

Joint Uncertainty Decoding for Noise Robust Speech Recognition

**Jasha Droppo, Alex Acero, and Li Deng
Microsoft Research, One Microsoft Way, Redmond,
Washington, USA**

Presented by
Howard

Outline

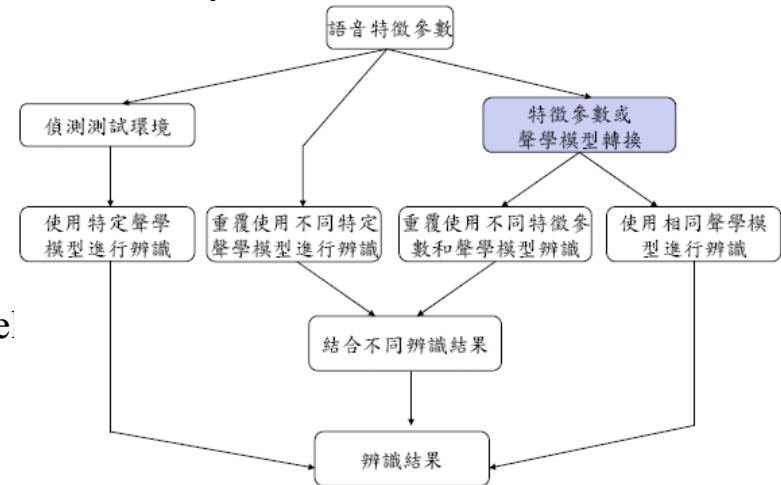
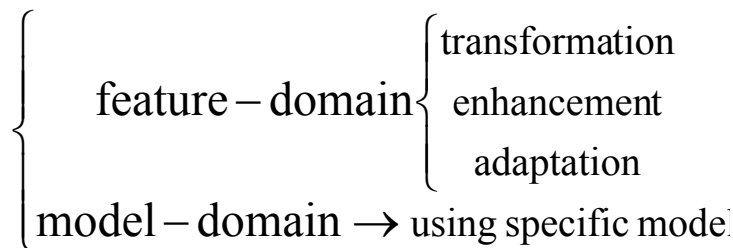
- What is robust?
- Feature-domain & model-domain
- Introduction to SPLICE
- What is uncertainty?
- Concept of uncertainty decoding
- Uncertainty Decoding with SPLICE
- Uncertainty with Joint Uncertainty Decoding

What is robust?

- We can say it is robust if it is hardly affected by extrinsic events.
 - Ex: A waterproof watch in water can still work as usual.
- For ASR
 - Speech recognition performance degrades in the presence of environmental noise, why?

