

ACOUSTIC  
PROCESSING  
FOR  
DUMMIES



# Machine recognition of speech

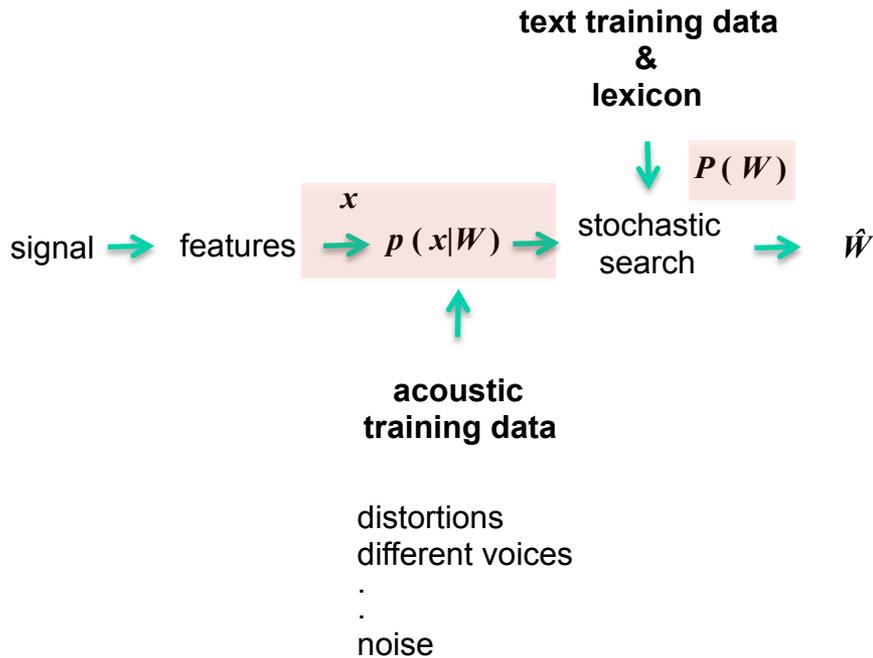
$$w = \underset{i}{\operatorname{argmax}} (P(M(w_i)|x))$$

The “best” model found through Bayes rule

$$w \propto \underset{i}{\operatorname{argmax}} (p(x | M(w_i))P(M(w_i)))$$

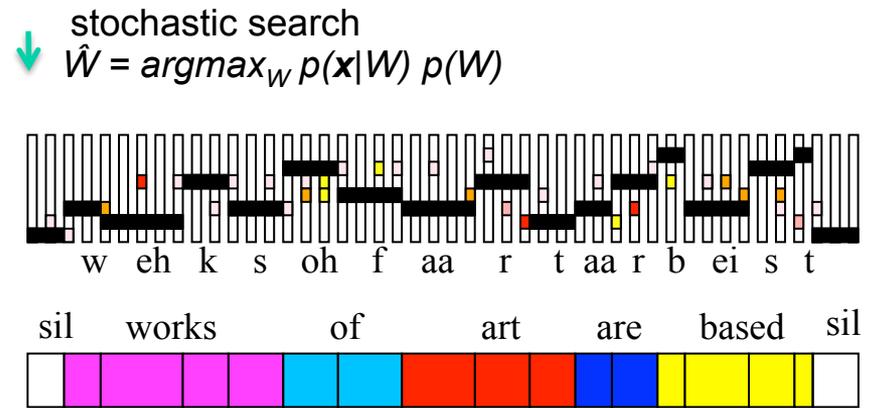
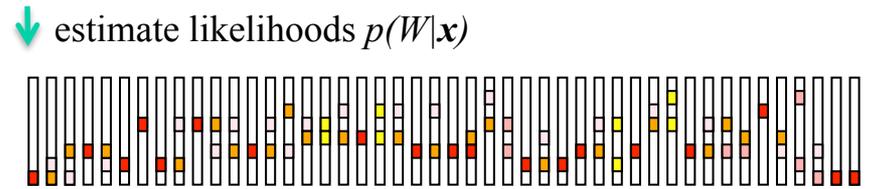
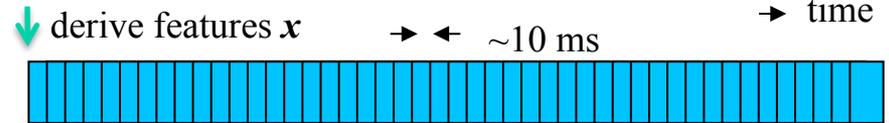
1. How to find  $w$  ?
2. What is the model  $M(w_j)$  ?
- 3. What is the data  $x$  ?**

# Machine Recognition of Speech



- Info lost in  $x$ , is lost forever
- Info left must be dealt with later

speech signal (message, speaker, environment,...)



# Data $x$ ?

## Speech signal ?



- Describes changes in acoustic pressure
  - original purpose is reconstruction of speech
  - rather high bit-rate
- additional processing is necessary to alleviate the irrelevant information
- besides information lost and retained, additional requirements on  $x$  may exist (Normal distributions, de-correlated,..)



