# Lecture 9:
# Speech Recognition: Front Ends

**1**  **Recognizing Speech**

**2**  **Feature Calculation**

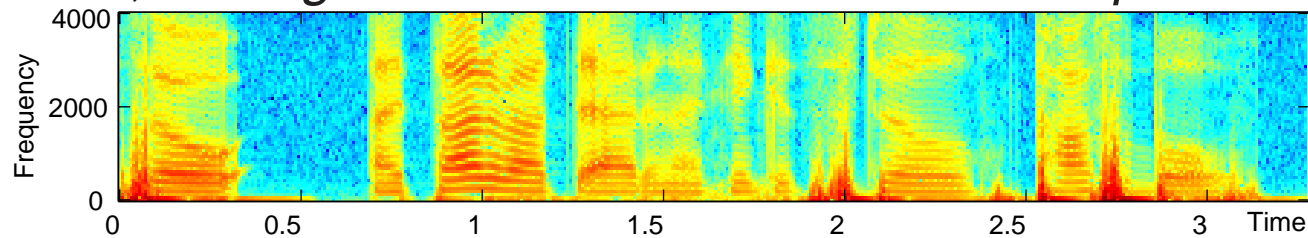Dan Ellis  <dpwe@ee.columbia.edu>
http://www.ee.columbia.edu/~dpwe/e6820/

Columbia University Dept. of Electrical Engineering
Spring 2003

# Recognizing Speech

*"So, I thought about that and I think it's still possible"*



- **What kind of information might we want from the speech signal?**
  - words
  - phrasing, 'speech acts' (prosody)
  - mood / emotion
  - speaker identity

- **What kind of processing do we need to get at that information?**
  - time scale of feature extraction
  - signal aspects to capture in features
  - signal aspects to *exclude* from features

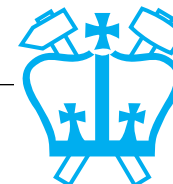# Speech recognition as Transcription

- **Transcription = "speech to text"**
  - find a word string to match the utterance

- **Best suited to small vocabulary tasks**
  - voice dialing, command & control etc.

- **Gives neat objective measure:**
  **word error rate (WER) %**
  - can be a sensitive measure of performance

- **Three kinds of errors:**

*Reference:* THE CAT SAT ON THE MAT

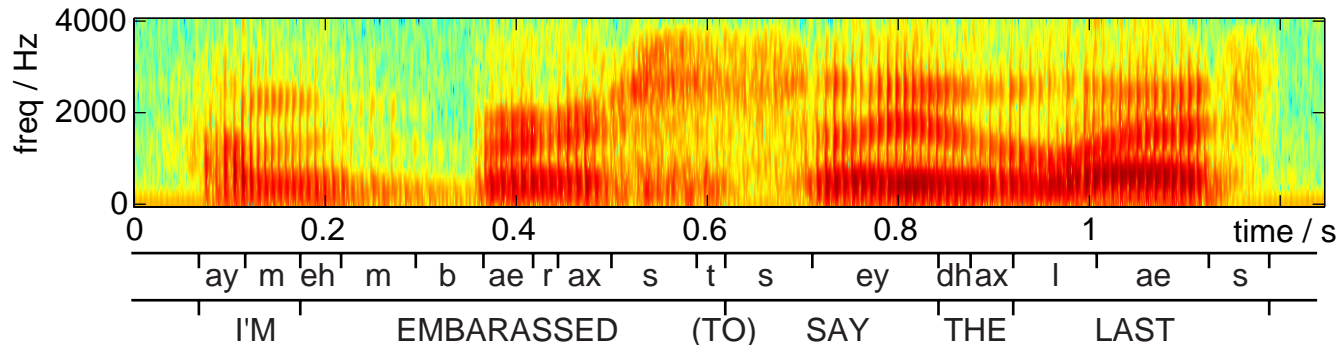*Recognized:* – CAT SAT AN THE A MAT

*Deletion*     *Substitution*     *Insertion*

- WER = (S + D + I) / N

# Limitations of the Transcription paradigm

- **Starts to fall down with 'natural' speech**
  - some "words" may not even exist



- **Word transcripts do not capture everything**
  - speaker changes, intonation, phrasing

- **Word error rate treats all errors as equal**
  - small words ("of") counted as big words
  - small differences ("company's" → "companies")
    vs. larger ("held police" → "health plans")

- **Move towards other measures**
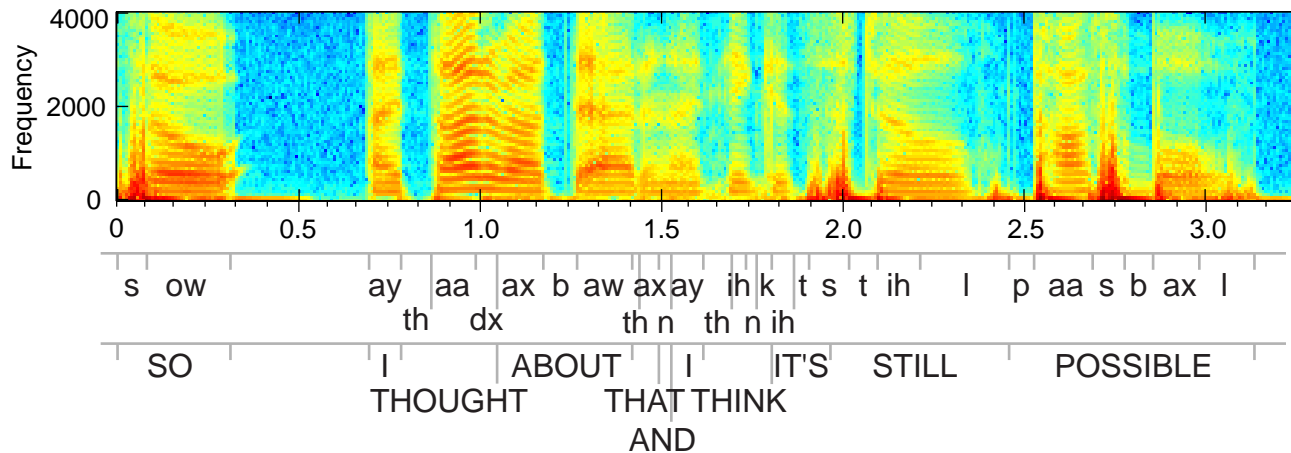  - e.g. task-defined:
    was the *meaning* recognized?

