boldsymbol{\mu}_y^{(r)}) = \mathbf{A}^{(r)} \mathbf{y}_t + \mathbf{b}^{(r)}$$

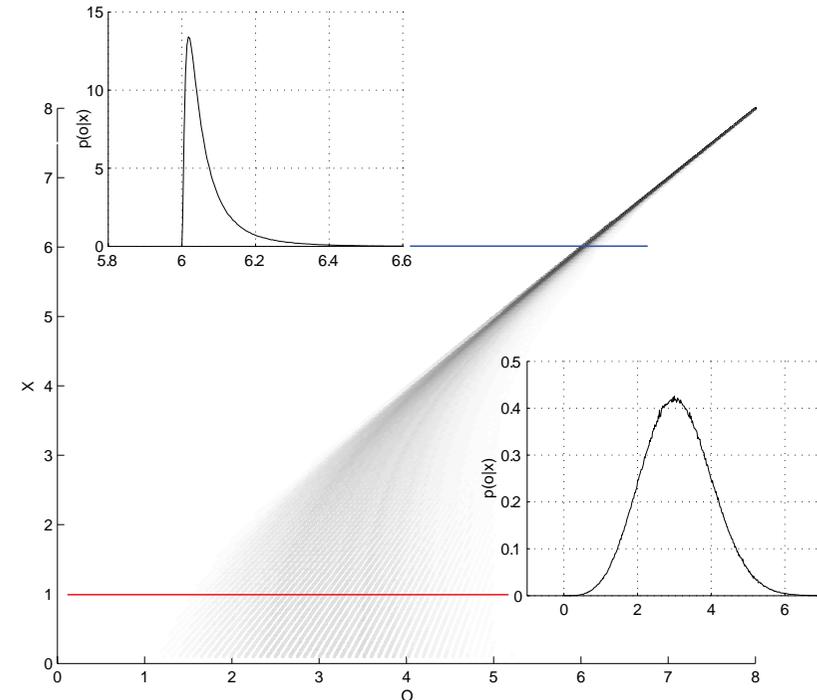
- joint distribution estimated using **stereo data**
- can be estimated using model-based compensation schemes [32, 35]
- various forms/variants possible: SPLICE [36], POF[37], VTS-based [32, 38]



Uncertainty Decoding



$$p(\mathbf{y}_t) = \int p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{n}_t) p(\mathbf{x}_t) p(\mathbf{n}_t) d\mathbf{n}_t d\mathbf{x}_t$$



- All the model-based approaches are computationally expensive
 - scales linearly with # components (100K+ for LVCSR systems)
- Need to model the conditional distribution $p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{n}_t)$ [39, 22, 33]
 - select form to allow efficient compensation/decoding (if possible)



Joint Uncertainty Decoding

- Rather than model $p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{n}_t)$ use [33]

$$p(\mathbf{y}_t|\mathbf{x}_t) = \int p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{n}_t)p(\mathbf{n}_t)d\mathbf{n}_t$$

- Simplest approach is to assume \mathbf{y}_t and \mathbf{x}_t jointly Gaussian (again)
 - to handle changes with acoustic-space make dependent on r
 - simple to derive conditional distribution $p(\mathbf{y}_t|\mathbf{x}_t, r)$
 - contrast to MMSE where $p(\mathbf{x}_t|\mathbf{y}_t, r)$ modelled
 - joint distribution estimated using VTS/PMC (stereo data can also be used)
- Product of Gaussians is an un-normalised Gaussian, so