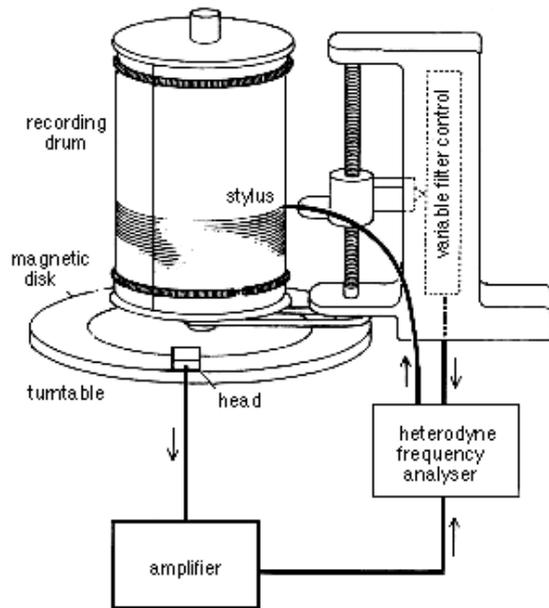
time



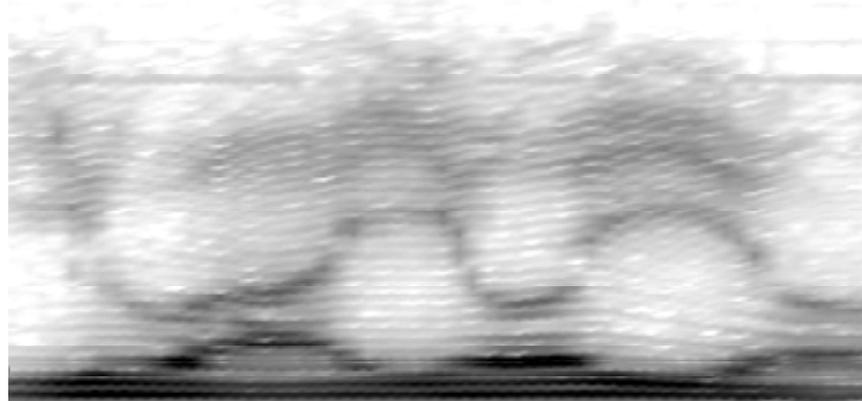
frequency



time

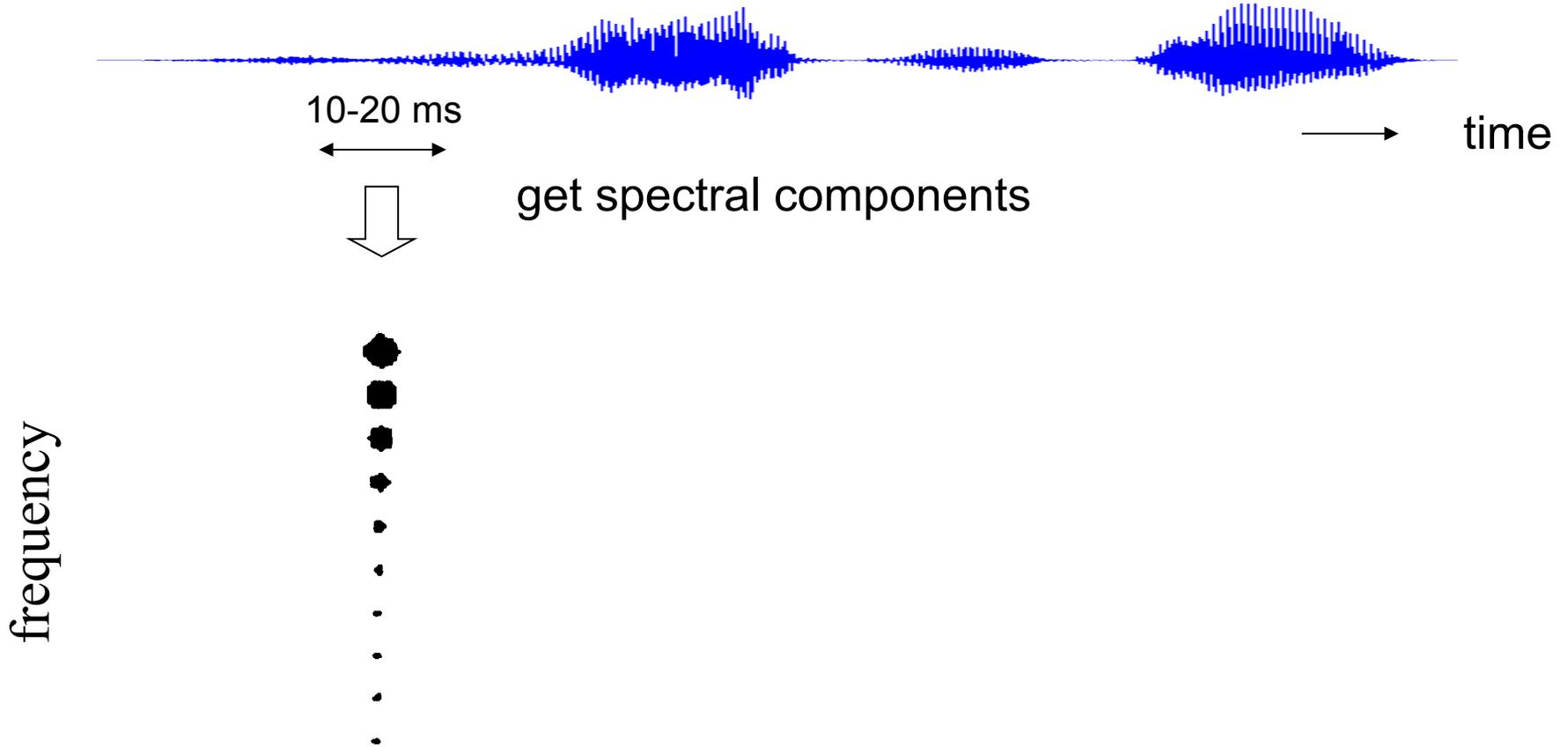


