

Introduction to Speech Data

Ville Hautamäki and Tomi Kinnunen

《一剪梅》

红藕香残玉簟秋。
轻解罗裳，独上兰舟。
云中谁寄锦书来，
雁字回时，月满西楼。
花自飘零水自流。
一种相思，两处闲愁。
此情无计可消除，
才下眉头，却上心头。

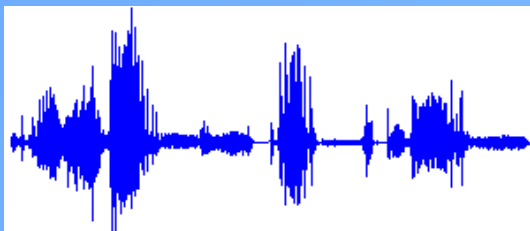


李清照

Content

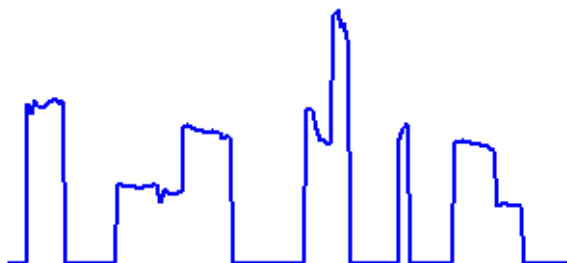
- Text-to-speech
- Speech-to-text
- 'High-level' speaker id

Speech

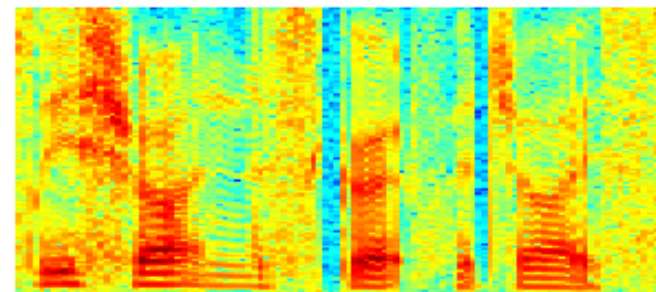


- Text-to-speech
- Voice conversion
- Automatic speaker verification

Prosody



Timbre



- Expressive TTS
- Singing synthesis
- Speaker verification

Speech production system

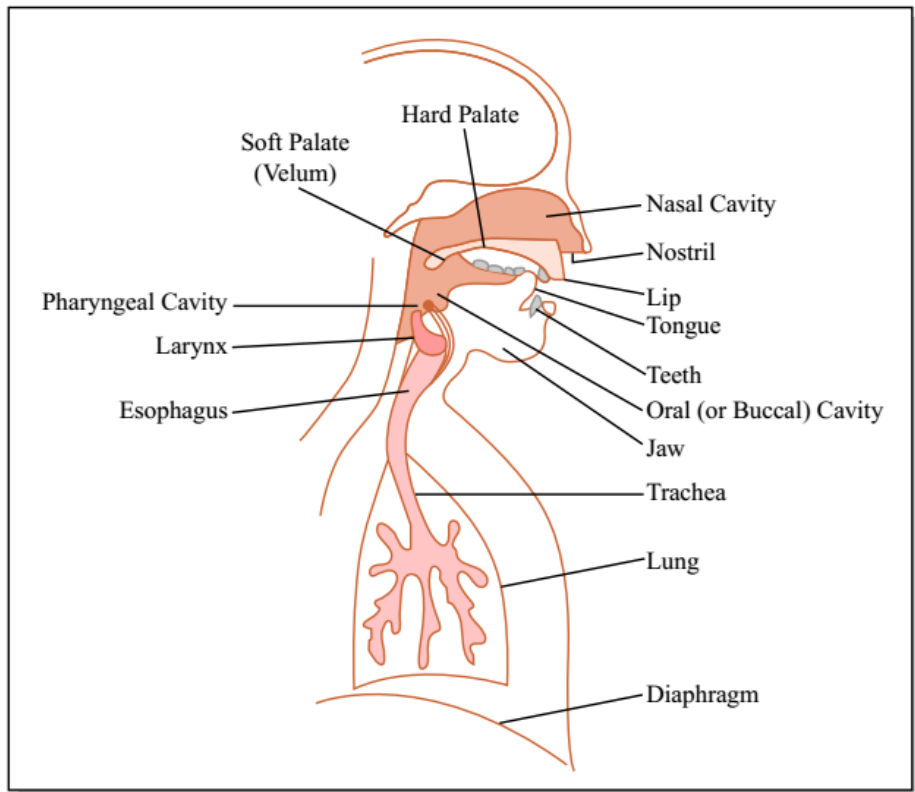


Image by MIT OpenCourseWare.

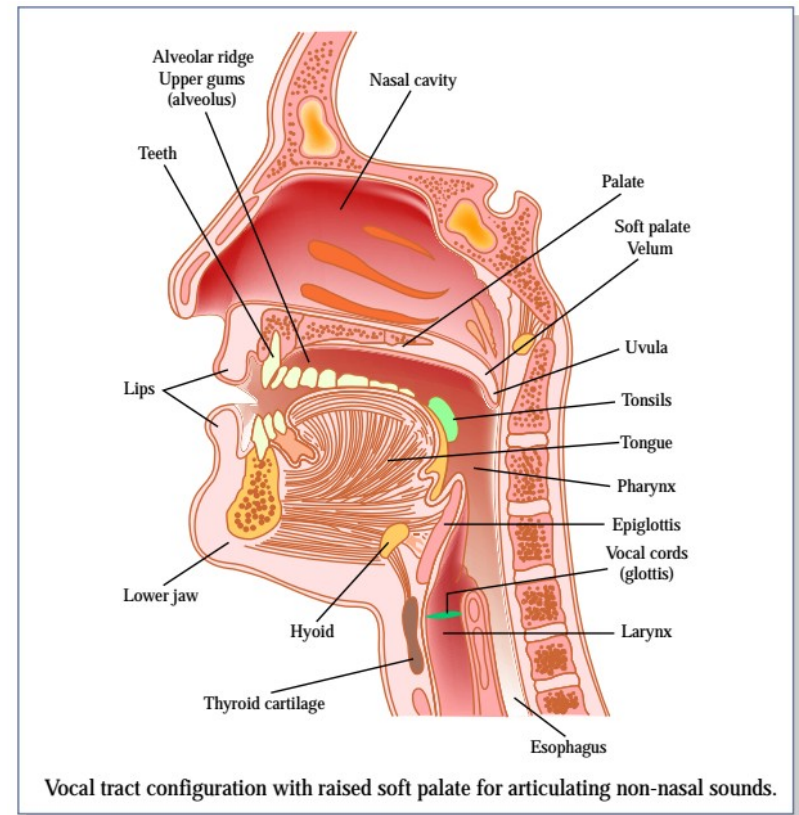


Image by MIT OpenCourseWare.