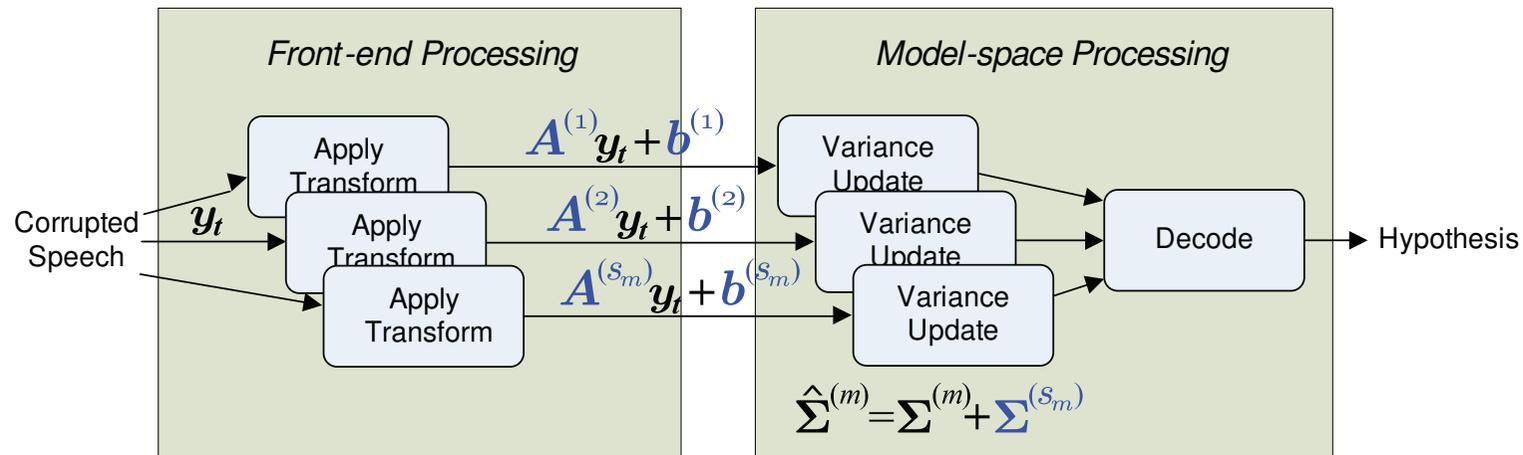
$$p(\mathbf{y}_t|m, r) = |\mathbf{A}^{(r)}| \mathcal{N}(\mathbf{A}^{(r)}\mathbf{y}_t + \mathbf{b}^{(r)}; \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}^{(m)} + \boldsymbol{\Sigma}_b^{(r)})$$

- r is normally determined by the component m [40]
- contrast to MMSE where GMM built in acoustic space to determine r



JUD versus (N)CMLLR

- For JUD compensation, PMC/VTS only required at regression class level
 - $\mathbf{A}^{(r)}$, $\mathbf{b}^{(r)}$ and $\Sigma_{\mathbf{b}}^{(r)}$ functions of noise parameters $\mu_{\mathbf{n}}$, $\mu_{\mathbf{h}}$, $\Sigma_{\mathbf{n}}$



- Similar to CMLLR however
 - JUD parameters estimated using noise models derived from data
 - CMLLR directly uses data to estimate parameters
 - JUD has a bias variance, found to be important for noise estimation
 - same form as NCMLLR, but estimated in a predictive fashion



Predictive Linear Transforms

- Consider a GMM, the corrupted/adapted distributions are

$$p(\mathbf{y}) = \sum_{m=1}^M c_y^{(m)} \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_y^{(m)}, \boldsymbol{\Sigma}_y^{(m)}); \quad \tilde{p}(\mathbf{y}) = \sum_{m=1}^M c_x^{(m)} |\mathbf{A}| \mathcal{N}(\mathbf{A}\mathbf{y} + \mathbf{b}; \boldsymbol{\mu}_x^{(m)}, \boldsymbol{\Sigma}_x^{(m)})$$

- how to estimate the “best” linear transform?
- Estimate should be based on minimising the KL-divergence

$$\mathcal{KL}(p||\tilde{p}) = \int p(\mathbf{y}) \log \left(\frac{p(\mathbf{y})}{\tilde{p}(\mathbf{y})} \right) d\mathbf{y}$$

- using the matched-bound approximation (K terms independent of \mathbf{A} , \mathbf{b})

$$\mathcal{KL}(p||\tilde{p}) \leq - \sum_{m=1}^M c_y^{(m)} \int \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_y^{(m)}, \boldsymbol{\Sigma}_y^{(m)}) \log \left(\mathcal{N}(\mathbf{A}\mathbf{y} + \mathbf{b}; \boldsymbol{\mu}_x^{(m)}, \boldsymbol{\Sigma}_x^{(m)}) \right) d\mathbf{y} + K$$