frequency

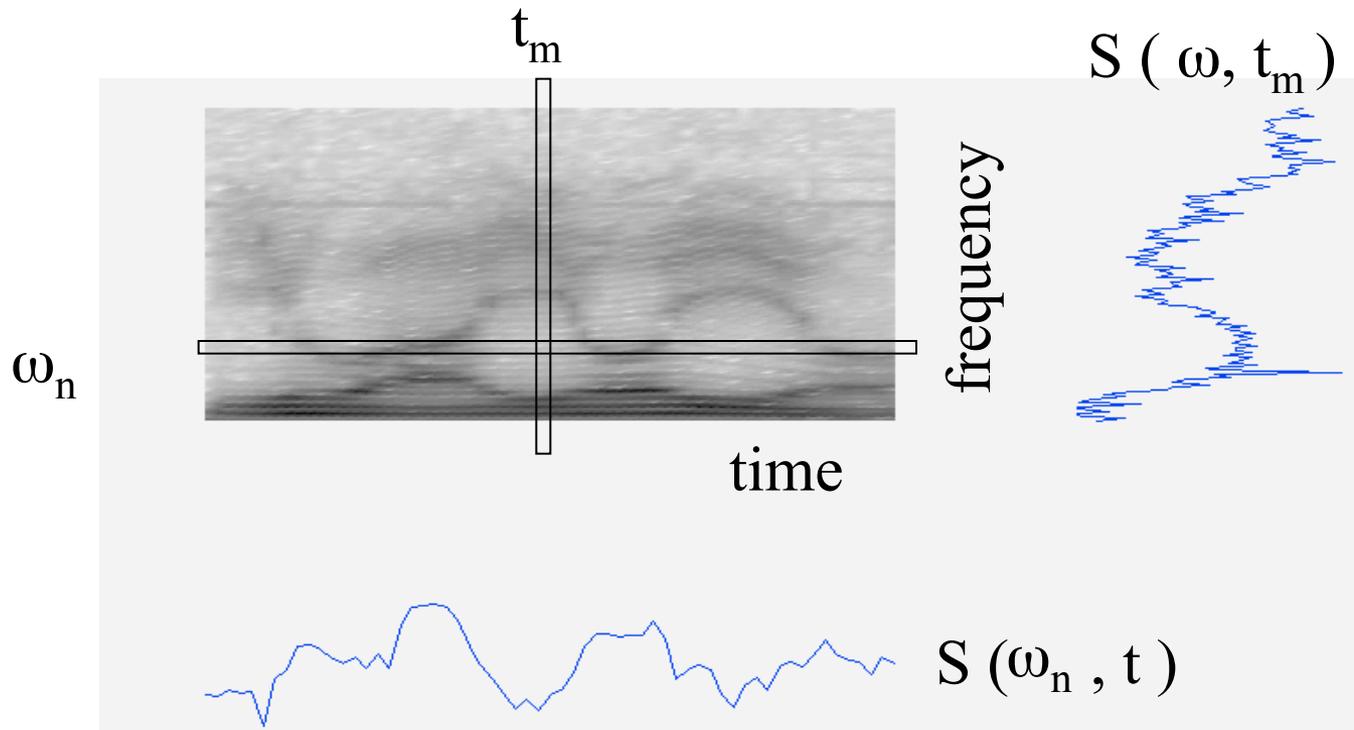


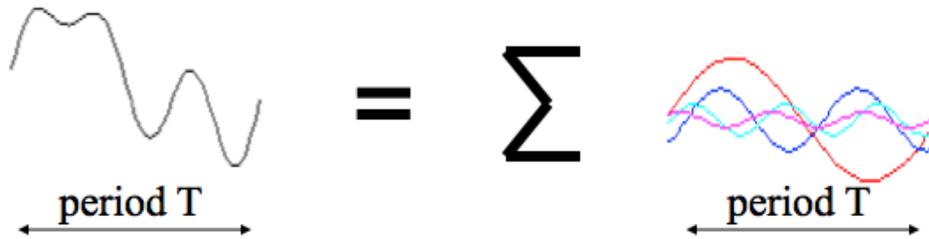
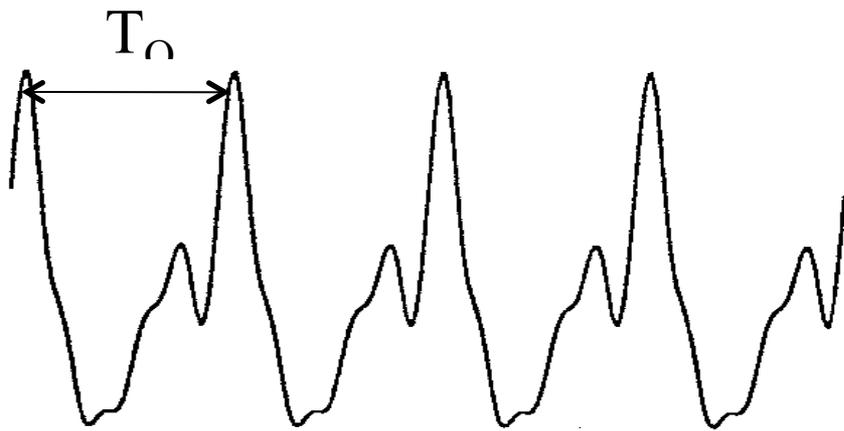
time