Content

Phone, word or N-gram
occupation counts

Categorical, sparse

- + Explicit/interpretable
- High error rate
- Difficult to extract

Prosody

'Stylized' F0 / energy contours,
Legendre polynomials,
multinomial subspace models

Timbre

Short-term spectrum: MFCC,
LFCC, LPCC, PNCC ... with $\Delta + \Delta^2$

Most practical (today)

- + High accuracy
- + Easy & fast to compute
- Sensitive to disturbances

Continuous, dense

Feature extraction at two levels

Level 1: short-term feature extraction

- Based on acoustic-phonetic knowledge of human speech production at millisecond time-scale
- Example: mel-frequency cepstral coefficients (MFCCs)

DIGITAL SIGNAL
PROCESSING

Level 2: recording-level feature extraction

- Extraction of higher-level features that are shared across all the frames in a recording (examples: speaker, language, recording device, age, gender, smoking habit - you name it)
- Example: identity vector (i-vector)

MACHINE
LEARNING

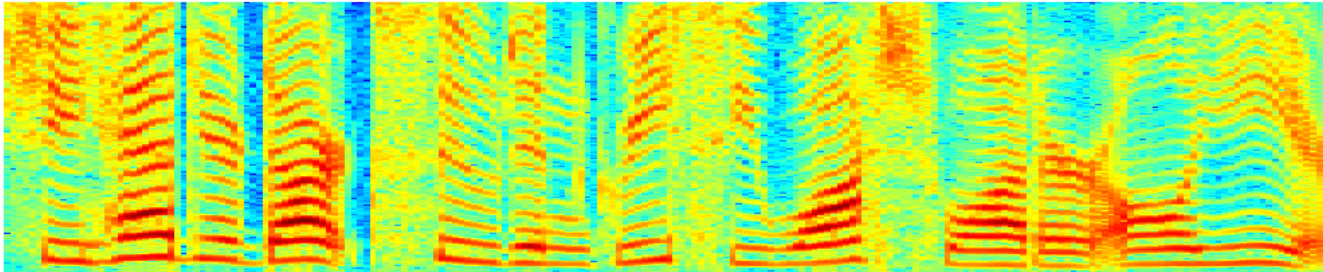
Short-term feature extraction

- Pretend that speech is a piecewise constant function over a short-time period of speech, known as a **frame**
- A feature extractor $f(\cdot)$ spits out a 'snapshot' of the short-term power spectrum \underline{x} a sequence of **feature vectors** $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$
 - Typical frame duration ~ 25 ms (1 ms = 1/1000th of a second)
 - Typical frame hop $\sim 50\%$ of the frame size
 - Typical dimension of the feature vector ~ 20 to 60
- Assume the observations to be **independent and identically distributed** (i.i.d.) samples from a recording-level statistical model, $\mathbf{x}_t \sim p(\mathbf{x} | \Theta_{\text{recording}})$
 - Usually, a Gaussian mixture model (GMM)

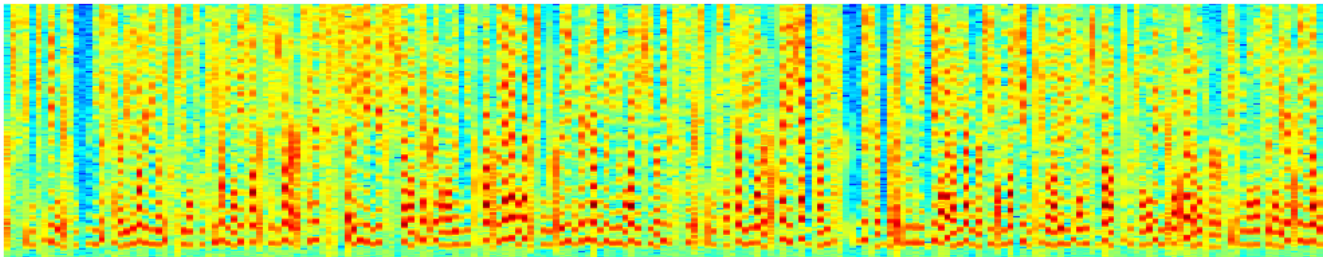
Bag-of-frames illustrated

Ordering (temporal) information gets destroyed

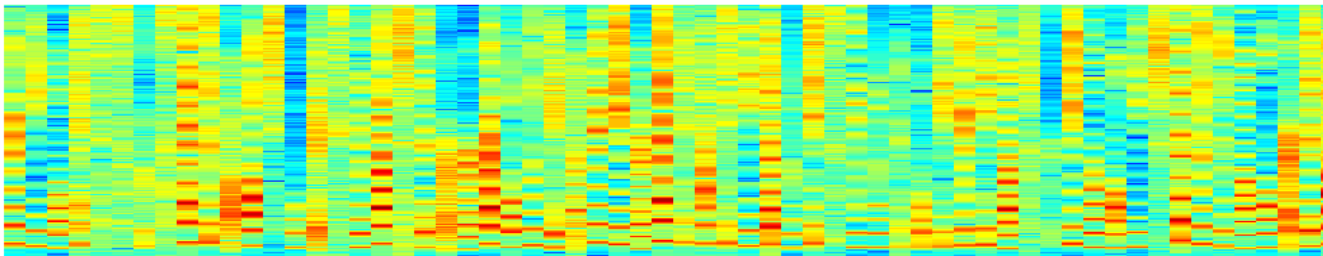
Original



Shuffled frames



Shuffled frames



Original



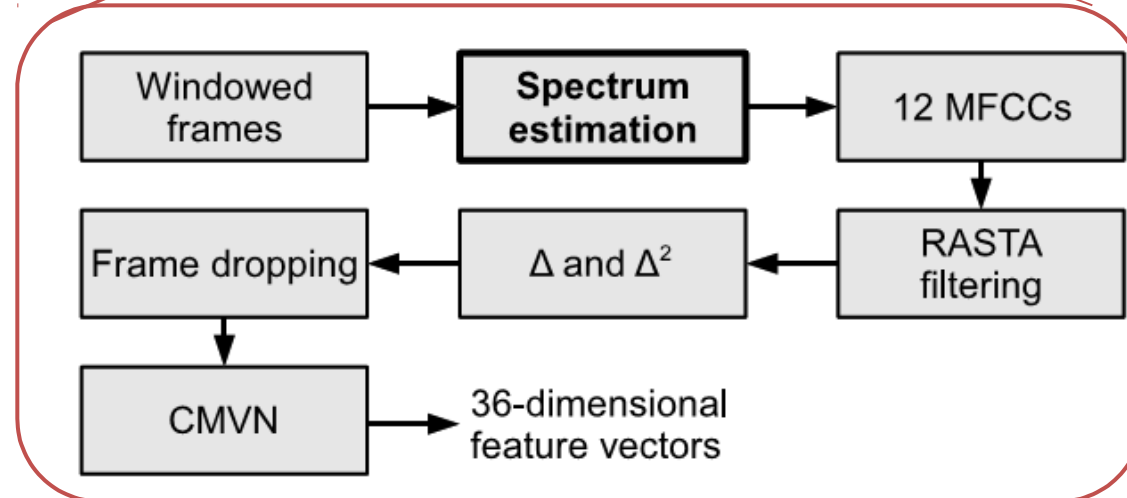
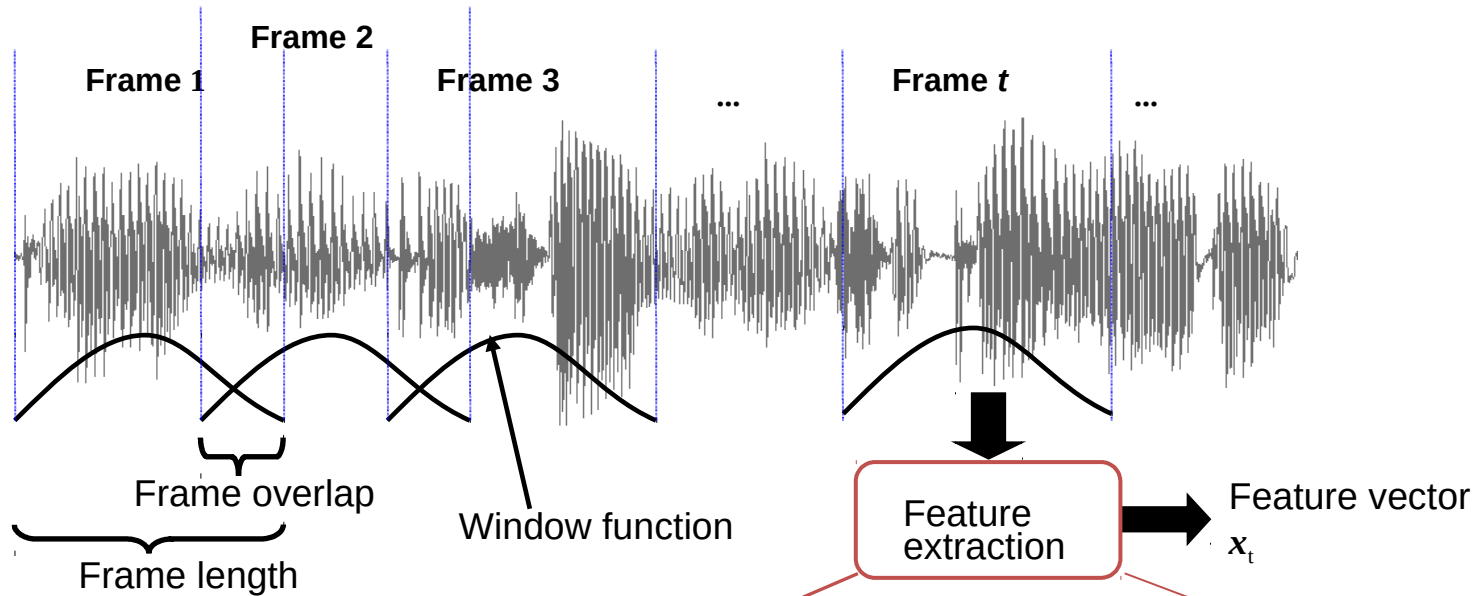
Frames shuffled, 30 ms frame