- a framework for estimating “predictive” linear transforms [41]



Predictive CMLLR

- For schemes like CMLLR required the “predictive” statistics are e.g.:

$$\mathbf{k}_{pci}^{(r)} = \sum_{m \in \mathbf{r}_r} \frac{\gamma^{(m)} \mu_{xi}^{(m)}}{\sigma_{xi}^{(m)2}} \begin{bmatrix} 1 \\ \mathcal{E}\{\mathbf{y}|m\} \end{bmatrix}$$

- normally the expectations obtained from observed data (**adaptive**)
- could also use model-compensation schemes to obtain values (**predictive**)
- $\gamma^{(m)}$ either based on observations $\gamma_{yt}^{(m)}$ or training data counts $\gamma_x^{(m)}$

$$\mathcal{E}\{\mathbf{y}|m\} = \frac{\sum_t \gamma_{yt}^{(m)} \mathbf{y}_t}{\sum_t \gamma_{yt}^{(m)}} \quad \text{or} \quad \mathcal{E}\{\mathbf{y}|m\} = \boldsymbol{\mu}_y^{(m)}$$

- Schemes such as JUD can be used to efficiently obtain “pseudo” statistics

$$\sum_{m \in \mathbf{r}_r} \frac{\gamma_x^{(m)}}{\sigma_{xi}^{(m)2}} \mathcal{E}\{\mathbf{y}|m\} = \mathbf{A}^{(r)-1} \left(\sum_{m \in \mathbf{r}_r} \frac{\gamma_x^{(m)} \boldsymbol{\mu}_x^{(m)}}{\sigma_{xi}^{(m)2}} \right) - \mathbf{A}^{(r)-1} \mathbf{b}^{(r)} \left(\sum_{m \in \mathbf{r}_r} \frac{\gamma_x^{(m)}}{\sigma_{xi}^{(m)2}} \right)$$



“Adaptive” vs “Predictive” Schemes

- Adaptive and predictive schemes complementary to one another

Adaptive	Predictive
general approach	applicable to noise
linear assumption	mismatch function required
- use many linear transforms	- may be inaccurate
transform parameters estimated	noise model estimated
- large numbers of parameters	- small number of parameters

- Obvious approach is to combine the two in a fashion similar to MAP [42]:
 - limited data predictive approaches used
 - increased data adaptive approaches used
- Count smoothing simple approach to use (parent transforms also possible)

$$\mathbf{k}_{pai}^{(r)} = \frac{\mathbf{k}_{pci}^{(r)}}{\sum_{m \in r_r} \gamma_x^{(m)}} + \tau_{sm} \mathbf{k}_i^{(r)}$$



PCMLLR vs MMSE Schemes

- Both MMSE and PCMLLR yield linear transforms of the feature-space

$$\hat{\mathbf{x}}_t = \mathbf{A}\mathbf{y}_t + \mathbf{b}$$

- both make use of the joint distribution between clean and corrupted speech
- **However** motivation for the two approaches very different
 - MMSE is the expected value of the **clean** speech
 - PCMLLR is the linear transform that minimises the **KL-divergence**
- **Theoretically** should use:
 - MMSE: when enhancing data for additional processing
 - PCMLLR: when transformed data directly for recognition
- Initial AURORA results show that PCMLLR out-performs MMSE