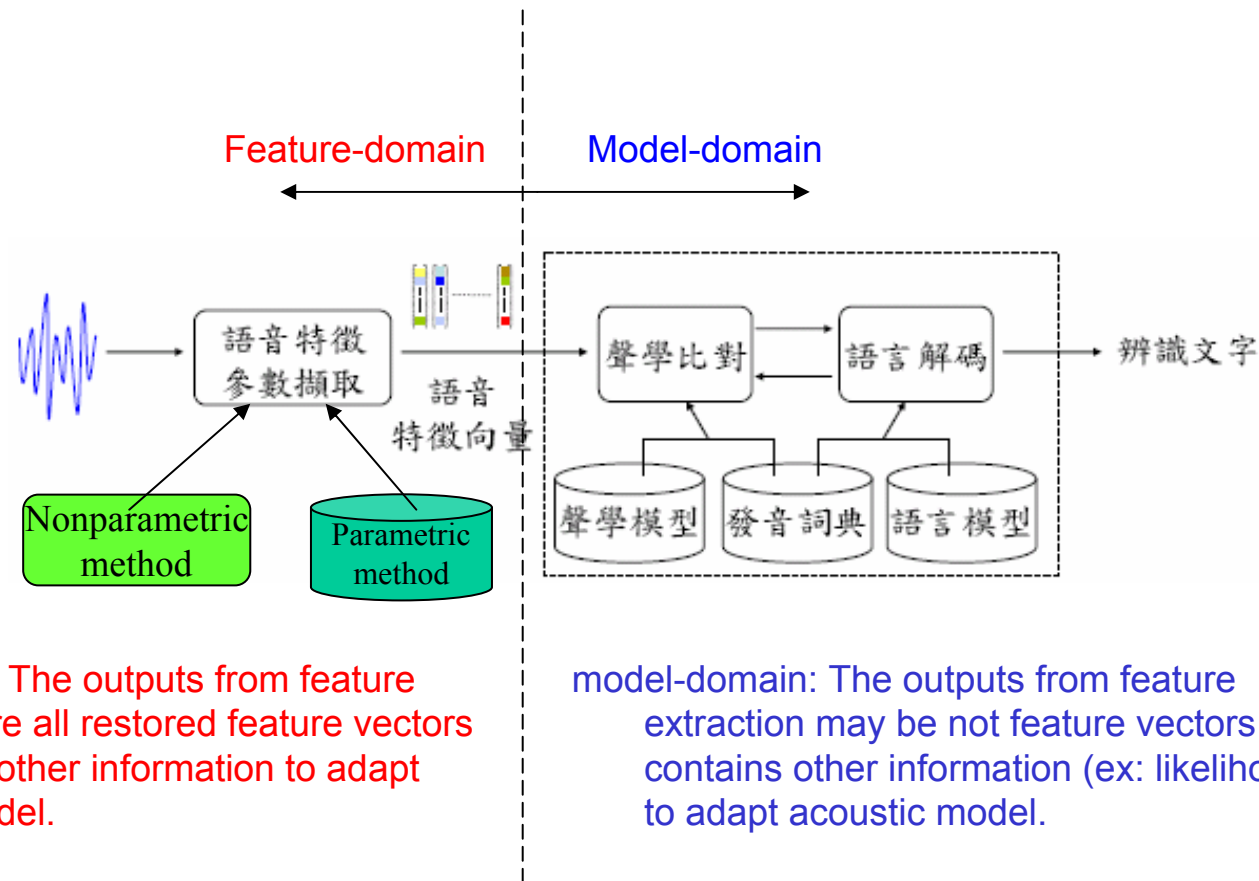
The answer is the mismatch between training and test condition.

- Solution
 - There are two main directions using different aspects to cut into this problem.



Feature-domain & Model-domain

- Where is part of feature-domain or model-domain?



Feature-domain: The outputs from feature extraction are all restored feature vectors without any other information to adapt acoustic model.

model-domain: The outputs from feature extraction may be not feature vectors but contains other information (ex: likelihood) to adapt acoustic model.

Introduction to SPLICE

- Stereo piecewise linear compensation for environment (SPLICE) takes advantage of seamlessly integrating into existing system, without a complete overhaul of existing code.

$$\tilde{y}_t = \tilde{x}_t = E[x_t | y_t] = \sum_k p(k | y_t) E[x_t | y_t, k]$$

- Assuming that the difference between clean data and corrupted data can be compensated by each single Gaussian providing a linear compensation.