Frames shuffled, 100 ms frame



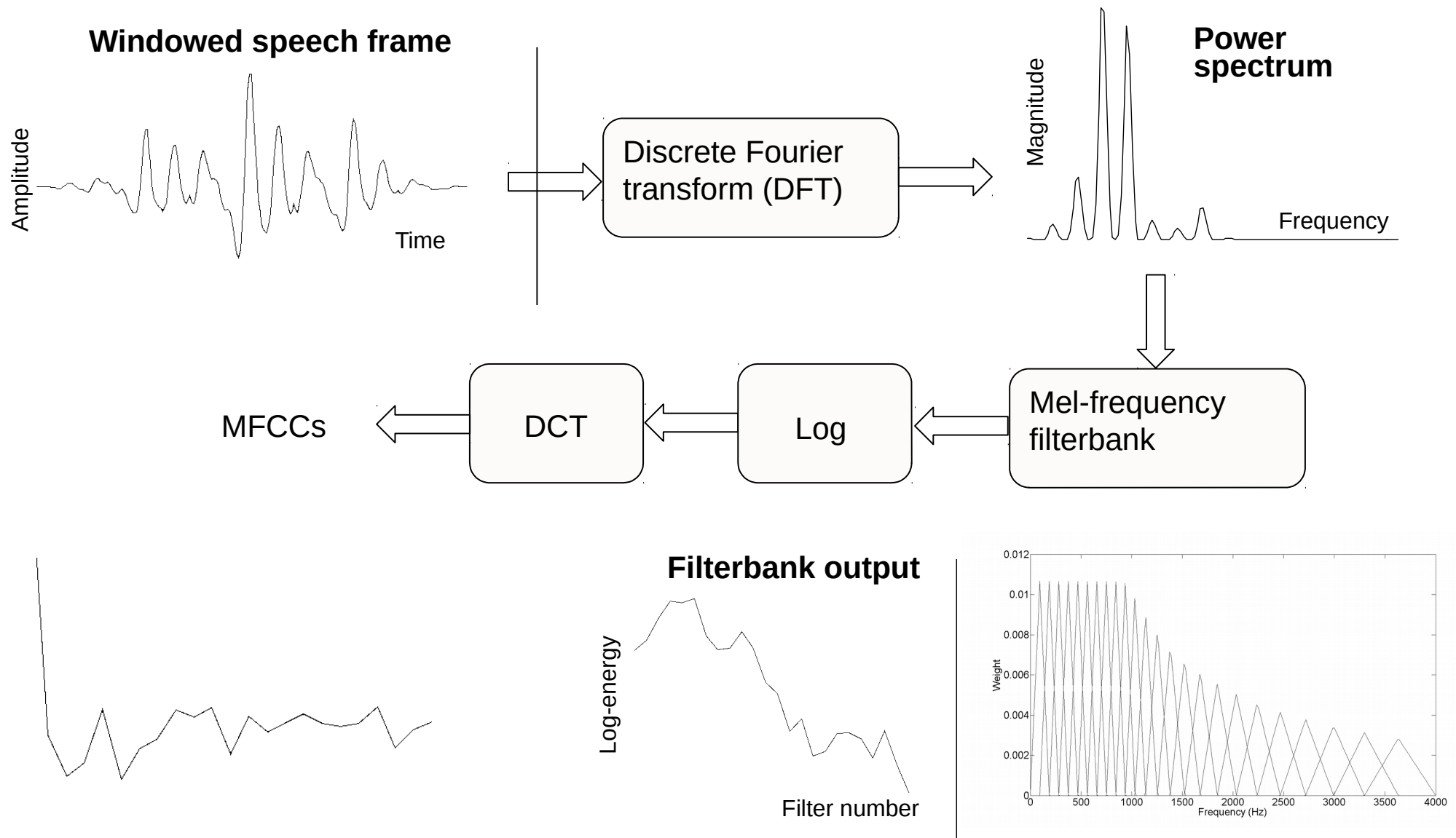
MFCC extraction



RASTA: RelAtive SpecTrA

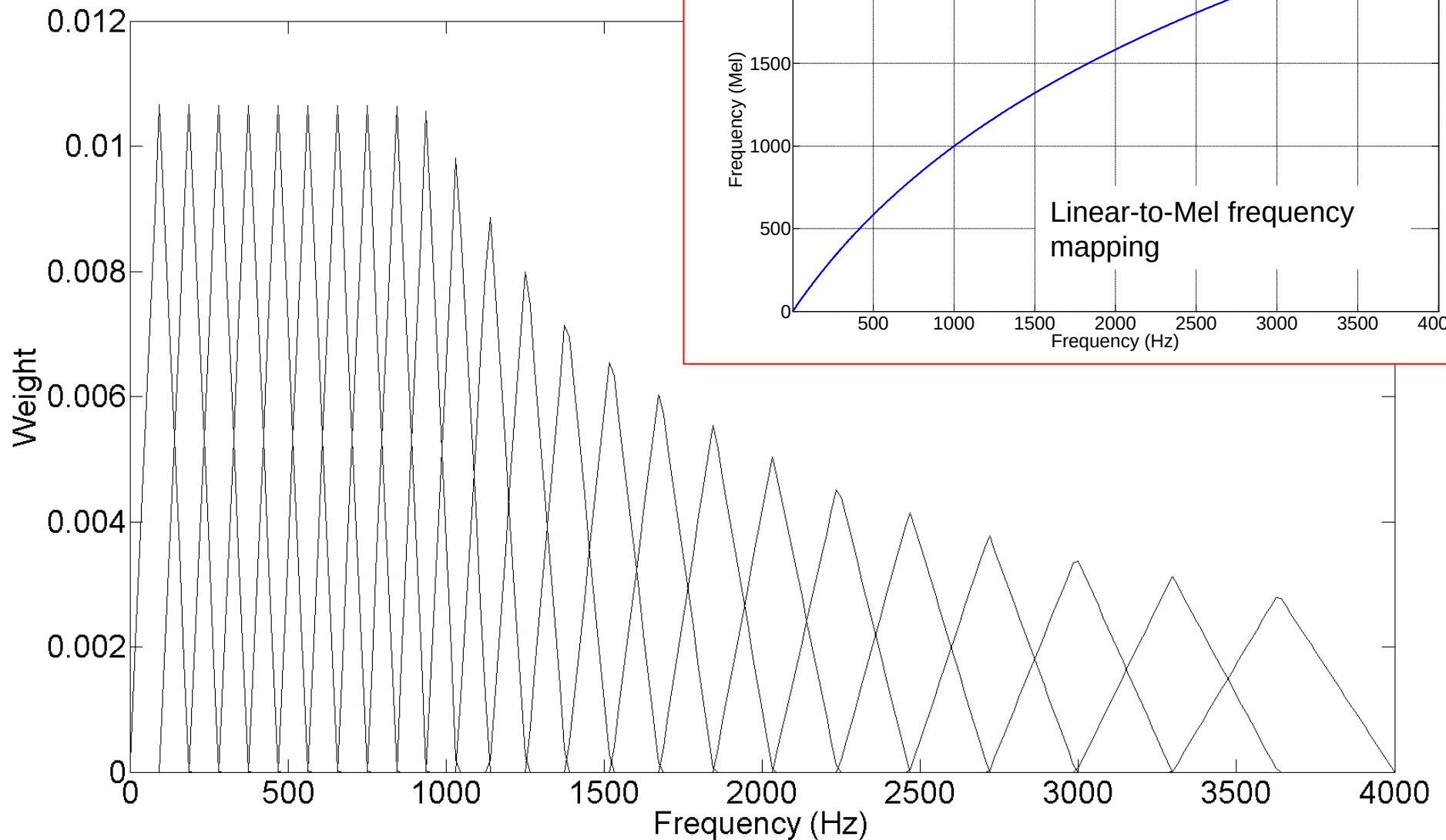
CMVN: Cepstral Mean and Variance Normalization