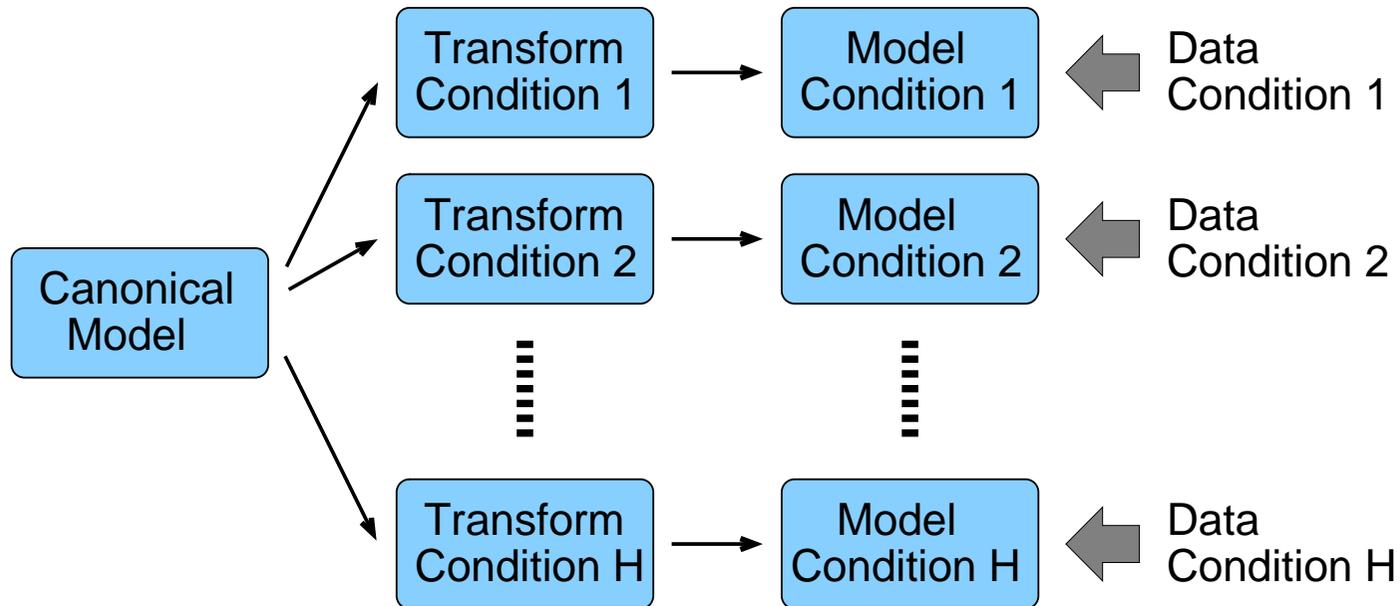


Adaptive Training

- In practice training data comes from multiple sources
 - various levels and sources of background noise
 - various speakers and channel conditions
- Multistyle (multi-environment) models required to represent **all** variabilities
 - for wide-range of noises models become very “broad”
 - previously seen issues with applying VTS/JUD to multi-style models
- **Adaptive training** one approach to handling this
 - adaptive training with various transforms previously investigated
 - generic transforms: MLLR, CMLLR, CAT
 - noise targeted transform: Noisy CMLLR
- **Perform adaptive training with VTS and JUD**
 - Joint Adaptive Training examined on Broadcast News transcription [43]
 - interested in applying VTS-adaptive training/JAT in lower SNR conditions



VTS/JOINT based Adaptive Training



- Same general framework as other adaptive training schemes
 - partition data into H homogeneous subsets
 - interleave updates of transform and canonical model
- Canonical model update (more) interesting with VTS/JUD transforms



VTS/JOINT Adaptive Training

- System trained by interleaving transform (VTS/JUD) and HMM estimates
 - VTS/JUD estimates usual approach using current canonical model
 - canonical model estimation based on second-order optimisation [43]
 - also possible to use EM-based [21]
- Derivative wrt to mean of canonical model given by