# Short-term spectrum



# Spectrogram





Joseph Fourier  
(1768-1830)

Student of Lagrange

Adviser of e.g. Dirichlet or Navier

- One of Fourier ideas

- Describe a periodic signal by an (infinite) sum of other well defined periodic signals (sines and cosines)

# Orthogonality

$$\int_0^T \sin n\omega t \cdot \cos m\omega t dt = 0$$

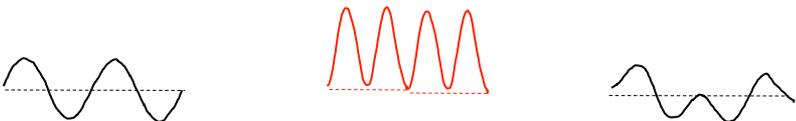
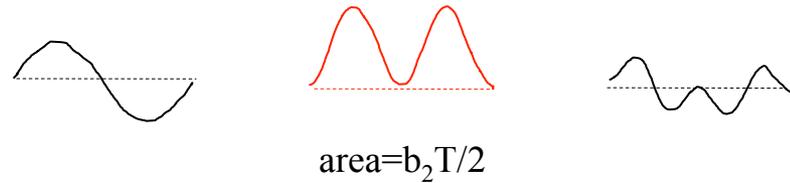
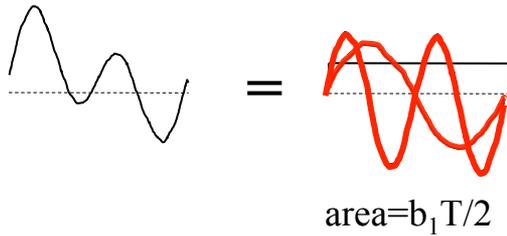
$$\int_0^T \cos n\omega t \cdot \cos m\omega t dt = 0 \text{ for } n \neq m \text{ and } \frac{T}{2} \text{ for } n = m$$

$$\int_0^T \sin n\omega t \cdot \sin m\omega t dt = 0 \text{ for } n \neq m \text{ and } \frac{T}{2} \text{ for } n = m$$

$$f(t) = DC + \sum_{i=1}^{\infty} \left[ a_i \cos\left(\frac{2\pi i t}{T}\right) + b_i \sin\left(\frac{2\pi i t}{T}\right) \right] = DC + a_1 \cos\left(\frac{2\pi t}{T}\right) + b_1 \sin\left(\frac{2\pi t}{T}\right) + a_2 \cos\left(\frac{4\pi t}{T}\right) + b_2 \sin\left(\frac{4\pi t}{T}\right) + a_3 \cos\left(\frac{6\pi t}{T}\right) + b_3 \sin\left(\frac{6\pi t}{T}\right) + \dots$$

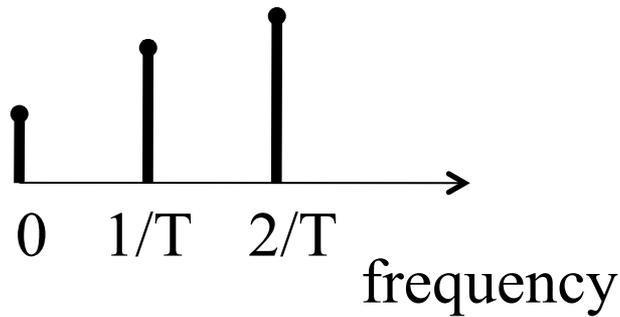
$$\int_0^T f(t) \sin\left(\frac{2\pi t}{T}\right) dt = \int_0^T \left\{ DC \sin\left(\frac{2\pi t}{T}\right) + a_1 \cos\left(\frac{2\pi t}{T}\right) \sin\left(\frac{2\pi t}{T}\right) + b_1 \sin\left(\frac{2\pi t}{T}\right) \sin\left(\frac{2\pi t}{T}\right) + a_2 \cos\left(\frac{4\pi t}{T}\right) \sin\left(\frac{2\pi t}{T}\right) + b_2 \sin\left(\frac{4\pi t}{T}\right) \sin\left(\frac{2\pi t}{T}\right) + \dots \right\} dt$$

$$0 \quad 0 \quad b_1 T/2 \quad 0 \quad 0 \dots$$

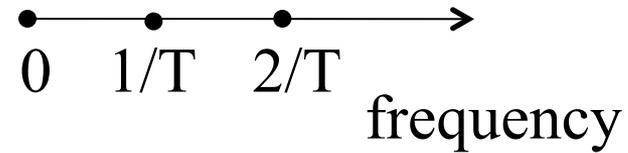


$$\int_0^T \sin^2\left(\frac{t}{T}\right) dt = \frac{T}{2}$$

Magnitude spectrum



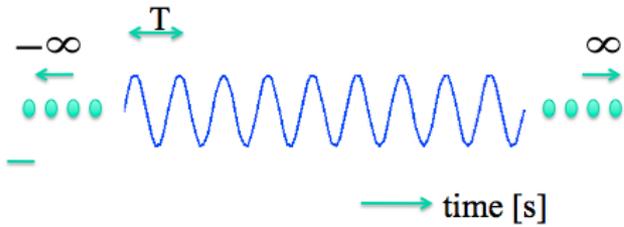
Phase spectrum



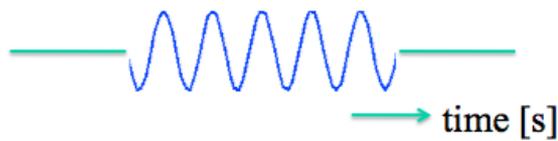
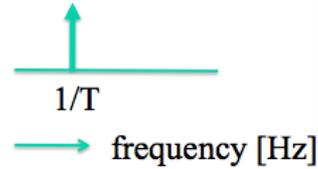
Spacing of spectral components is  $1/T$