Mel-frequency cepstral coefficient (MFCC) features



Mel-frequency filterbank

[Generated using 'RASTAmat' package of Dan Ellis]

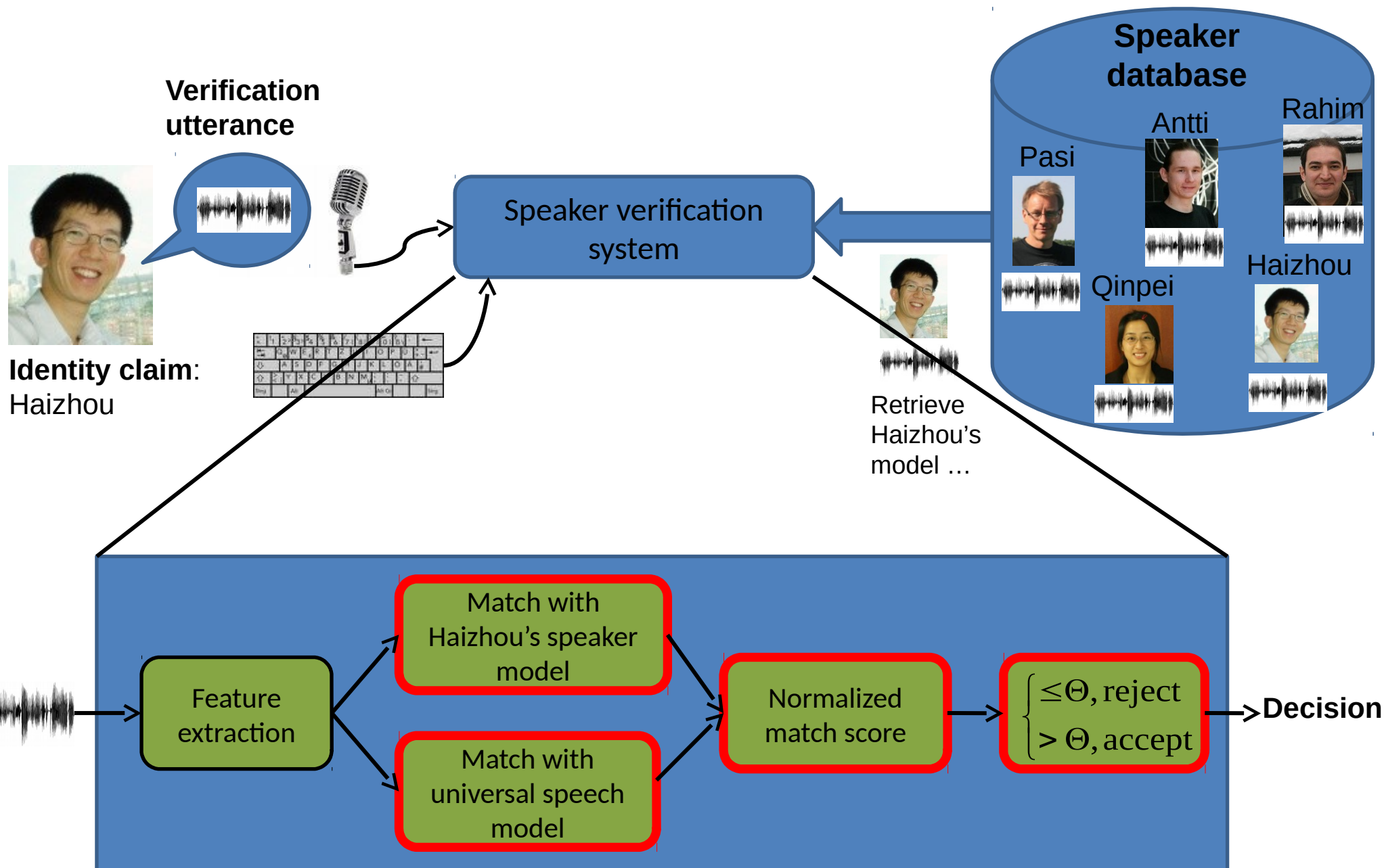


Speaker verification: science and art of data-driven modeling

“... from the speaker-recognition research trend in the last decade, it seems that improving feature robustness beyond a certain level (for a variety of degradations) is extremely difficult—or, in other words, **data-driven modeling techniques have been more successful in improving robustness compared to new features**”

[John H.L. Hansen and Taufiq Hasan, Speaker Recognition by Machines and Humans: A Tutorial Review, IEEE Signal Processing Magazine, Nov 2015]

Speaker modeling and comparison



History of speaker modeling

1970s and before

- Long-term feature averaging

1980s and 1990s

- Dynamic time warping (DTW), vector quantiz. (VQ)
- Hidden Markov Models
- Early neural net models

- All rooted on GMMs
- 1996 onwards: NIST SREs
- Focus on text-independent models

~1995 to ~2005

- Gaussian mixture models (GMMs)
- Universal background model

2005—today

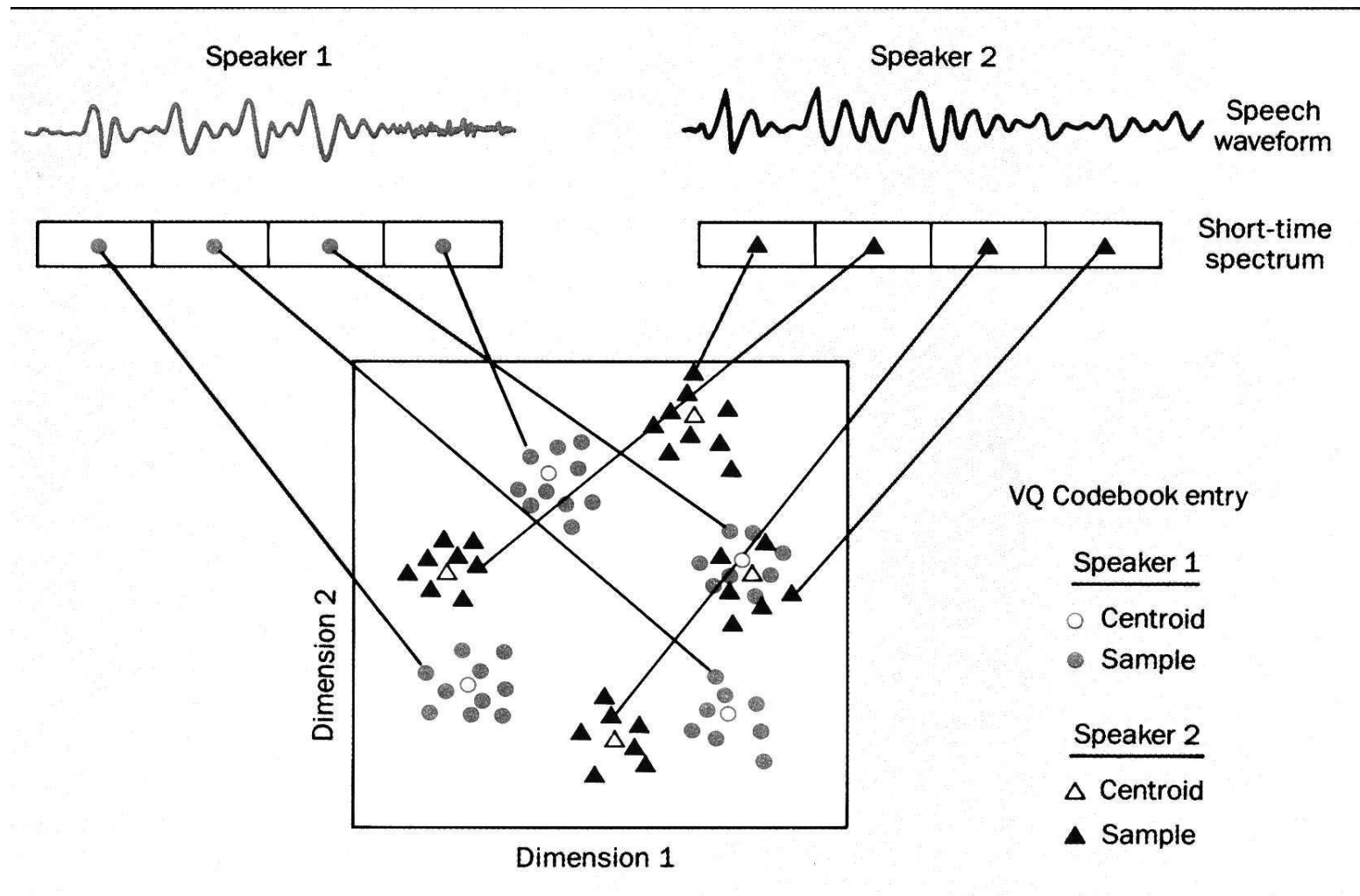
- GMM supervectors
- Joint factor analysis (JFA)
- i-vectors
- Probabilistic linear discr. analysis (PLDA) scoring
- Deep neural nets