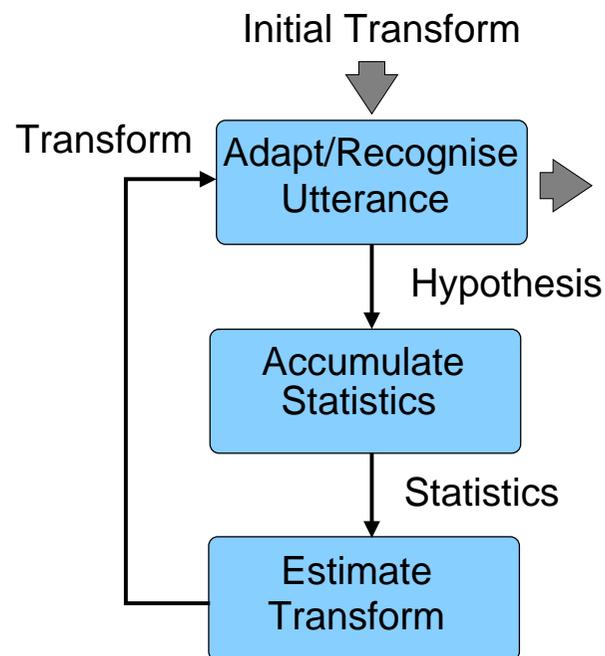
$$\frac{\partial Q_j}{\partial \mu_{xi}^{(m)}} = \sum_{h=1}^H \sum_{t=1}^{T^{(h)}} \gamma_{yt}^{(m)} \left(\frac{\hat{\mathbf{a}}_i^{(rh)} \mathbf{y}_t + \hat{b}_i^{(rh)} - \mu_{xi}^{(m)}}{\sigma_{xi}^{(m)2} + \hat{\sigma}_{bi}^{(rh)2}} \right)$$

- $\hat{\sigma}_{bi}^{(rh)2}$ is larger in low-SNR regions
 - impact of observation on derivative **decreases** for low-SNR
 - **agrees with intuition**
- VAT implementation is JAT with #regression classes = #system components



Incremental Noise Estimation

- Batch-mode adaptation introduces latency into decoding process
 - for some tasks, e.g. in-car command/control, need to minimise latency
 - many tasks require multiple interactions over a short period
- Incremental adaptation introduces no latency-



- generate hypothesis using current transform
- accumulate statistics $\mathcal{O}_i^{(m)}$ using hypothesis

$$\mathcal{O}_i^{(m)} = \sum_t \gamma_{yt}^{(m)} \mathbf{y}_t^{(i)} + \alpha \mathcal{O}_{i-1}^{(m)}$$

- 0 \Rightarrow no smoothing
- 1 \Rightarrow “complete” smoothing
- estimate transform for next utterance



Combined Incremental/Adaptive Processing

- Adaptive training requires a test-condition transform for good performance
 - normally requires a multi-pass system
 - multiple models may be required (where to get initial hypothesis)
 - sensitivity to initial hypothesis/transform
- Noise parameters can be estimated using a single utterance
- **Incremental adaptation** a good framework for VTS/JUD adaptive systems
 - no need for multiple-passes
 - output can be generated in a causal fashion
 - hypothesis/initial transform may be “good” (depending on form of noise)
 - only an adaptively trained system needed
- **BUT** still need an initial transform for first utterance to get things started



Toshiba In-Car Task

- TREL-CRL04 small/medium sized recognition task
 - Speech collected in the office and in vehicles (enon, city, highway)
 - phone numbers (PH) task used for initial evaluations
 - * 30 English speakers (15 male, 15 female) uttering 30 sentences each
 - * 35, 18 SNR averages for the enon, highway condition, respectively
 - 4 digits (4D), command & control (CC) and city names (CN) also used

#	PH	4D	CC	CN
utt	861	757	1916	958
words/utt	9.5	4.0	5.5	1.2
secs/utt	6.9	3.6	4.7	3.2
vocabulary	11	11	119	544

averaged between enon/hway (no city)

- Range of lengths (total not speech) and vocabularies



System Configuration

- Acoustic training data - 486 hours of data
 - mixture of real in-car data and clean data artificially corrupted
 - approx 283 hours artificially corrupted data (from WSJ data)
 - approx 203 hours “real” data
- Acoustic model characteristics
 - MFCC-parameters plus delta/delta-delta features (39-dimensional)
 - ≈ 650 states, 12 components/state, ≈ 7800 components
 - acoustic models: decision tree clustered states, cross-word triphones
- compact system (embedded market is possible target domain)
- Both multi-style (multi) and adaptively trained (adapt) system built