# Why is Speech Recognition hard?

- **Why not match against a set of waveforms?**
  - waveforms are never (nearly!) the same twice
  - speakers minimize information/effort in speech

- **Speech variability comes from many sources:**
  - speaker-dependent (SD) recognizers must handle within-speaker variability
  - speaker-independent (SI) recognizers must also deal with variation between speakers
  - all recognizers are afflicted by background noise, variable channels

- → **Need recognition models that:**
  - generalize i.e. accept variations in a range, and
  - adapt i.e. 'tune in' to a particular variant

# Within-speaker variability

- **Timing variation:**
  - word duration varies enormously



|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| s | ow |  | ay | aa | ax | b | aw | ax | ay | ih | k | t | s | t | ih | l | p | aa | s | b | ax l |
|   |   |   |   | th | dx |   |   | th | n |   | th n ih |   |   |   |   |   |   |   |   |   |   |

SO — I THOUGHT — ABOUT — I THAT THINK AND — IT'S STILL — POSSIBLE

  - fast speech 'reduces' vowels

- **Speaking style variation:**
  - careful/casual articulation
  - soft/loud speech

- **Contextual effects:**
  - speech sounds vary with context, role:
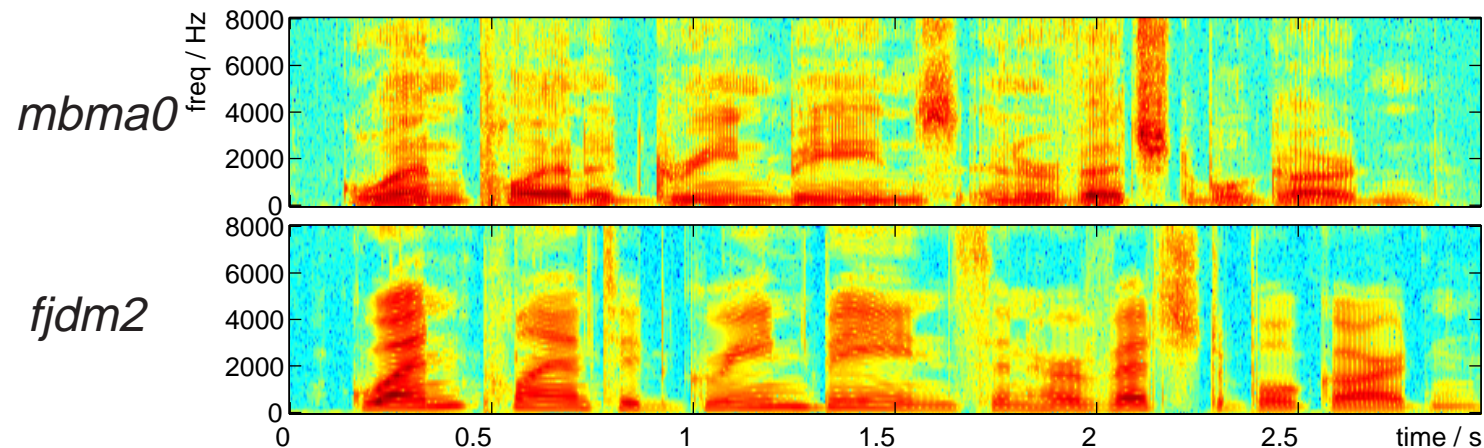    "How **do** you **do**?"

# Between-speaker variability

- **Accent variation**
  - regional / mother tongue

- **Voice quality variation**
  - gender, age, huskiness, nasality

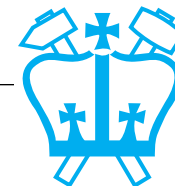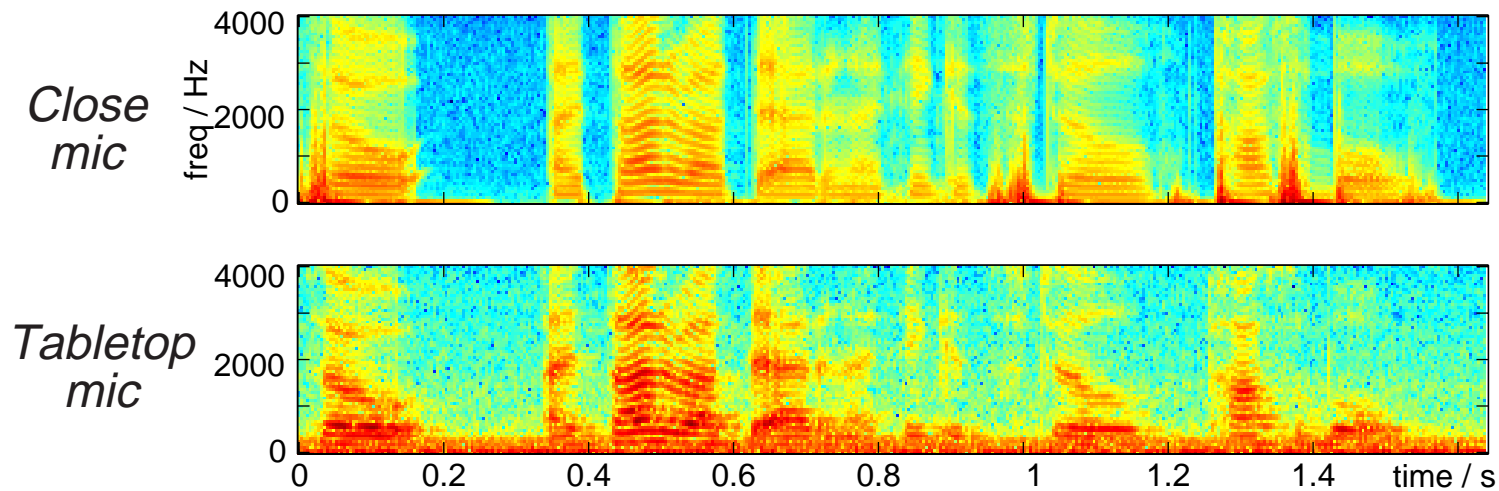- **Individual characteristics**
  - mannerisms, speed, prosody