Pre-history

Modern era

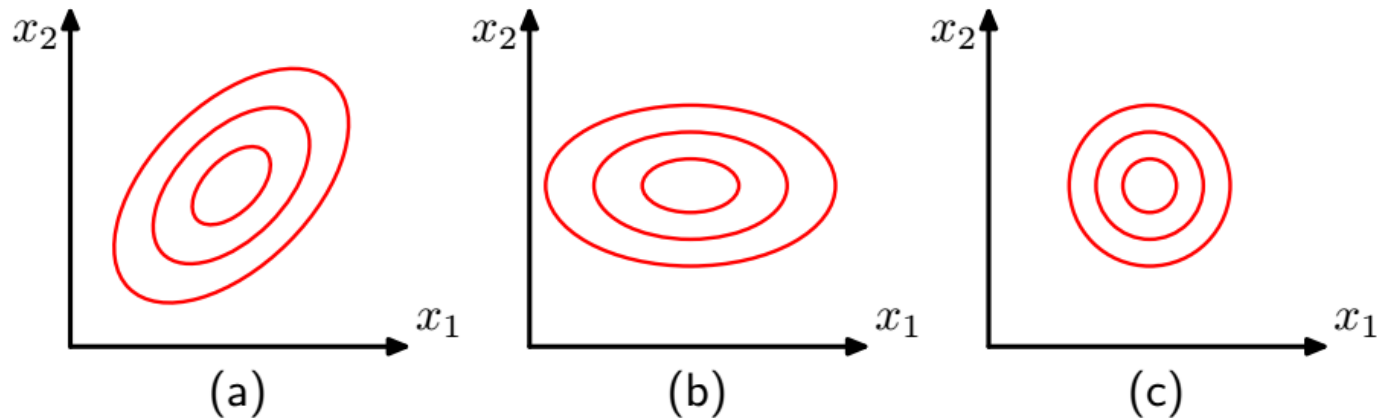
Speaker / language recognition is about modeling the "feature clouds"

[F. K. Soong, A. E. Rosenberg, L. R. Rabiner and B. H. Juang, "A Vector Quantization Approach to Speaker Recognition," *AT&T Technical Journal*, Vol. 66, pp. 14-26, Mar/Apr 1987]



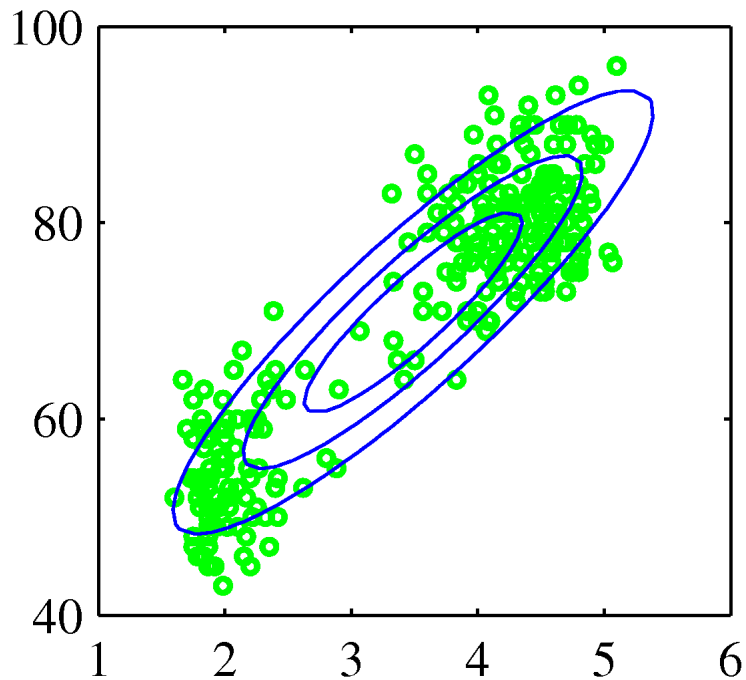
Single Gaussian

Bi-variate normal distribution

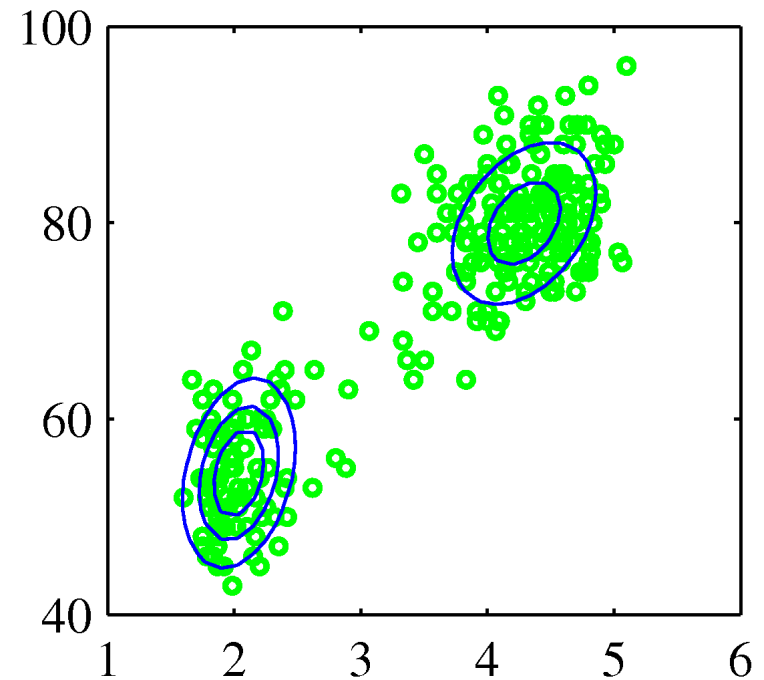


- ▶ (a) Covariance matrix Σ is full, $\text{rank } \Sigma = 2$, (b) $\text{diag}(\Sigma) = (a, b)^t$ and (c) $\Sigma = a * I$.
- ▶ How many parameters we need to estimate in (a), (b) , (c) ?