$$E[x_t | y_t, k] \approx y_t + r_k$$

$$\tilde{y}_t = y_t + \sum_k p(k | y_t) \cdot r_k$$

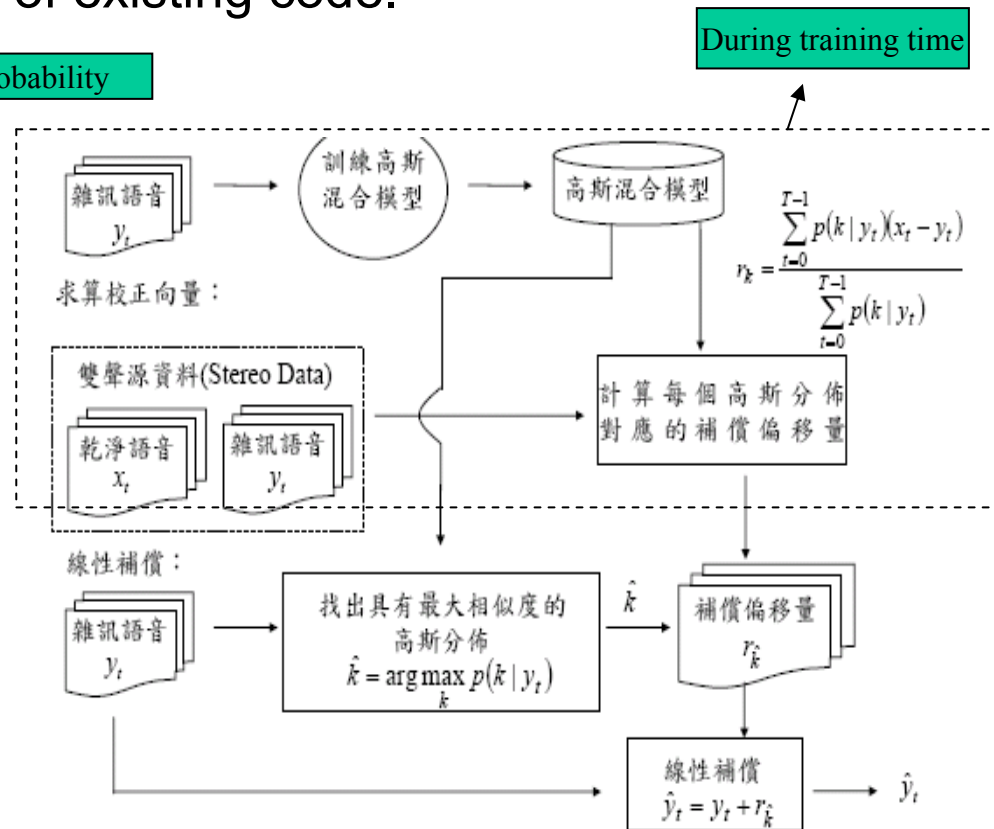
$$\hat{k} = \arg \max_k p(k | y_t)$$

$$p(k | y_t) = \frac{p(y_t | k)p(k)}{\sum_{k'=0}^{k-1} p(y_t | k')p(k')}$$

Probability of k are sum up to 1

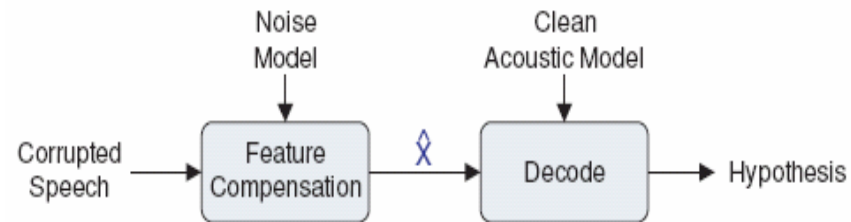
Total probability

During training time

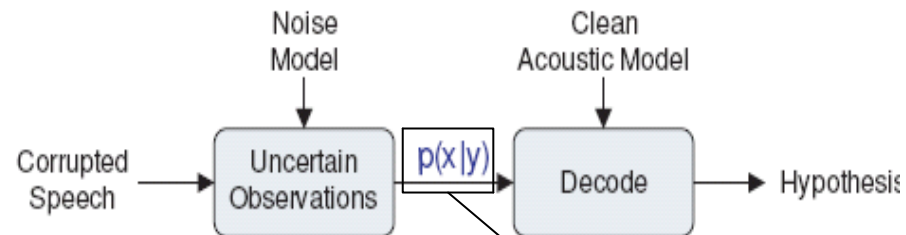


What is uncertainty?

- Feature compensation **without** uncertainty
 - The corrupted speech is restored by compensation and sent into decoder. The \hat{x} is viewed as the clean feature, is that right?