# Environment variability

- **Background noise**
  - fans, cars, doors, papers

- **Reverberation**
  - 'boxiness' in recordings

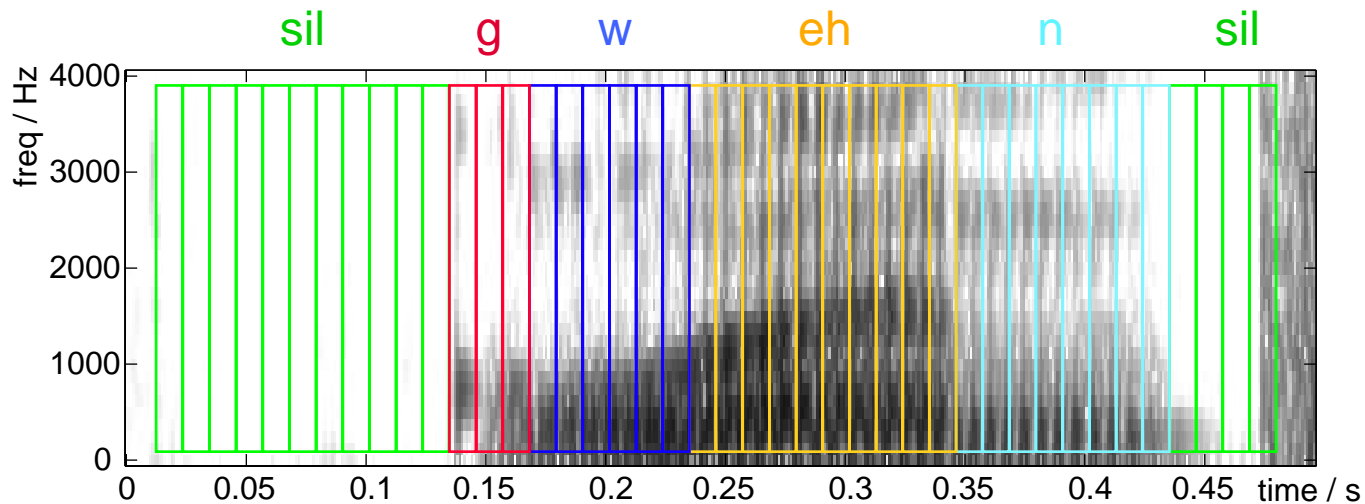- **Microphone channel**
  - huge effect on relative spectral gain

# How to recognize speech?

- **Cross correlate templates?**
  - waveform?
  - spectrogram?
  - time-warp problems

- **Match short-segments & handle time-warp later**
  - model with slices of ~ 10 ms
  - pseudo-stationary model of words:



  - other sources of variation...

# Which segments to use?

- **Assume words can be broken down into pseudo-stationary segments**
  - not a perfect fit, but worth a try

- **Linguists offer phonemes or phones**
  - phonemes are the minimal set needed to disambiguate words
  - phones are realizations of phonemes

- **Other possibilities:**
  - data-clustering techniques to define segments 'intrinsically'
  - lesson from synthesis: transitions as important or more important than steady portions?

  ...but how to model?

# Probabilistic formulation

- **Probability** **that segment label is correct**

  - gives standard form of speech recognizers:

- **Feature calculation**

  transforms signal into easily-classified domain

  $$s[n] \rightarrow X_m \quad \left( m = \frac{n}{H} \right)$$
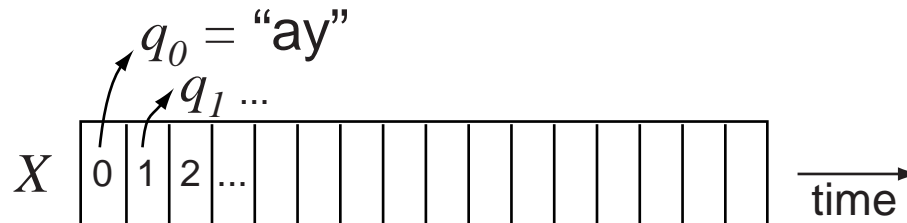
- **Acoustic classifier**