Gaussian mixture model (GMM)

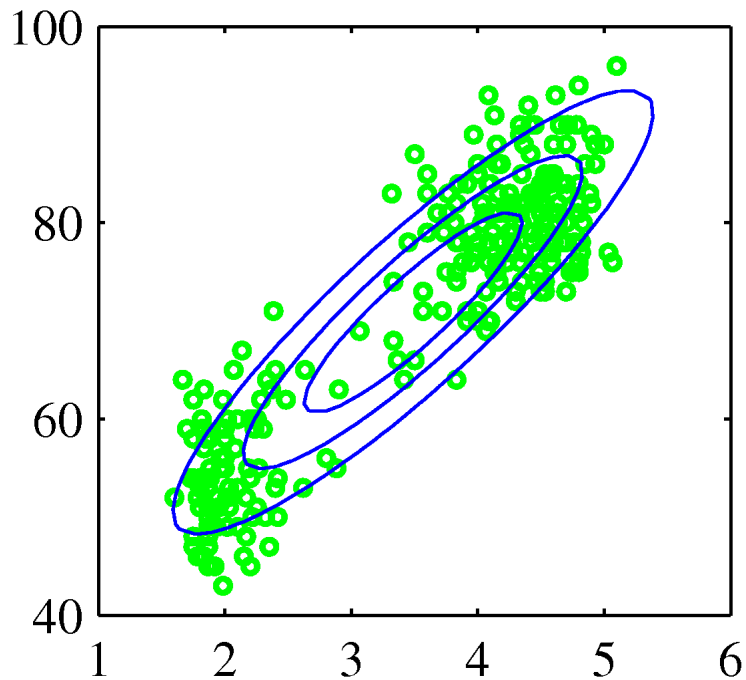


1 Gaussian not good enough ?

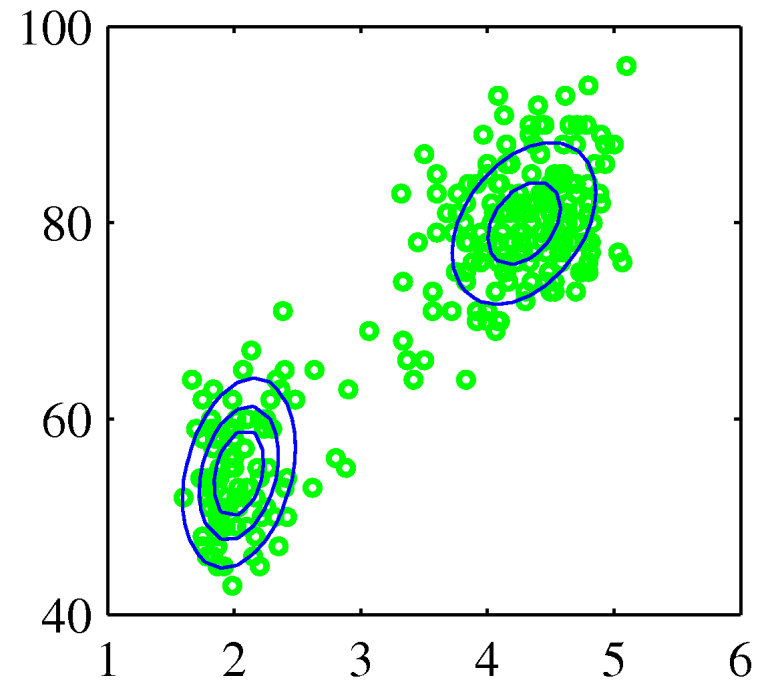


Looks better, yes? ^

Gaussian mixture model (GMM)



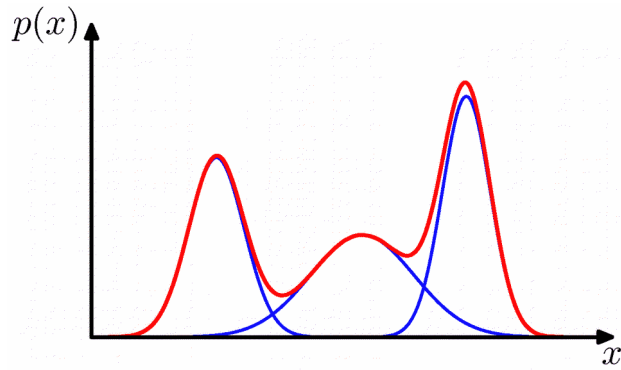
1 Gaussian not good enough ?



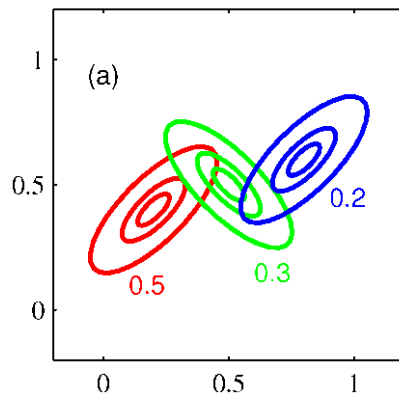
Looks better, yes? ^

GMM as a density estimator

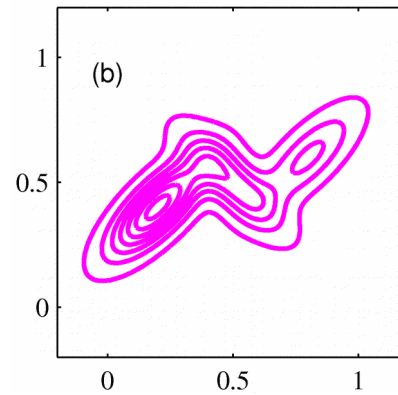
1-d data



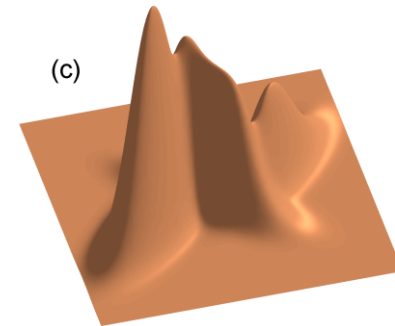
2-d data



Individual components

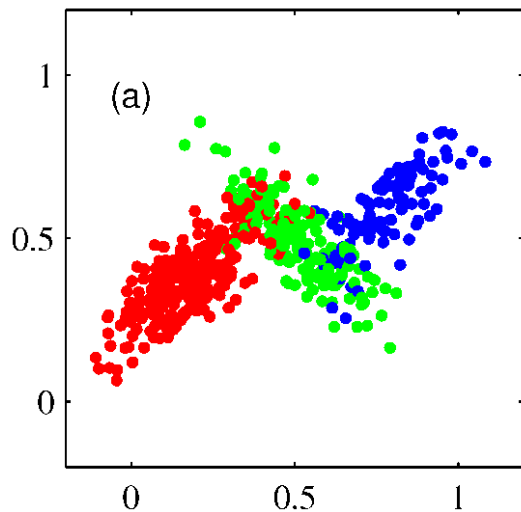


Density contour

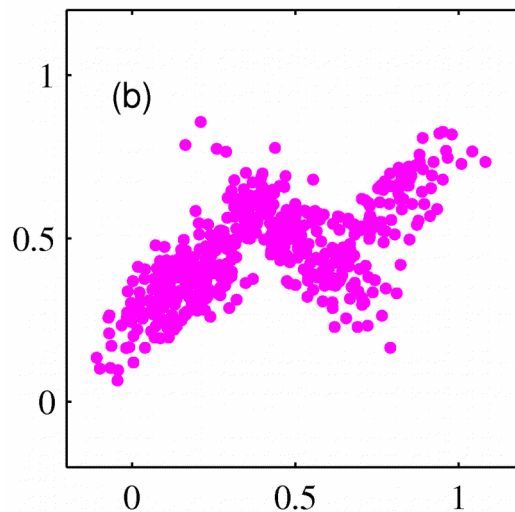


... another visualization

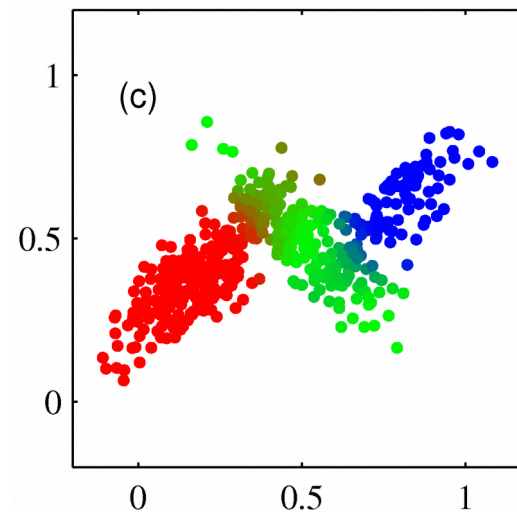
GMM as a latent (hidden) variable model



"Oracle": Assignment of data points to components is known (red, green or blue)