Periodicity in one domain (here time) implies discrete representation in the dual domain (here frequency)

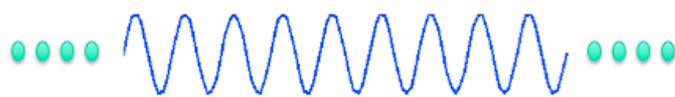
Sinusoidal signal (pure tone)



Its spectrum



Truncated signal

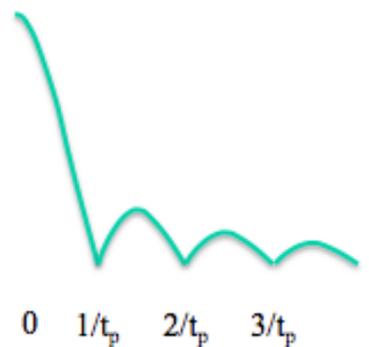
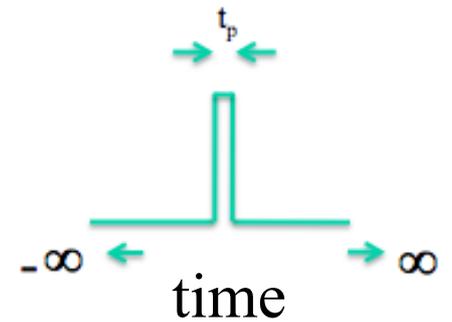


Infinite signal

*multiplied by*



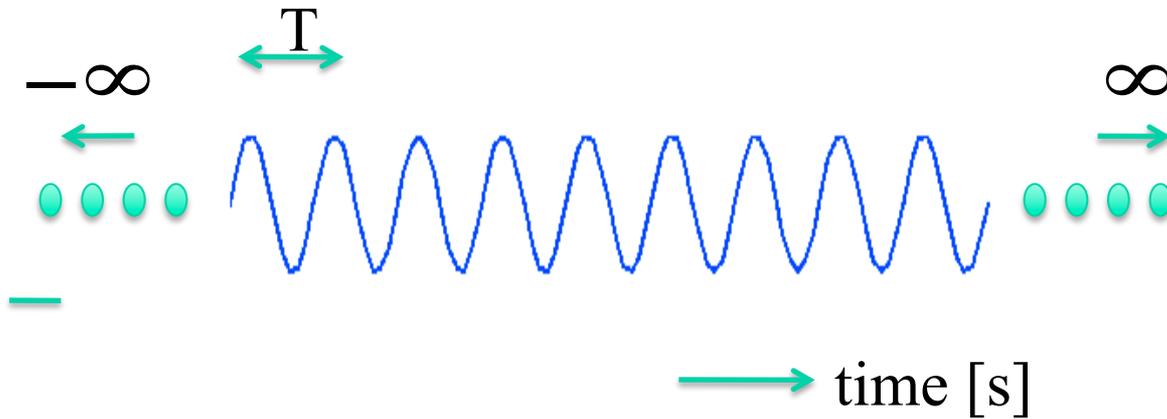
square window



Multiplication in one (time) domain is convolution in the dual (frequency) domain

frequency

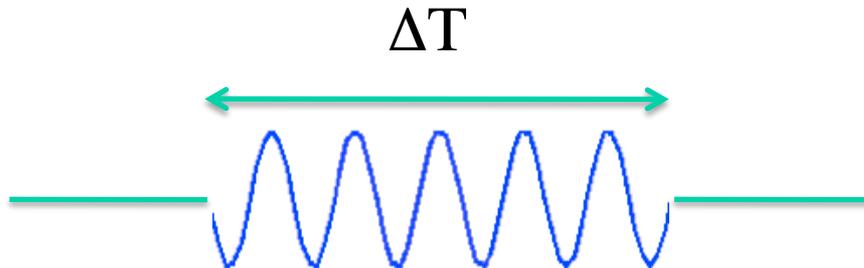
Sinusoidal signal (pure tone)



Its spectrum



Truncated sinusoidal signal

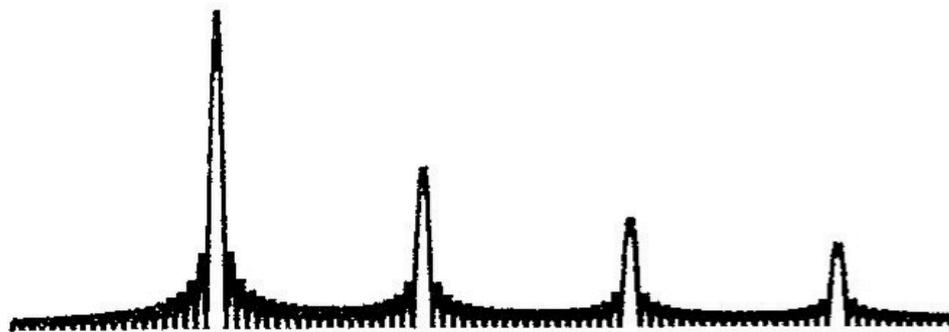


Its spectrum

?