Batch Multi-Style vs VTS Adaptive Training (phone-numbers)

Iteration	ENON		HWY	
	multi	adapt	multi	adapt
0	1.1		4.5	
1	1.5	0.5	3.2	1.6
2	1.4	0.5	2.4	1.5

- VTS adaptation applied to both multi-style and adaptively trained systems
 - update hypothesis and transform at each iteration
- Multi-style performance degraded by VTS for high SNR conditions
 - mismatch function not suitable for multi-style training
- VTS Adaptive training consistent gains
 - better than multi-style training, worked at higher SNR conditions



Adaptively-Trained VTS Performance Summary

		ENON				HWY				Avg
		PH	4D	CC	CN	PH	4D	CC	CN	
iter	0	1.1	1.0	0.9	3.9	4.5	3.0	2.0	14.5	3.86
	1	0.5	0.3	0.7	3.6	1.6	1.5	1.4	13.6	2.90
	2	0.5	0.2	0.7	3.7	1.5	1.3	1.3	11.0	2.53
ETSI-adv		1.2	1.1	0.9	4.5	3.1	1.8	1.2	8.2	2.75

- Batch VTS adaptation evaluated on full range of tasks
 - compared to ETSI-advanced front-end with same training data
- For ENON consistent gains for all conditions
- For HWY mixed results are more mixed
 - City-Names (CN) very poor performance
 - related to sensitivity to initial hypothesis/noise estimate
 - can get similar performance to ETSI advanced front-end eventually



Incremental vs Batch with Adaptively Trained System

System	ENON				HWY				Avg
	PH	4D	CC	CN	PH	4D	CC	CN	
VAT-batch	0.5	0.2	0.7	3.7	1.5	1.3	1.3	11.0	2.53
VAT-INC	0.6	0.3	0.8	3.8	2.0	1.9	1.5	6.8	2.21
JAT-INC	1.2	0.7	0.8	3.6	2.5	2.3	1.7	6.9	2.46
ETSI-adv	1.2	1.1	0.9	4.5	3.1	1.8	1.2	8.2	2.75

- Incremental adaptation applied to adaptively trained systems
 - smoothing factor of $\alpha = 0.6$, 2-iteration batch results
 - also used JAT - highly efficient noise estimation/adaptation
 - initial hypothesis from multi-style system ...
- Slight degradation from batch to incremental for ENON
 - issues with City-Names addressed for HWY
- Incremental adaptation yields good overall performance



Predictive and Adaptive Incremental Adaptation

System	HWY			
	PH	4D	CC	CN
VAT-INC	2.0	1.9	1.5	6.8
+CMLLR	1.7	2.0	1.1	6.7
ETSI-adv	3.1	1.8	1.2	8.2

- Predictive with Adaptive incremental adaptation
 - use VTS compensated models to get prior statistics for CMLLR
- Combination of predictive and adaptive provides gains
 - problem with command & controls task “fixed”
 - Phone numbers also improved using combined compensation approach
- Only preliminary results - need real data to test schemes



Summary Model-Based Noise Robustness

- Model-based compensation approaches
 - good theoretical motivation
 - but requires a mismatch function
 - slow compared to feature-enhancement
- Range of extensions to standard approaches
 - joint uncertainty decoding for faster compensation
 - predictive linear transforms - additional flexibility
 - adaptive training - handle multi-style data
 - incremental adaptation
- Interesting area - still problems
 - efficiency still needs to be improved
 - improved compensation (full covariance matrices)
 - improved use of adaptive training