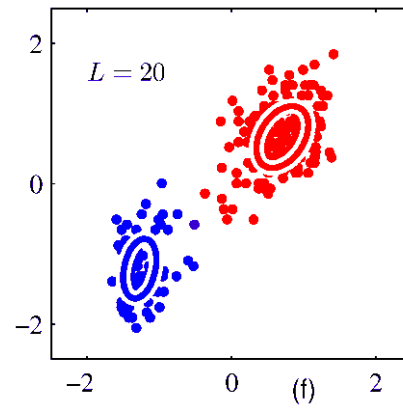
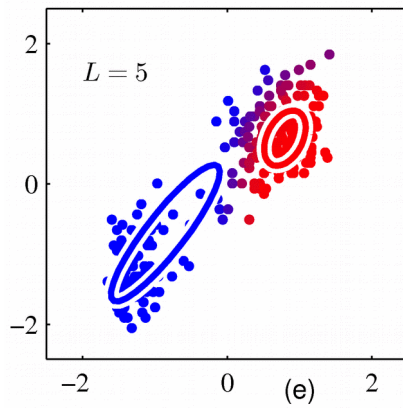
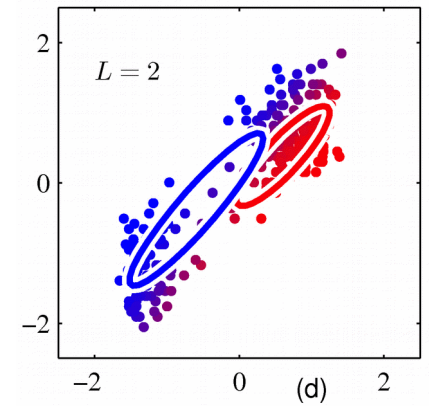
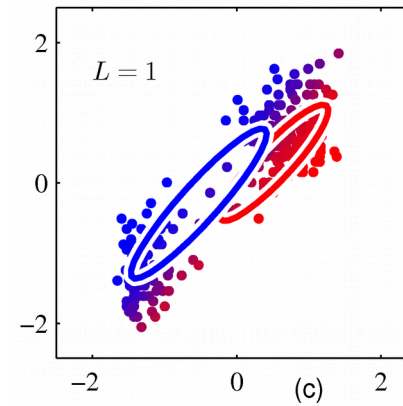
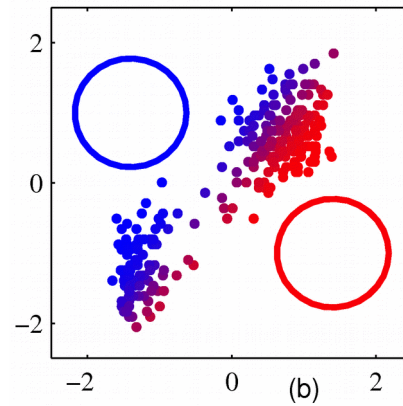
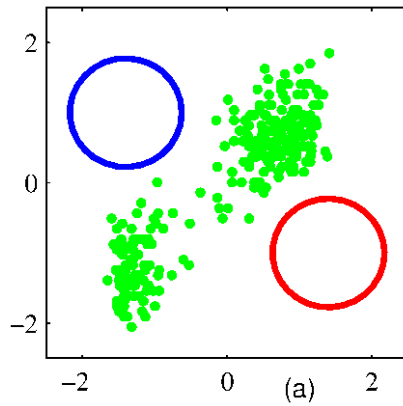


... but this is how our training sample is given: no labels!



strength of each color = posterior probability that the given observation was generated by the corresponding component

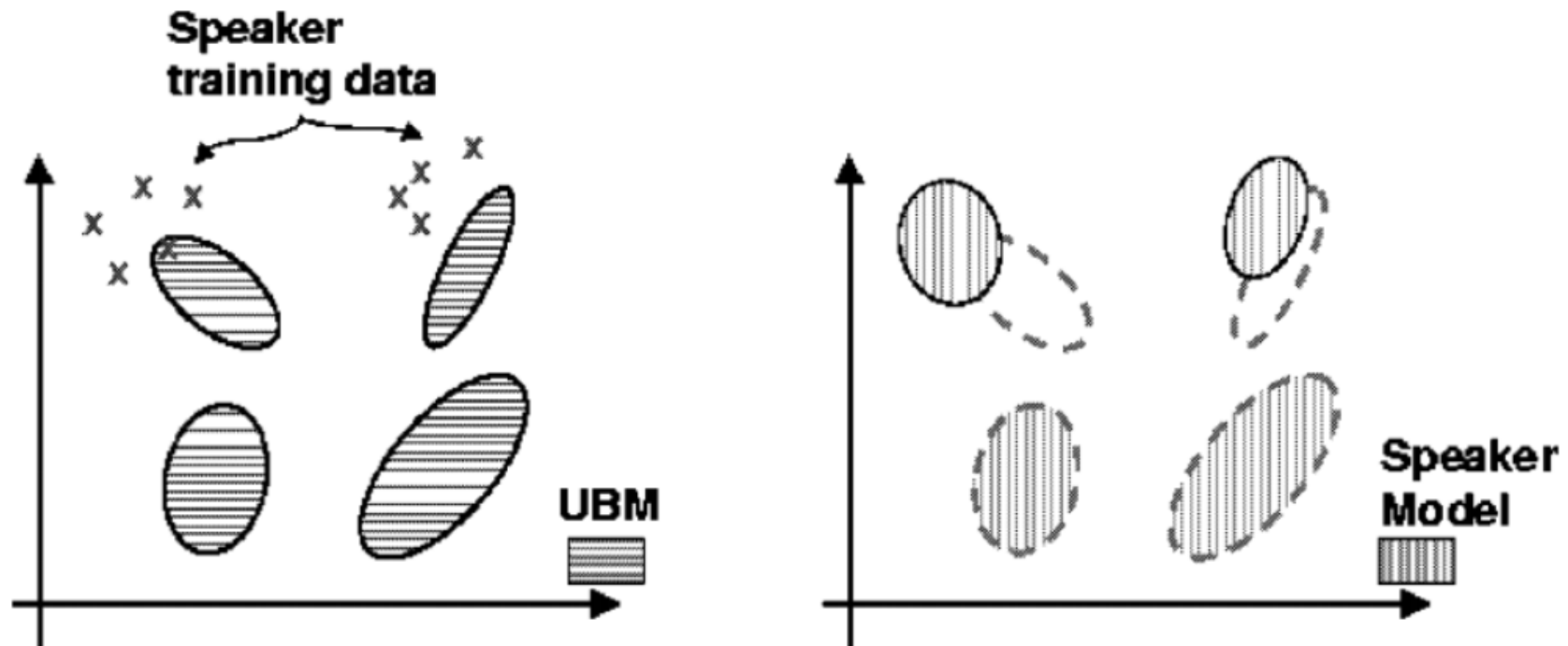
Progress of EM algorithm



GMM as a 'data compressor'