$$\Delta t = \infty$$



$$\Delta t = 100 \text{ ms}$$

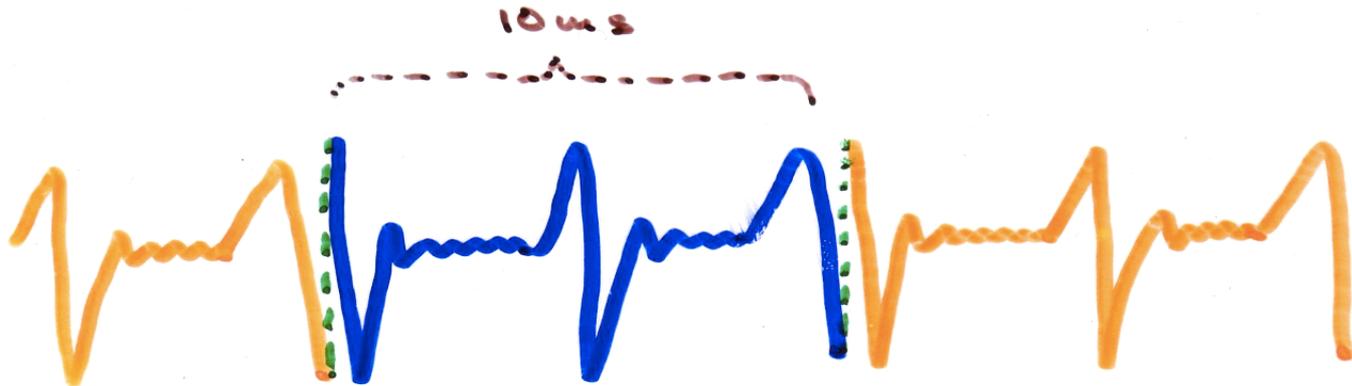


$$\Delta t = 13 \text{ ms}$$

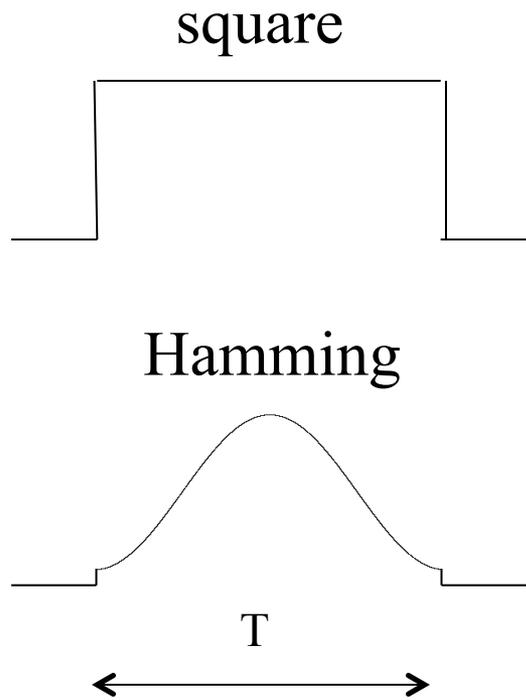
0

850 Hz

Non-stationary turns into periodic

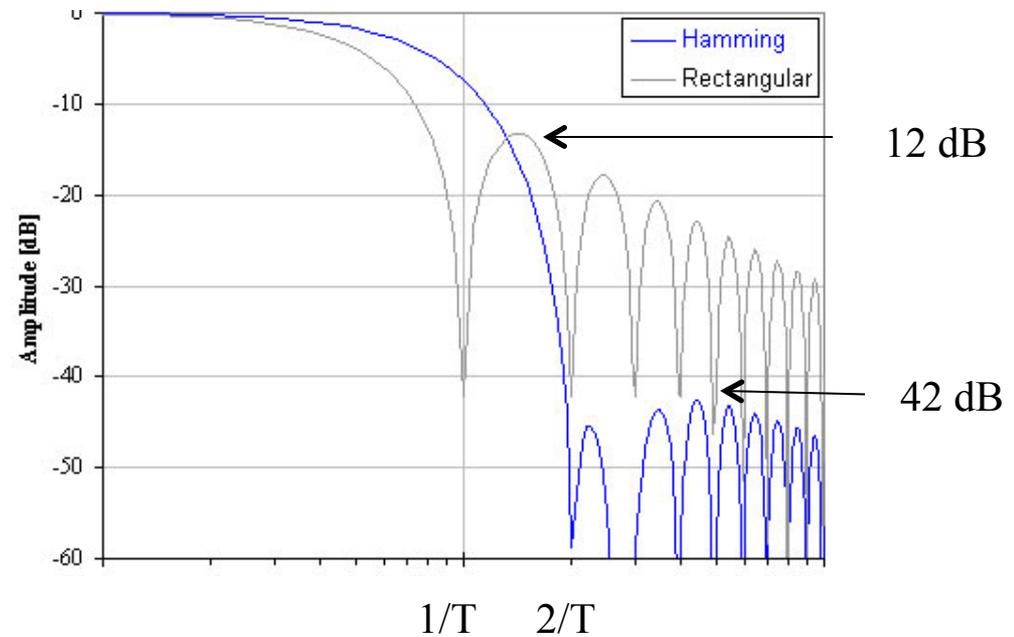