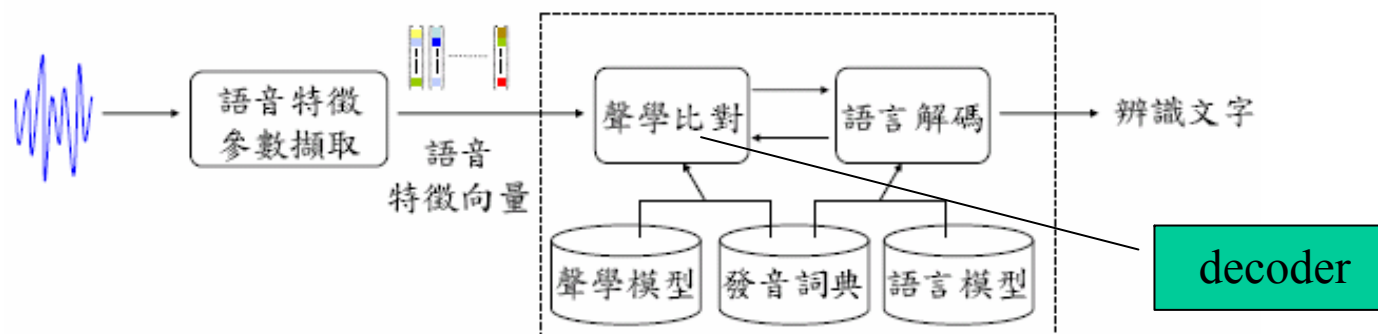


- Feature compensation **with** uncertainty
 - It is intuitively reasonable to incorporate with **uncertain observation**.

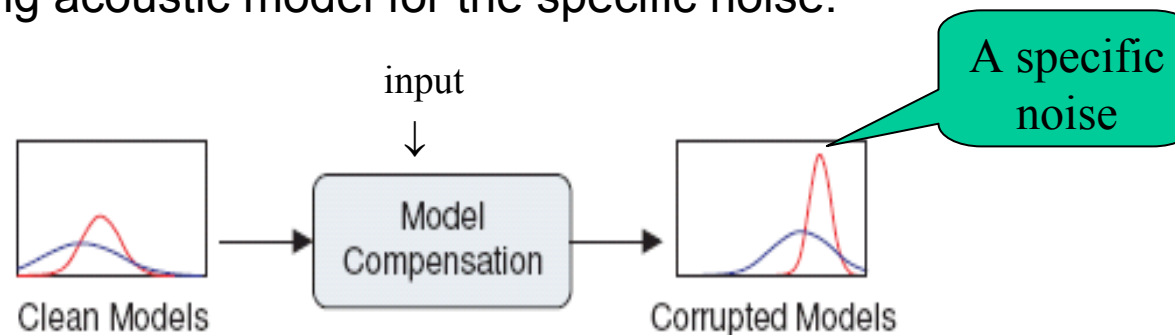


Noticing here, it is the key idea adjusted by SPLICE and JUN uncertainty decoding to make process efficient.

Concept of uncertainty decoding



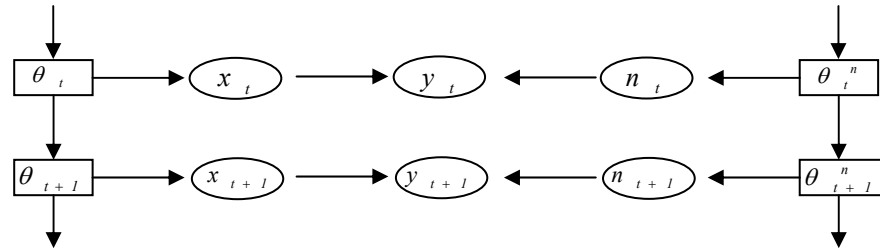
- Model-compensation
 - Renewing acoustic model for the specific noise.



- The input is either the corrupted speech data or the data combined clean and corrupted speech to achieve this goal.
- Computationally expensive

How to design uncertainty?

- Noise robustness DBN



- Corrupted speech likelihood given by

The key point

$$p(y_t | M, \check{M}, \theta_t) = \int p(y_t | x_t, \check{M}) p(x_t | M, \theta_t) dx_t \quad (1)$$

$$p(y_t | x_t, \check{M}) = \int p(y_t | x_t, n_t) p(n_t | \check{M}, \theta_t^n) dn_t \quad (2)$$

- Only $p(y_t | x_t, \check{M})$ depend on noise.
- Efficient approximation emerges from above formulation.
 - Independent of clean model complexity.
 - Appropriate form for integration.