References

- [1] P. C. Woodland, "Speaker adaptation for continuous density HMMs: A review," in *ISCA Adaptation Workshop*, 2001.
- [2] J. L. Gauvain and C.-H. Lee, "Maximum a-posteriori estimation for multivariate Gaussian mixture observations of Markov chains," *IEEE Transactions Speech and Audio Processing*, vol. 2, pp. 291–298, 1994.
- [3] R. Kuhn, P. Nguyen, J.-C. Junqua, L. Goldwasser, N. Niedzielski, S. Fincke, K. Field, and M. Contolini, "Eigenvoices for speaker adaptation," in *Proceedings ICSLP*, 1998, pp. 1771–1774.
- [4] M. J. F. Gales, "Cluster adaptive training of hidden Markov models," *IEEE Transactions Speech and Audio Processing*, vol. 8, pp. 417–428, 2000.
- [5] C. J. Leggetter and P. C. Woodland, "Maximum likelihood linear regression for speaker adaptation of continuous density HMMs," *Computer Speech and Language*, vol. 9, pp. 171–186, 1995.
- [6] M. J. F. Gales and P. C. Woodland, "Mean and variance adaptation within the MLLR framework," *Computer Speech and Languages*, vol. 10, pp. 249–264, 1996.
- [7] M. J. F. Gales, "Maximum likelihood linear transformations for HMM-based speech recognition," *Computer Speech and Language*, vol. 12, pp. 75–98, 1998.
- [8] S. J. Young, G. Evermann, M. J. F. Gales, T. Hain, X. Liu, D. Kershaw, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, and P. C. Woodland, *The HTK Book, version 3.4*. Cambridge, UK: Cambridge University Engineering Department, 2006.
- [9] W. Chou, "Maximum a-posterior linear regression with elliptical symmetric matrix variate priors," in *Proceedings Eurospeech*, 1999, pp. 1–4.
- [10] T. Anastasakos, J. McDonough, R. Schwartz, and J. Makhoul, "A compact model for speaker-adaptive training," in *Proceedings ICSLP*, 1996, pp. 1137–1140.
- [11] P. C. Woodland, D. Pye, and M. J. F. Gales, "Iterative unsupervised adaptation using maximum likelihood linear regression," in *Proc. Int. Conf. Spoken Lang. Process.*, Philadelphia, 1996, pp. 1133–1136.
- [12] L. F. Uebel and P. C. Woodland, "Speaker adaptation using lattice-based MLLR," in *Proc. ITRW on Adaptation Methods for Speech Recognition*, 2001.
- [13] K. Yu and M. Gales, "Bayesian adaptive inference and adaptive training," *IEEE Transactions Speech and Audio Processing*, vol. 15, no. 6, pp. 1932–1943, August 2007.



- [14] F. Wallhof, D. Willett, and G. Rigoll, "Frame-discriminative and Confidence-driven Adaptation for LVCSR," in *Proc ICASSP'00*, Istanbul, 2000, pp. 1835–1838.
- [15] S. Tsakalidis, V. Doumptotis, and W. Byrne, "Discriminative Linear Transforms for Feature Normalisation and Speaker Adaptation in HMM Estimation," *IEEE Trans Speech and Audio Processing*, vol. 13, no. 3, pp. 367–376, 2005.
- [16] L. Wang and P. Woodland, "Discriminative Adaptive Training using the MPE Criterion," in *Proc ASRU*, St Thomas, US Virgin Islands, 2003.
- [17] K. Yu, M. Gales, and P. Woodland, "Unsupervised discriminative adaptation using discriminative mapping transforms," in *Proc. ICASSP*, 2008.
- [18] C. Raut, K. Yu, and M. Gales, "Adaptive training using discriminative mapping transforms," in *Proc. InterSpeech*, 2008.
- [19] D. Kim and M. Gales, "Noisy CMLLR for noise-robust speech recognition," University of Cambridge, Tech. Rep. TR611, February 2009, available from <http://mi.eng.cam.ac.uk/~mjfg/>.
- [20] R. Gopinath, B. Ramabhadran, and S. Dharanipragada, "Factor analysis invariant to linear transformations of data," in *Proc. Int. Conf. Spoken Lang. Process.*, 1998, pp. 397–400.
- [21] Y. Hu and Q. Huo, "Irrelevant variability normalization based HMM training using VTS approximation of an explicit model of environmental distortions," in *Proc. InterSpeech*, 2007.
- [22] T. Kristjansson, "Speech recognition in adverse environments: a probabilistic approach," Ph.D. dissertation, University of Waterloo, Waterloo, Canada, 2002.
- [23] L. Deng, J. Droppo, and A. Acero, "Enhancement of log mel power spectra of speech using a phase sensitive model the acoustic environemnt and sequential estimation of the corrupting noise," *Proc. IEEE Transactions on Speech and Audio Processing*, 2004.
- [24] J. Li, L. Deng, Y. Gong, and A. Acero, "HMM adaptation with joint compensation of additive and convolutive distortions via vector Taylor series," in *ASRU 2007*, Kyoto, Japan, 2007.
- [25] M. Gales and F. Flego, "Discriminative classifiers with generative kernels for noise robust speech recognition," Cambridge University, Tech. Rep. CUED/F-INFENG/TR605, August 2008, available from: <http://mi.eng.cam.ac.uk/~mjfg>.
- [26] S. Furui, "Speaker independent isolated word recognition using dynamic features of speech spectrum," *IEEE Transactions ASSP*, vol. 34, pp. 52–59, 1986.
- [27] R. A. Gopinath, M. J. F. Gales, P. S. Gopalakrishnan, S. Balakrishnan-Aiyer, and M. A. Picheny, "Robust speech recognition in noise - performance of the IBM continuous speech recognizer on the ARPA noise spoke task," in *ARPA Workshop on Spoken Language System Technology*, 1995, pp. 127–130.