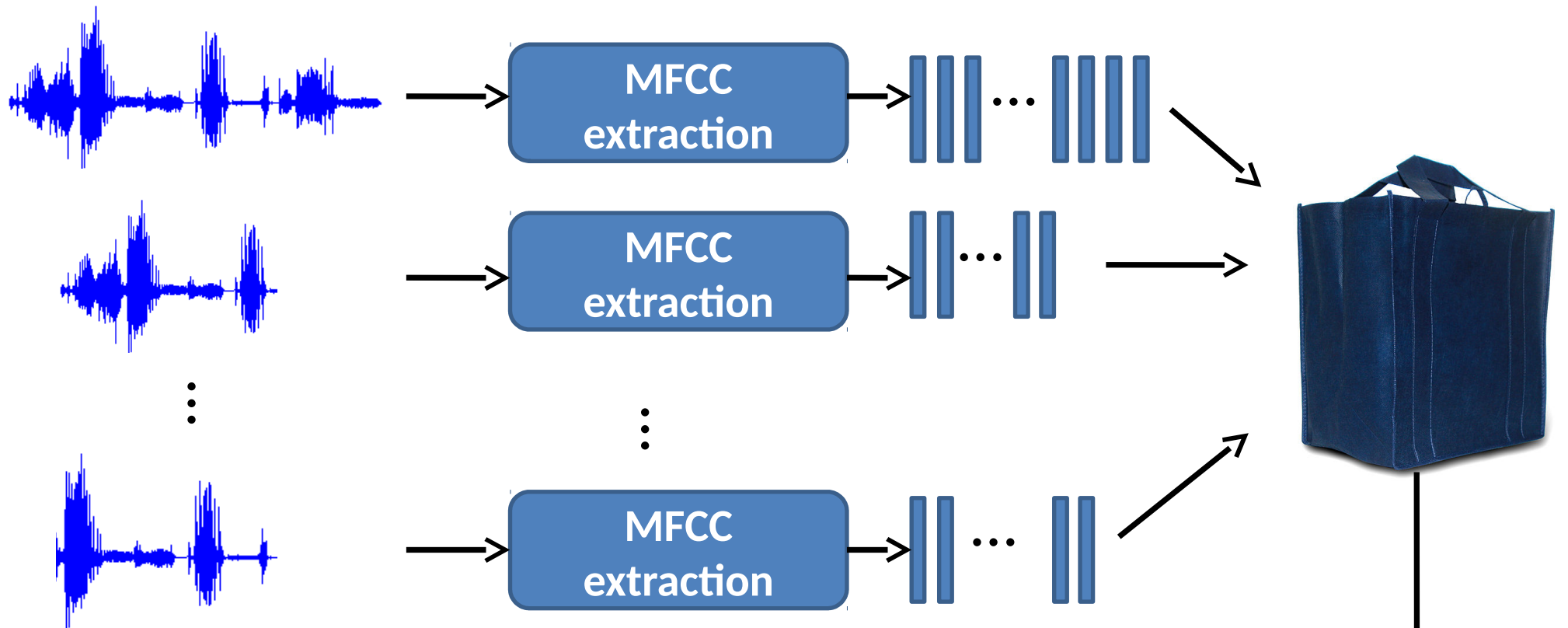
- Assume 1 minute recording of audio sampled with 8 kHz.
Raw data has $8000 \times 60 = \mathbf{480,000}$ **samples**
- Assume 12 MFCCs extracted every 10 ms $\underline{12}$ 12 coeffs x 100 frames/sec x 60 sec = **72,000 feature values**
- A GMM with 128 Gaussians & diagonal covariance matrices has
 - 128 mean vectors (dimensionality = 12)
 - 128 variance vectors (dimensionality = 12)
 - 128 mixing weights
 - Total $128 \times 12 + 128 \times 12 + 128 = \mathbf{3200}$ **parameters**

Gaussian mixture model – universal background model (GMM-UBM)



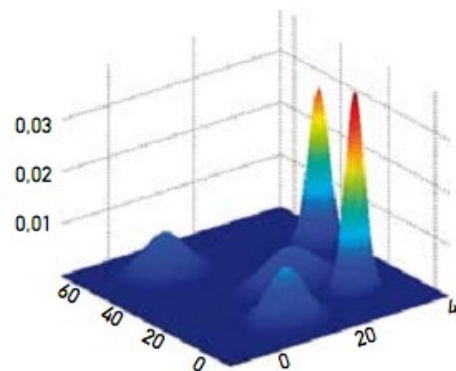
Douglas A. Reynolds and Thomas F. Quatieri and Robert B. Dunn, "Speaker Verification Using Adapted Gaussian Mixture Models", *Digital Signal Processing*, 2000.

Step 1: training the UBM



Utterances from a large number of different speakers

- No need for speaker identity, transcripts or other metadata
- Sometimes, gender information can be useful

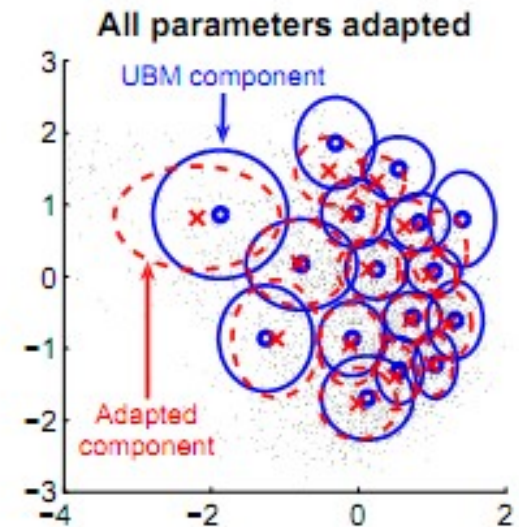
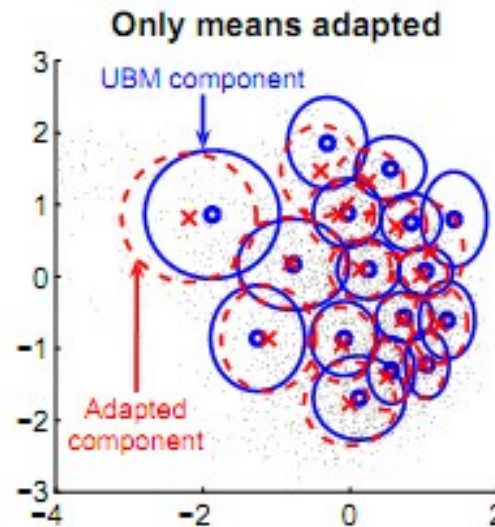


Expectation maximization (EM) algorithm

Step 2: speaker enrollment via maximum a posteriori (MAP) adaptation

$$\mu'_k = \alpha_k \tilde{\mathbf{x}}_k + (1 - \alpha_k) \mu_k$$

μ'_k : Adapted mean vector
 $\tilde{\mathbf{x}}_k$: Mean of training data
 μ_k : Prior mean from univ. background model (UBM)



$$\alpha_k = \frac{n_k}{n_k + r}$$

Adaptation coefficient
(r = relevance factor)

$$\tilde{\mathbf{x}}_k = \frac{1}{n_k} \sum_{t=1}^T P(k|\mathbf{x}_t) \mathbf{x}_t$$

Mean of training data assigned to k th Gaussian

$$n_k = \sum_{t=1}^T P(k|\mathbf{x}_t)$$

Soft count of vectors assigned to k th Gaussian

$$P(k|\mathbf{x}_t) = \frac{P_k \mathcal{N}(\mathbf{x}_t | \mu_k, \Sigma_k)}{\sum_{m=1}^K P_m \mathcal{N}(\mathbf{x}_t | \mu_m, \Sigma_m)}$$

Posterior probability of the k th Gaussian for one feature vector

Step 3: speaker comparison

- $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ is a sequence of test utterance feature vectors and $\boldsymbol{\theta}_s$ is a GMM of speaker s , claimed to have 'generated' \mathbf{X} , adapted from the UBM as explained above.
- We compute log-likelihood ratio
$$\text{LLR}(\mathbf{X}) = \log p(\mathbf{X} | \boldsymbol{\theta}_s) - \log p(\mathbf{X} | \boldsymbol{\theta}_{\text{ubm}})$$
- The larger the LLR value, the stronger the evidence for the "same speaker" hypothesis. Binary decision is made by comparing $\text{LLR}(\mathbf{X})$ against a pre-selected threshold (such as 0).
- Intuition: $\log p(\mathbf{X} | \boldsymbol{\theta}_{\text{ubm}})$ represents the uninformative part common to all speakers (such as speech content)

GMM supervectors

[W. M. Campbell, D. E. Sturim, D. A. Reynolds, "Support vector machines using GMM supervectors for speaker verification", *IEEE Signal Proc Lett* 2006]