Appendix A for (1)

marginalise

$$p(A | D) = \int p(A, B | D) d_B$$

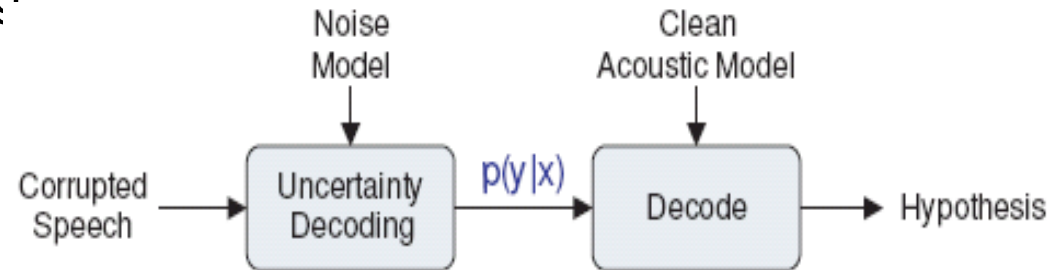
$$\begin{aligned} p(y_t | M, \check{M}, \theta_t) &= \int p(y_t, x_t | M, \check{M}, \theta_t) d_{x_t} \\ &= \int p(y_t | x_t, M, \check{M}, \theta_t) p(x_t | M, \check{M}, \theta_t) d_{x_t} \\ &= \int p(y_t | x_t, \check{M}) p(x_t | M, \theta_t) d_{x_t} \end{aligned}$$

Appendix B for (2)

$$\begin{aligned} P(y_t | x_t, \check{M}) &= \int p(y_t, n_t | x_t, \check{M}) d_{n_t} \\ &= \int p(y_t | n_t, x_t, \check{M}) p(n_t | x_t, \check{M}) d_{n_t} \\ &= \int p(y_t | n_t, x_t) p(n_t | \check{M}, \theta_t^n) d_{n_t} \end{aligned}$$

What's difference of decoding between SPLICE & JUD

- Passing conditional probability to decoding



- Passing conditional probability to decoding
- Two form of uncertainty decoding
 - Splice with uncertainty $\longrightarrow p(y_t|x_t, \tilde{M})$ by Bayes' rule
 - Joint distribution $\longrightarrow p(y_t|x_t, \tilde{M})$ by joint probability
- Both are based on Gaussian mixture model
 - Using different approximation to make process efficient

Uncertainty decoding with SPLICE

- Splice with uncertainty decoding uses Bayes' rule to write GMM as

$$p(y_t | x_t, \check{M}) = \sum_{n=1}^N \left(\frac{p(x_t | y_t, \check{s}_n, \check{M}) p(y_t | \check{s}_n, \check{M}) \check{c}_n}{p(x_t | \check{M})} \right) \quad (3)$$

- $p(x_t | y_t, \check{s}_n, \check{M})$ related to standard SPLICE estimate
- Denominator $p(x_t | \check{M})$ is a GMM – simplify using a single Gaussian

Appendix C for (3)

$$\begin{aligned} p(y_t | x_t, \check{M}) &= \sum_{n=1}^N p(y_t | \check{s}_n, x_t, \check{M}) p(\check{s}_n | x_t, \check{M}) && \leftarrow \text{Total probability} \\ &= \sum_{n=1}^N \frac{p(x_t, y_t | \check{s}_n, \check{M}) p(\check{s}_n | x_t, \check{M})}{p(x_t | \check{s}_n, \check{M})} && \leftarrow \text{Bayes' rule} \\ &= \sum_{n=1}^N \frac{p(x_t | \check{s}_n, y_t, \check{M}) p(y_t | \check{s}_n, \check{M}) p(\check{s}_n | x_t, \check{M})}{p(x_t | \check{s}_n, \check{M})} && \leftarrow \text{Using single Gaussian} \\ &= \frac{\sum_{n=1}^N p(x_t | \check{s}_n, y_t, \check{M}) p(y_t | \check{s}_n, \check{M}) p(\check{s}_n | x_t, \check{M})}{p(x_t | \check{M})} && \leftarrow \text{Replaced with prior} \\ &= \frac{\sum_{n=1}^N p(x_t | \check{s}_n, y_t, \check{M}) p(y_t | \check{s}_n, \check{M}) \check{c}_n}{p(x_t | \check{M})} \end{aligned}$$