- [28] R. van Dalen and M. Gales, "Extended VTS for noise-robust speech recognition," in *Proc ICASSP*, 2009.
- [29] M. Gales, "Model-based Techniques for Noise Robust Speech Recognition," Ph.D. dissertation, Cambridge University, 1995.
- [30] A. Acero, L. Deng, T. Kristjansson, and J. Zhang, "HMM Adaptation using Vector Taylor Series for Noisy Speech Recognition," in *Proc ICSLP*, Beijing, China, 2000.
- [31] A. Varga and R. Moore, "Hidden Markov model decomposition of speech and noise," in *Proc ICASSP*, 1990, pp. 845–848.
- [32] P. Moreno, "Speech Recognition in Noisy Environments," Ph.D. dissertation, Carnegie Mellon University, 1996.
- [33] H. Liao, "Uncertainty Decoding For Noise Robust Speech Recognition," Ph.D. dissertation, Cambridge University, 2007.
- [34] X. Huang, A. Acero, and H.-W. Hon, *Spoken Language Processing*. Prentice Hall, 2001.
- [35] V. Stouten, H. V. Hamme, and P. Wambacq, "Accounting for the uncertainty of speech estimates in the context of model-based feature enhancement," in *Proc. ICSLP*, vol. 1, Jeju Island, Korea, Oct. 2004, pp. 105–108.
- [36] J. Droppo, L. Deng, and A. Acero, "Evaluation of the SPLICE Algorithm on the Aurora 2 Database," in *Proc. Eurospeech*, Aalborg, Denmark, 2001, pp. 217–220.
- [37] L. Neumeyer and M. Weintraub, "Probabilistic Optimum Filtering for Robust Speech Recognition," in *Proc. ICASSP*, Adelaide, 1994.
- [38] V. Stouten, "Robust speech recognition in time-varying environments," Ph.D. dissertation, Universitet Leuven, 2006.
- [39] J. Droppo, A. Acero, and L. Deng, "Uncertainty decoding with SPLICE for Noise Robust Speech Recognition," in *Proc ICASSP 02*, Orlando, Florida, 2002.
- [40] H. Liao and M. Gales, "Issues with Uncertainty Decoding for Noise Robust Speech Recognition," in *Proc. ICSLP*, Pittsburgh, PA, 2006.
- [41] M. J. F. Gales and R. C. van Dalen, "Predictive linear transforms for noise robust speech recognition," in *Proc. ASRU*, 2007, pp. 59–64.
- [42] F. Flego and M. Gales, "Incremental predictive and adaptive noise compensation," in *Proc ICASSP*, 2009.
- [43] H. Liao and M. Gales, "Adaptive training with joint uncertainty decoding for robust recognition of noisy data," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2007.