  calculates probabilities of each mutually-exclusive state $q^i$
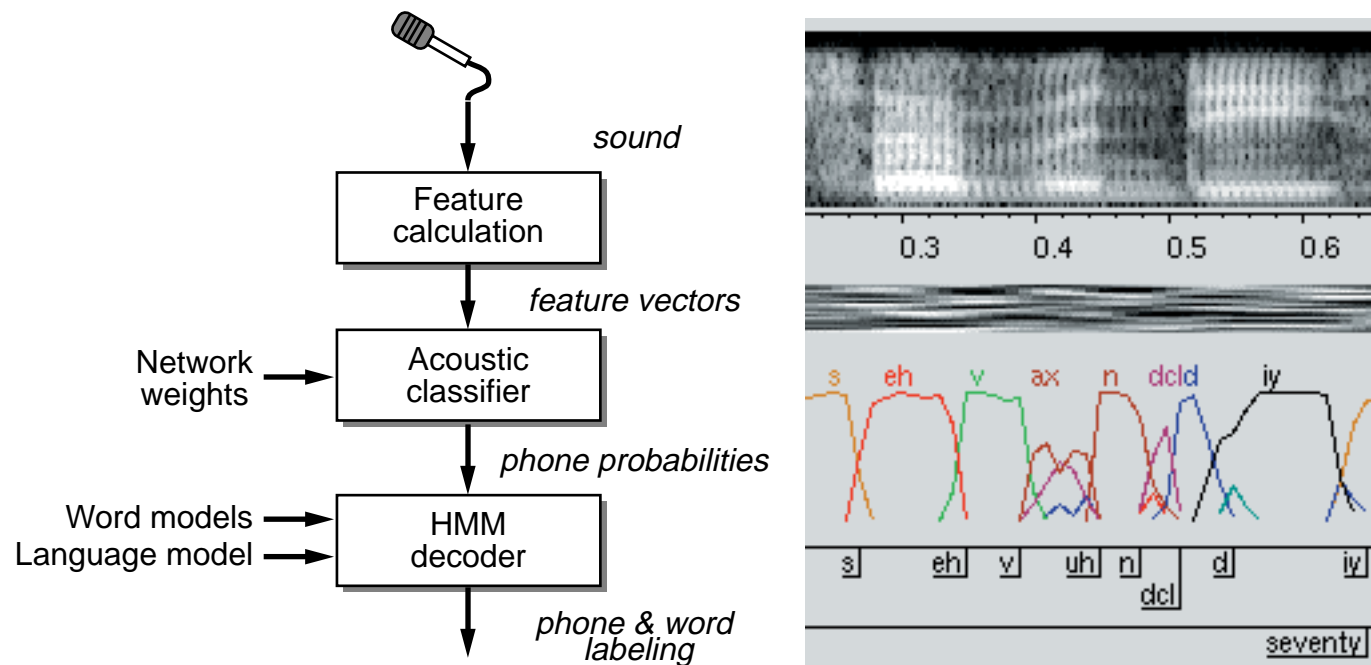
  $$p(q^i | X)$$

- **'Finite state acceptor' (i.e. HMM)**

  $$\hat{Q} = \underset{\{q_0, q_1, \dots q_L\}}{\mathrm{argmax}} p(q_0, q_1, \dots q_L | X_0, X_1 \dots X_L)$$

  MAP match of allowable sequence to probabilities:

  $q_0 = \text{"ay"}$

  $q_1 \dots$

  $X$  | 0 | 1 | 2 |... |   |   |   |   |   |   |   |   |   |   |   |  time →

# Standard speech recognizer structure



- **Questions:**
  - what are the best features?
  - how do we do the acoustic classification?
  - how do we find/match the state sequence?

# Outline

**1**    **Recognizing Speech**

**2**    **Feature Calculation**

- Spectrogram, MFCCs & PLP
- Improving robustness

## 2     Feature Calculation

- **Goal: Find a representational space most suitable for classification**
    - waveform: voluminous, redundant, variable
    - spectrogram: better, still quite variable
    - ...?

- **Pattern Recognition:**
  **Representation is upper bound on performance**
    - maybe we *should* use the waveform...
    - or, maybe the representation can do *all* the work

- **Feature calculation is intimately bound to classifier**
    - pragmatic strengths and weaknesses

- **Features develop by slow evolution**
    - current choices more historical than principled

# Desired characteristics for features

- **Provide the 'right' information**
  - extract signal information for classification task
  - suppress irrelevant information

- **Be compatible with acoustic classifier**
  - relatively low dimensionality
  - uncorrelated dimensions?

- **Be practical**
  - applicable in 'all' circumstances
  - relatively inexpensive to compute

- **Be robust**
  - so far as possible, exclude nonspeech information

→ **How to evaluate features?**
  - normally: just put them in a recognizer

# Features (1): Spectrogram

- **Plain STFT as features e.g.**

$$X_m[k] = S[mH, k] = \sum_n s[n + mH] \cdot w[n] \cdot e^{-(j2\pi kn)/N}$$

- **Consider examples:**



*Feature vector slice*

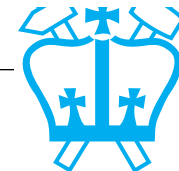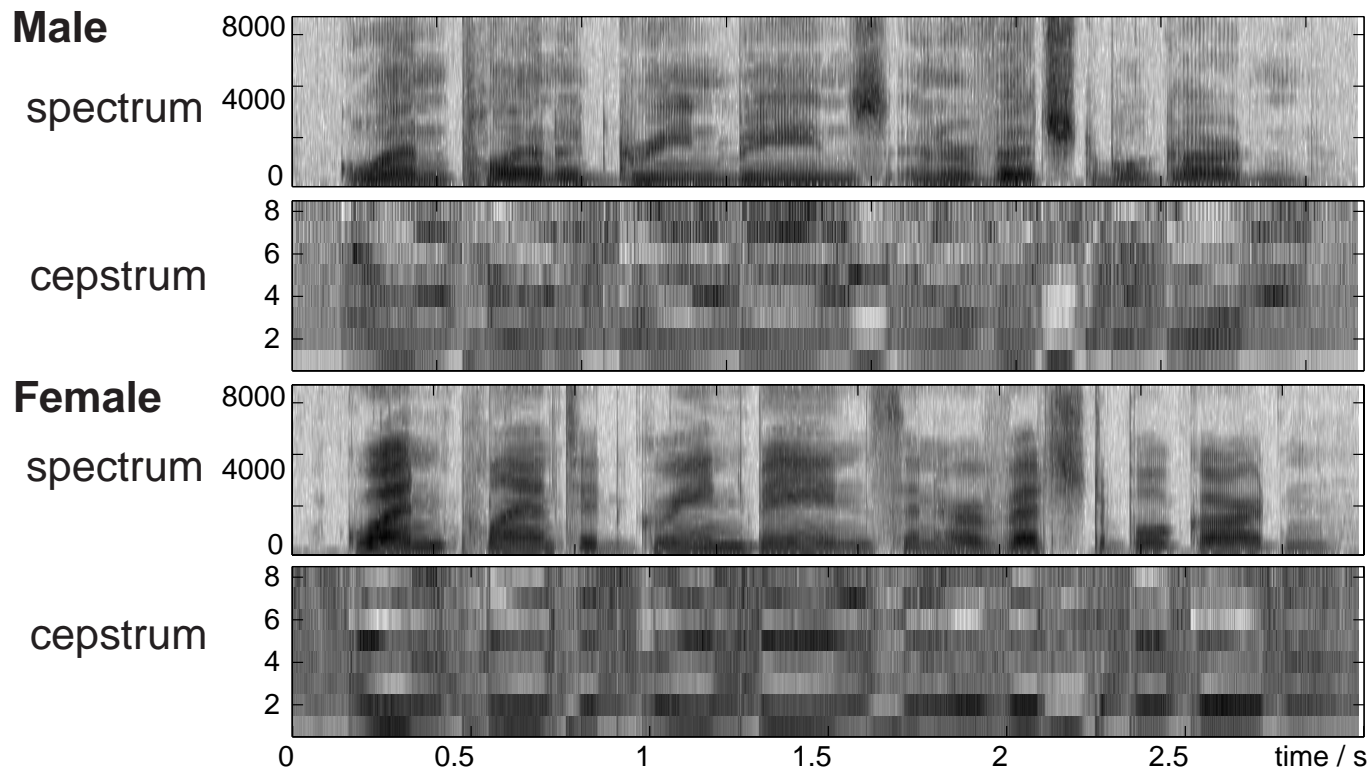- **Similarities between corresponding segments**
  - but still large differences