Uncertainty with SPLICE

- Standard SPLICE uses

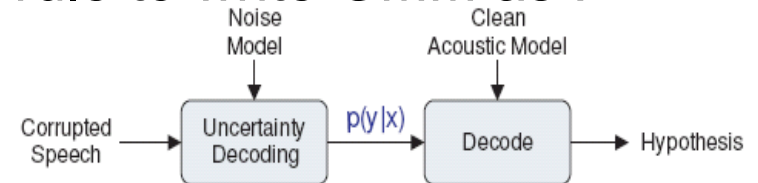
$$\hat{x}_t = E[x_t | y_t] = \sum_k P(k | y_t) E[x_t | y_t, k]$$

Replace k with s

$$= \sum_{n=1}^N P(\tilde{s}_n | y_t, \tilde{M}) \int_{x_t} P(x_t | y_t, \tilde{s}_n, \tilde{M}) d_{x_t}$$

- Uncertainty with SPLICE uses Bayes' rule to write GMM as :

$$P(y_t | x_t, \tilde{M}) = \sum_{n=1}^N \left(\frac{P(x_t | y_t, \tilde{s}_n, \tilde{M}) p(y_t | \tilde{s}_n, \tilde{M}) \tilde{c}_n}{p(x_t | \tilde{M})} \right)$$



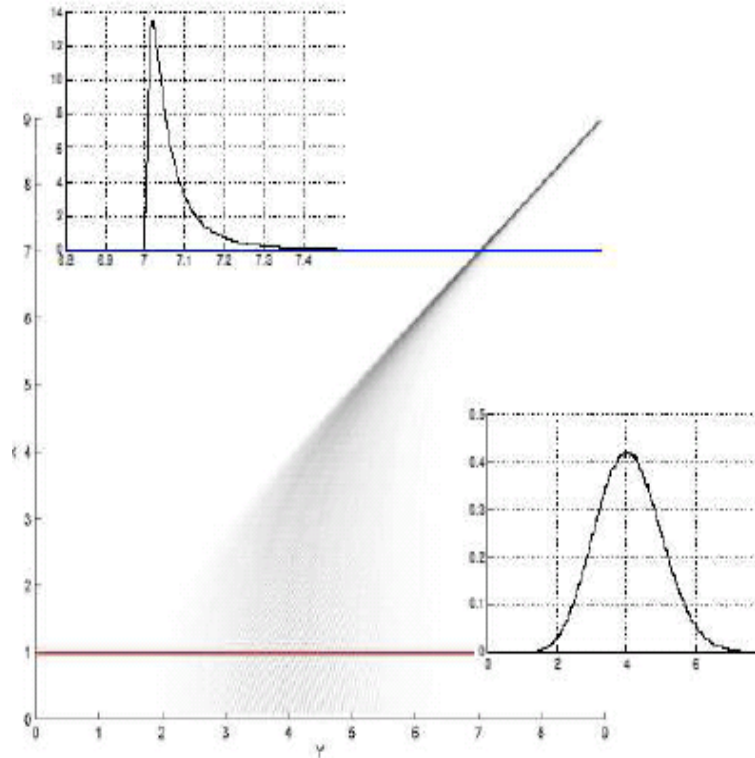
$$p(y_t | x_t, \tilde{M}, \tilde{s}_n) = f(y_t, \tilde{s}_n) N(A^{(n)} y_t + b^{(n)}; x_t, \sum_b^{(n)})$$

$$\tilde{s}_{n^*} = \arg \max_{\tilde{s}_n} \left(\frac{\tilde{c}_n p(y_t | \tilde{s}_n, \tilde{M})}{\sum_{i=1}^N \tilde{c}_i p(y_t | \tilde{s}_i, \tilde{M})} \right)$$

$$p(y_t | M, \tilde{M}, \theta_t) \propto \sum_{m \in \theta_t} c_m N(A^{(n^*)} y_t + B^{(n^*)}; \mu^{(m)}, \sum^{(m)} + \sum_b^{(n^*)})$$