# Features (2): Cepstrum

- **Idea: Decorrelate, summarize spectral slices:**

$$X_m[l] = IDFT\{\log|S[mH, k]|\}$$

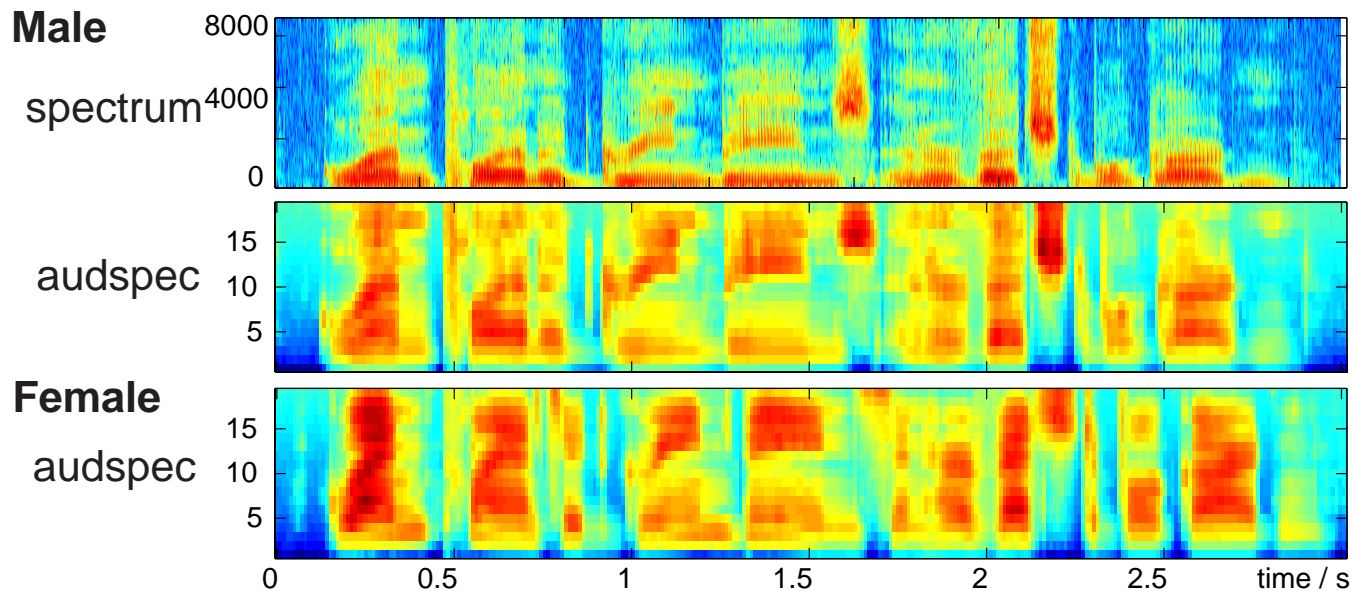  - good for Gaussian models
  - greatly reduce feature dimension

# Features (3): Frequency axis warp

- **Linear frequency axis gives equal 'space' to 0-1 kHz and 3-4 kHz**
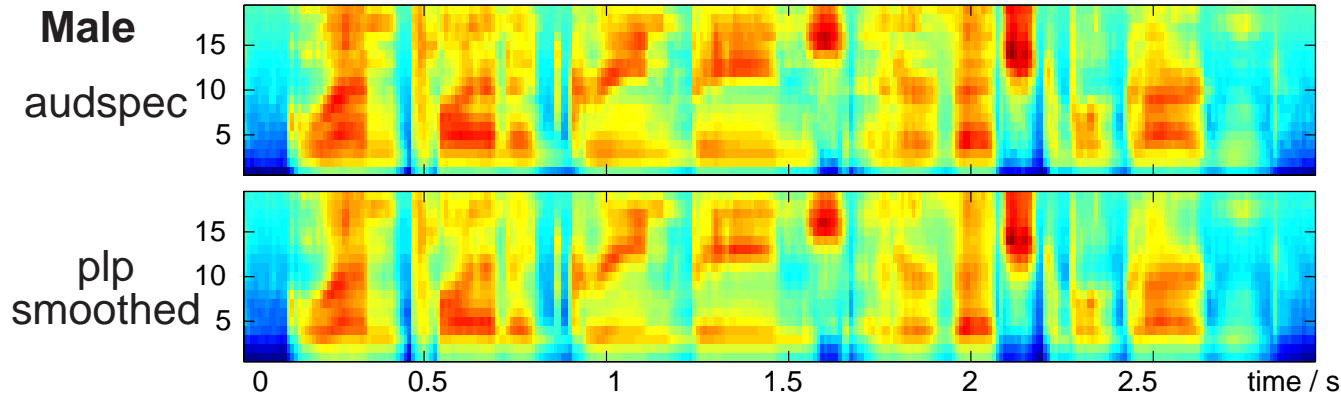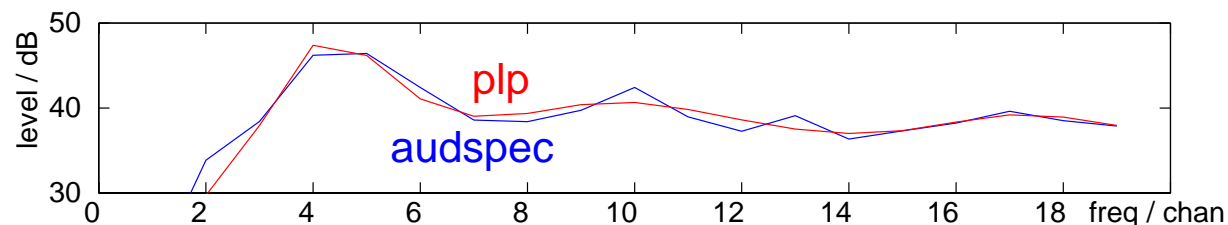
  - but perceptual importance very different

- **Warp frequency axis closer to perceptual axis:**

$$X[c] = \sum_{k=l_c}^{u_c} |S[k]|^2$$

  - mel, Bark, constant-Q ...



**Male** spectrum

**Female** audspec

# Features (4): Spectral smoothing

- **Generalizing across different speakers is helped by smoothing (i.e. *blurring*) spectrum**

- **Truncated cepstrum is one way:**
  - MSE approx to $\log |S[k]|$

- **LPC modeling is a little different:**
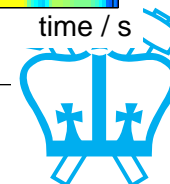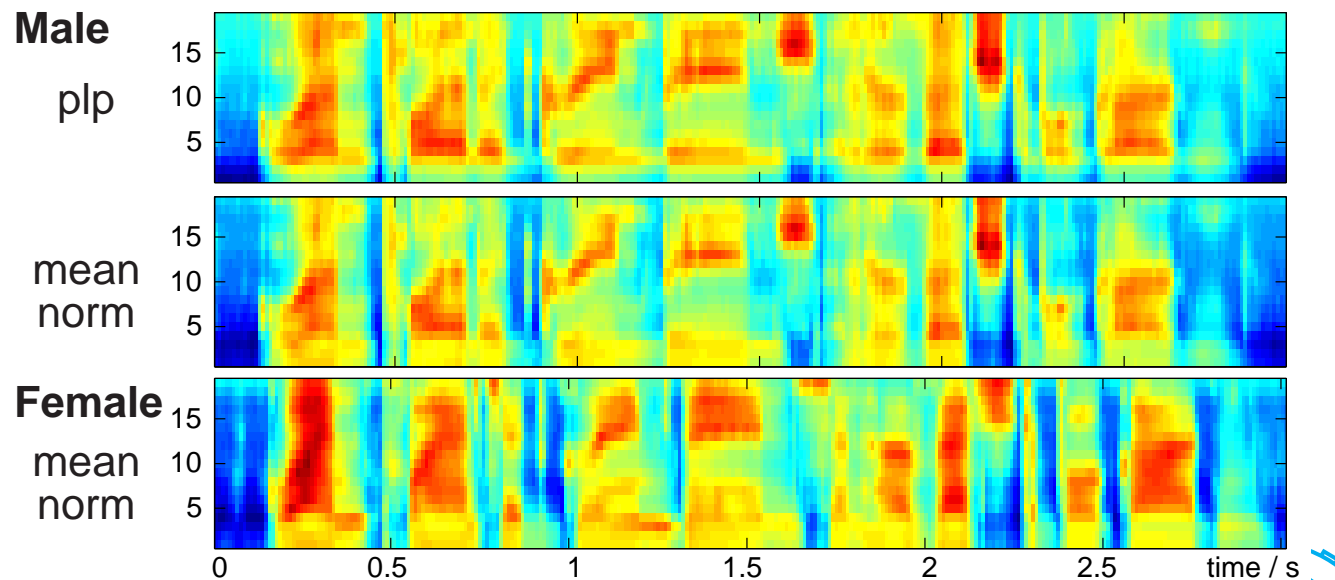  - MSE approx to $|S[k]| \rightarrow$ prefers detail at peaks