Uncertainty decoding with JUD

- Joint distribution $p(x, y)$



When SNR high, the conditional is deterministic.
When SNR low, the conditional is Gaussian

Uncertainty decoding with JUN

- GMM is a standard approach to handle complex distribution
 - It's simple to marginalise tow Gaussians
- Using approximation front-end compensation model \check{M}

$$p(y_t | x_t, \check{M}) \approx \sum_{n=1}^N p(\check{s}_n | x_t, \check{M}) p(y_t | x_t, \check{s}_n, \check{M})$$

- Only \check{M} is a function of noise.
- Some issues need to be handled with
 - Component posterior $p(\check{s}_n | x_t, \check{M})$ is a function of clean speech
 - Component compensation parameters $p(y_t | x_t, \check{s}_n, \check{M})$
 - Direct use increases number of components

Uncertainty decoding for JUD

- Joint uncertainty decoding uses the GMM directly,

$$\begin{aligned} p(y_t | x_t, \tilde{M}) &= \sum_{n=1}^N p(y_t | \tilde{s}_n, x_t, \tilde{M}) p(\tilde{s}_n | x_t, \tilde{M}) \\ &= \sum_{n=1}^N \frac{p(x_t, y_t | \tilde{s}_n, \tilde{M}) p(\tilde{s}_n | x_t, \tilde{M})}{p(x_t | \tilde{s}_n, \tilde{M})} \end{aligned}$$

but

- Approximates the component posterior of clean speech, using the corrupted speech:

$$p(\tilde{s}_n | x_t, \tilde{M}) \approx p(\tilde{s}_n | y_t, \tilde{M})$$

- This decouples the front-end distribution from being dependent on the acoustic model through the clean speech variable
- conditional probability derived from the joint distribution

$$p(x_t, y_t | \tilde{s}_n, \tilde{M}) = N\left(\begin{bmatrix} x_t \\ y_t \end{bmatrix}; \begin{bmatrix} \tilde{\mu}_x^{(n)} \\ \tilde{\mu}_y^{(n)} \end{bmatrix}, \begin{bmatrix} \sum^{(n)}_{xx} & \sum^{(n)}_{xy} \\ \sum^{(n)}_{xy} & \sum^{(n)}_{yy} \end{bmatrix}\right)$$

covariance matrix is usually made diagonal for efficiency

Uncertainty decoding for JUD

- Both uncertainty decoding schemes yield same decoding form:

$$p(y_y | M, \tilde{M}, \theta_t) \approx \sum_{m=1}^M \sum_{n=1}^N \alpha^{(mn)} \mathbf{N}(\mathbf{A}^{(n)} y_t + \mathbf{b}^{(n)}; \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}^{(m)} + \boldsymbol{\Sigma}^{(n)})$$

- Form of $\mathbf{A}^{(n)}$, $\mathbf{b}^{(n)}$ and $\boldsymbol{\Sigma}^{(n)}$ differ in the two cases

- For JUN that would be

$$\mathbf{A}^{(n)} = \sum_x^{(n)} \sum_{yx}^{(n)-1}$$

$$\mathbf{b}^{(n)} = \boldsymbol{\mu}_x^{(n)} - \mathbf{A}^{(n)} \boldsymbol{\mu}_y^{(n)}$$

$$\boldsymbol{\Sigma}_b^{(n)} = \mathbf{A}^{(n)} \sum_y^{(n)} \mathbf{A}^{(n)T} - \sum_x^{(n)}$$

- To improve efficiency only a single front-end component selected, for Joint based on $p(\tilde{s}_n | y_t, \tilde{M})$
- Compared to model-based compensation computational cost is:
 - only a function of the N,
 - Not the number of components in clean speech model through variance bias must be applied