# Features (5): Normalization along time

- **Idea: feature variations, not absolute level**

- **Hence: calculate average level & subtract it:**

$$X[k] = S[k] - \text{mean}\{S[k]\}$$

- **Factors out fixed channel frequency response:**

$$s[n] = h[n] * e[n]$$

$$\log|S[k]| = \log|H[k]| + \log|E[k]|$$



**Male**

plp

15 10 5

mean norm

15 10 5

**Female** mean norm

15 10 5

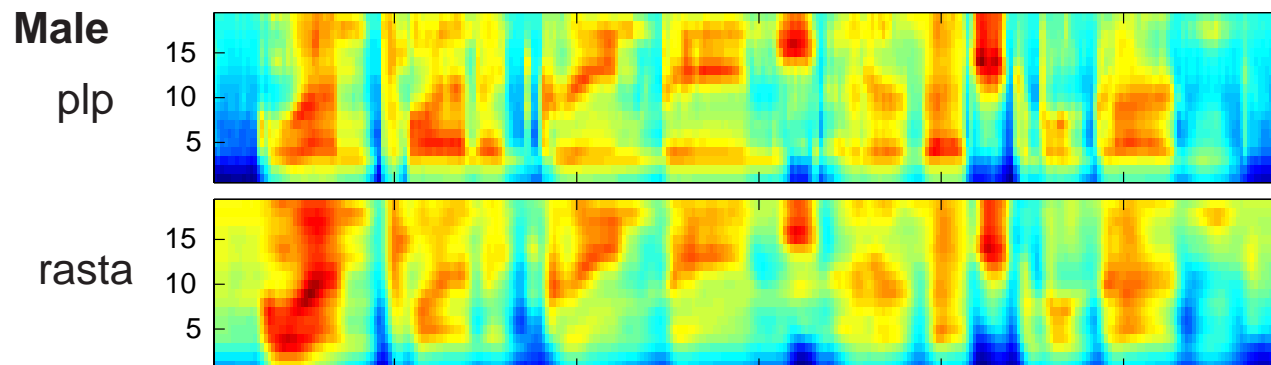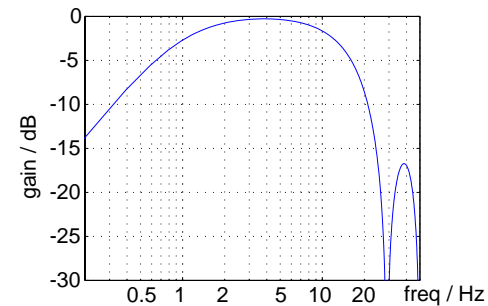0    0.5    1    1.5    2    2.5    time / s

# Features (6): RASTA filtering

- **Mean subtraction $\approx$ high-pass filtering along time in log-spectral domain**

$$X[k] = S[k] - \mathbf{lpf}\{S[k]\}$$

- **+ smooth along time for more blurring**

- $\rightarrow$ **Bandpass filter in time**

  - relates to 'modulation sensitivity' in hearing?



**Male**

plp

rasta
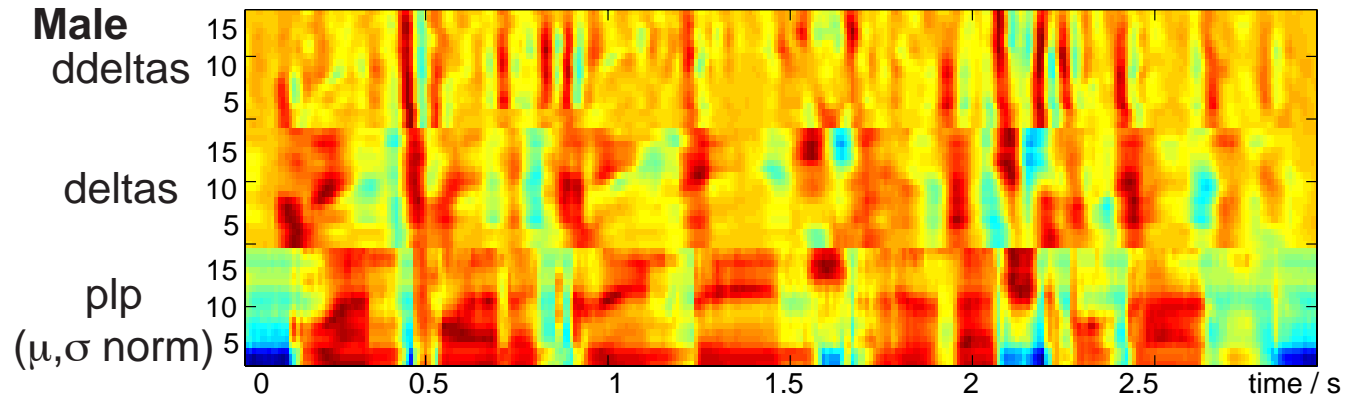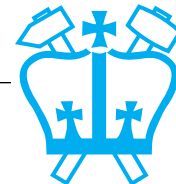
# Delta features

- **Want each segment to have 'static' feature vals**
  - but some segments intrinsically dynamic!

  →calculate their derivatives - maybe steadier?

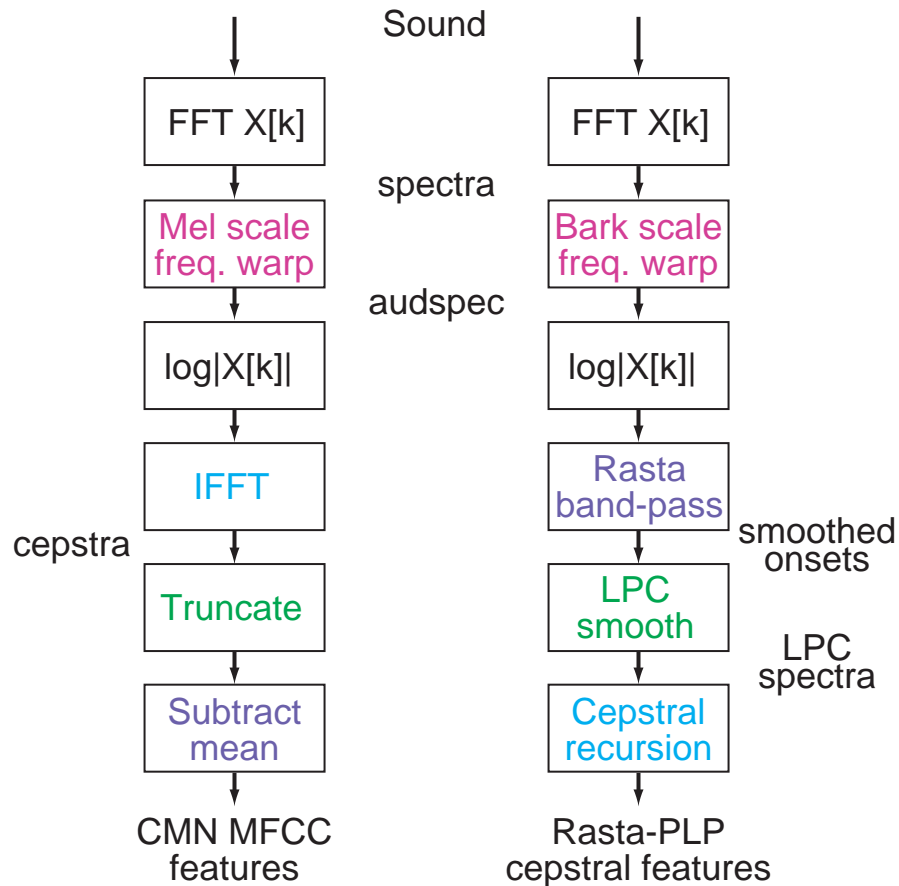- **Append $dX/dt$ (+ $d^2X/dt^2$) to feature vectors**



**Male** ddeltas / deltas / plp ($\mu,\sigma$ norm)

- **Relates to onset sensitivity in humans?**

# Overall feature calculation

- **MFCCs and/or RASTA-PLP**

Sound

| FFT X[k] | FFT X[k] |

spectra

| Mel scale freq. warp | Bark scale freq. warp |

audspec

| log\|X[k]\| | log\|X[k]\| |

| IFFT | Rasta band-pass |

cepstra | | smoothed onsets |

| Truncate | LPC smooth |

| | LPC spectra |

| Subtract mean | Cepstral recursion |

CMN MFCC features      Rasta-PLP cepstral features

- **Key attributes:**
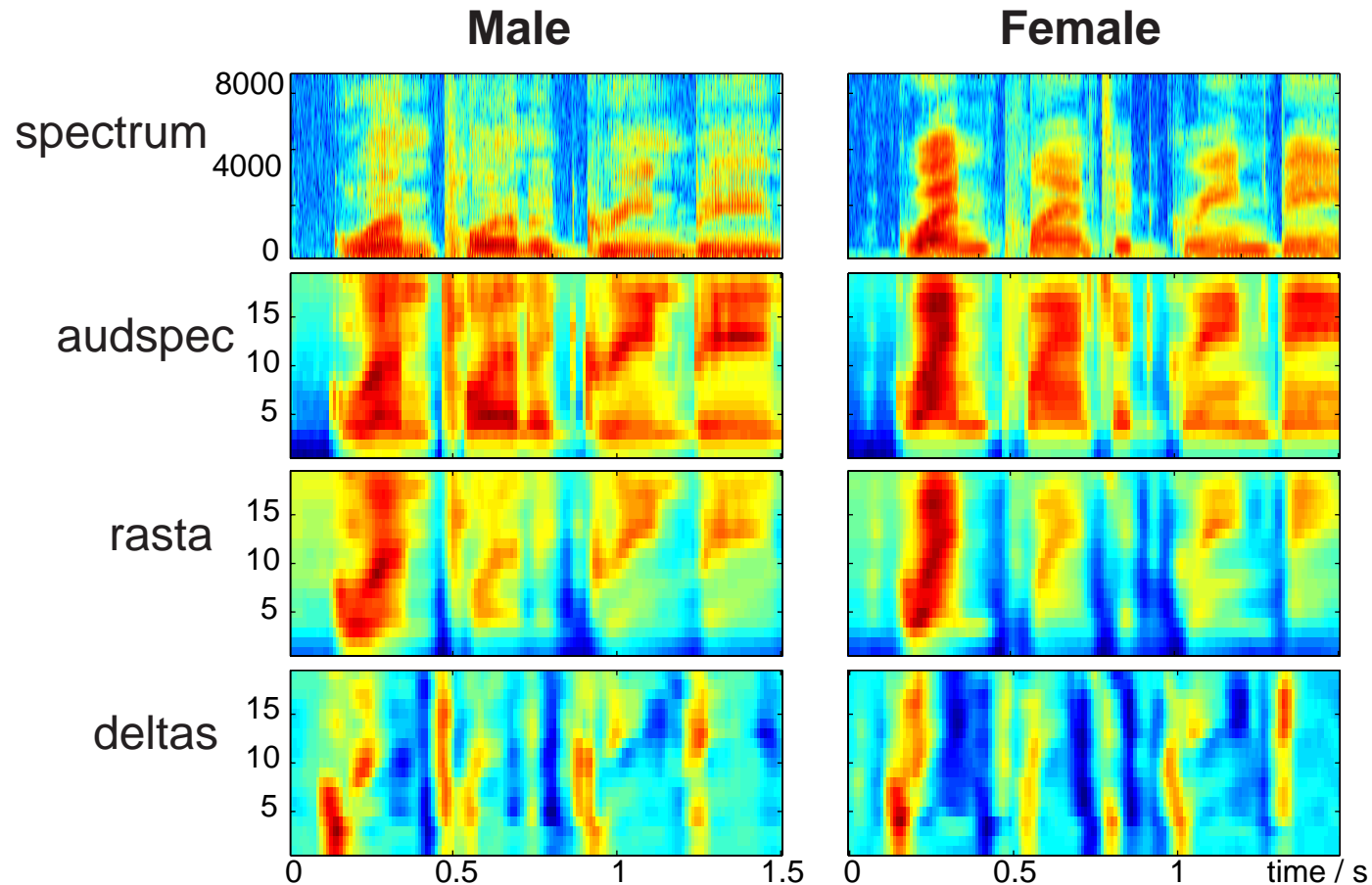  - spectral, auditory scale
  - decorrelation
  - smoothed (spectral) detail
  - normalization of levels

# Features summary
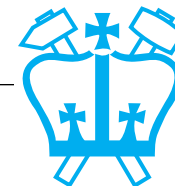
**Male**                    **Female**



- **Normalize same phones**

- **Contrast different phones**

# Summary

- **Speech recognition as word transcription**
  - neat definition, but limited
  - hard because of variability

- **Feature calculation extracts information**
  - smoothed, decorrelated spectral parameters
  - long evolution to match classifiers

**How to actually recognize feature sequences?**