## Time domain

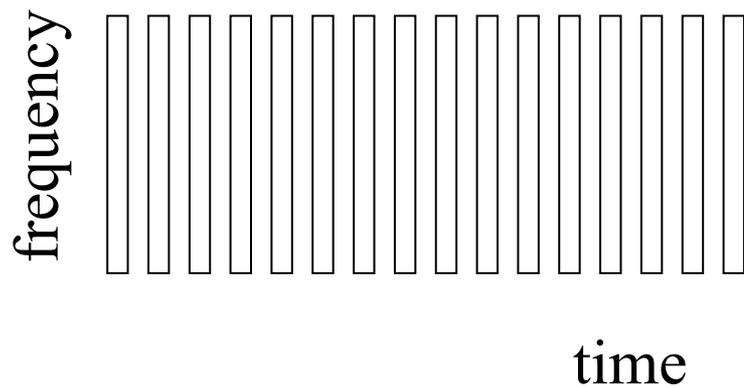


**Multiplication**  
with signal

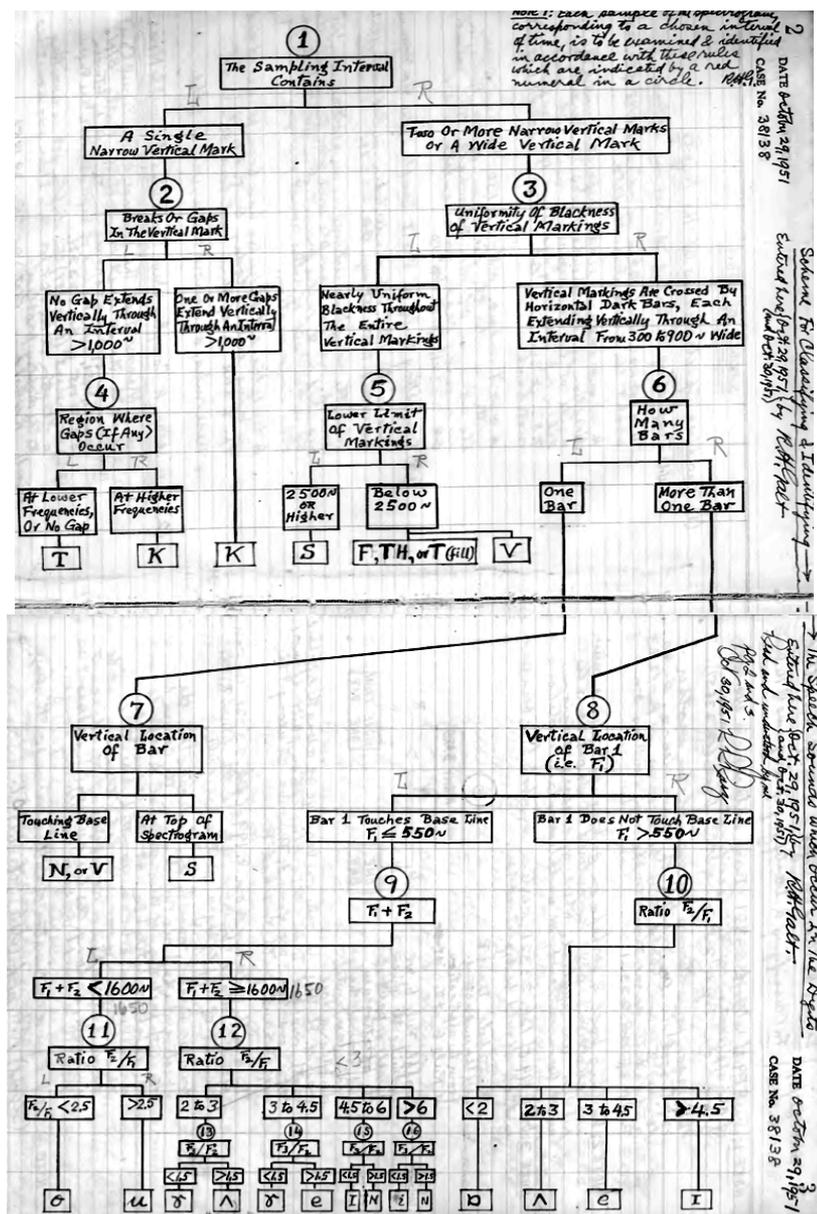
## Frequency domain



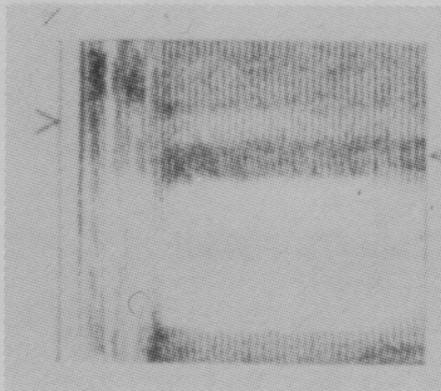
**Convolution**  
with signal spectrum



# Concept of the first "real" automatic speech recognizer (R.H. Galt 1951)

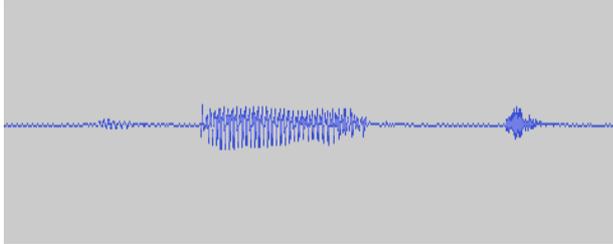


/k/

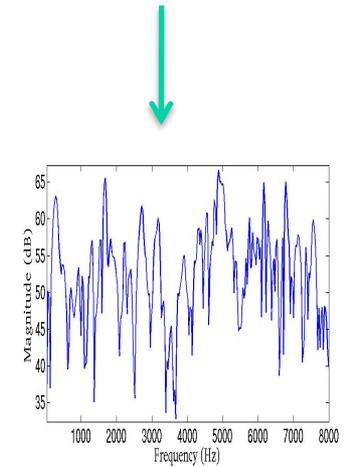
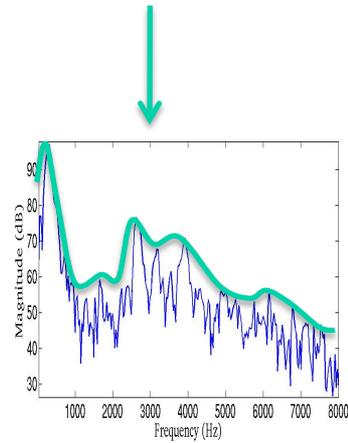
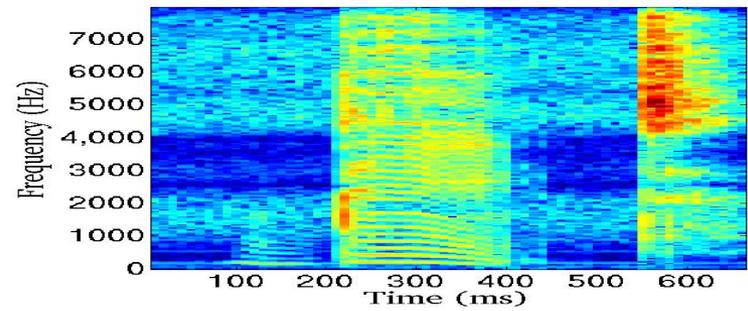
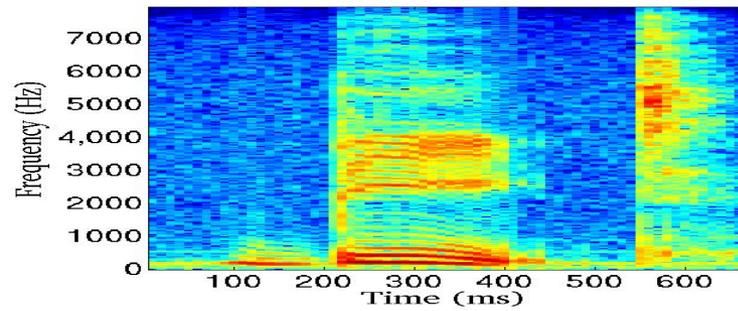


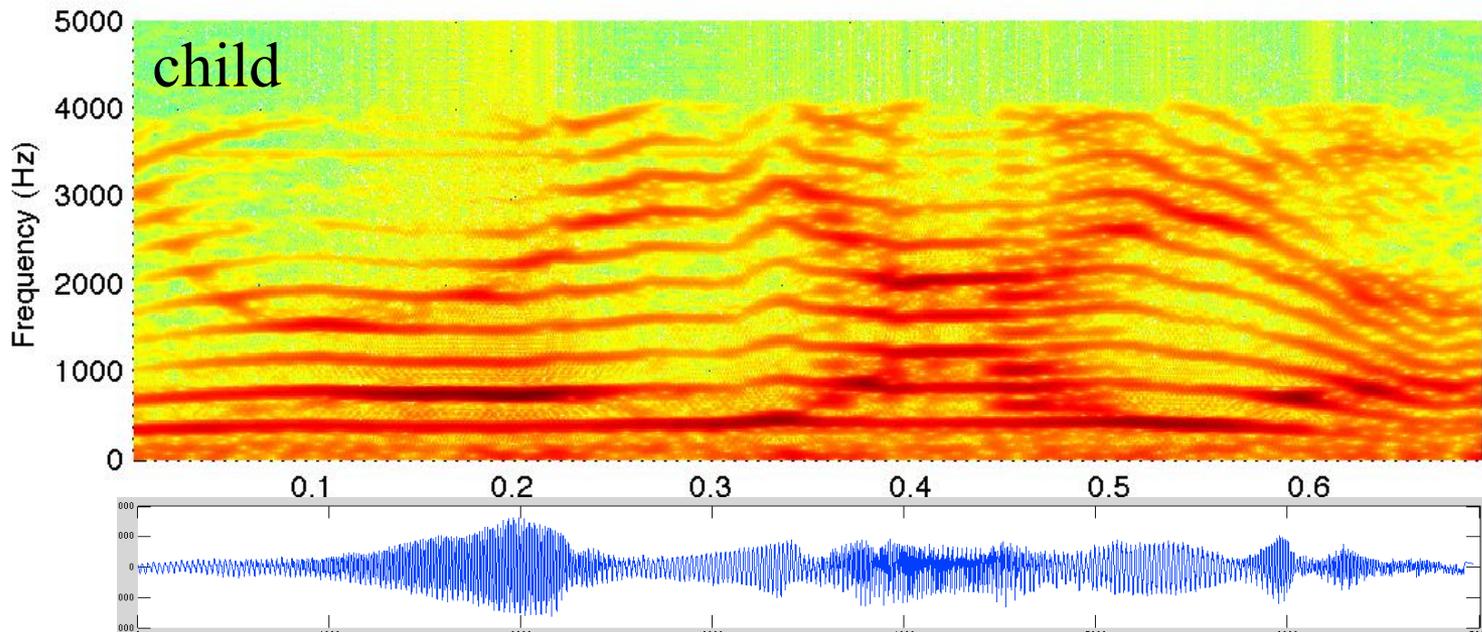
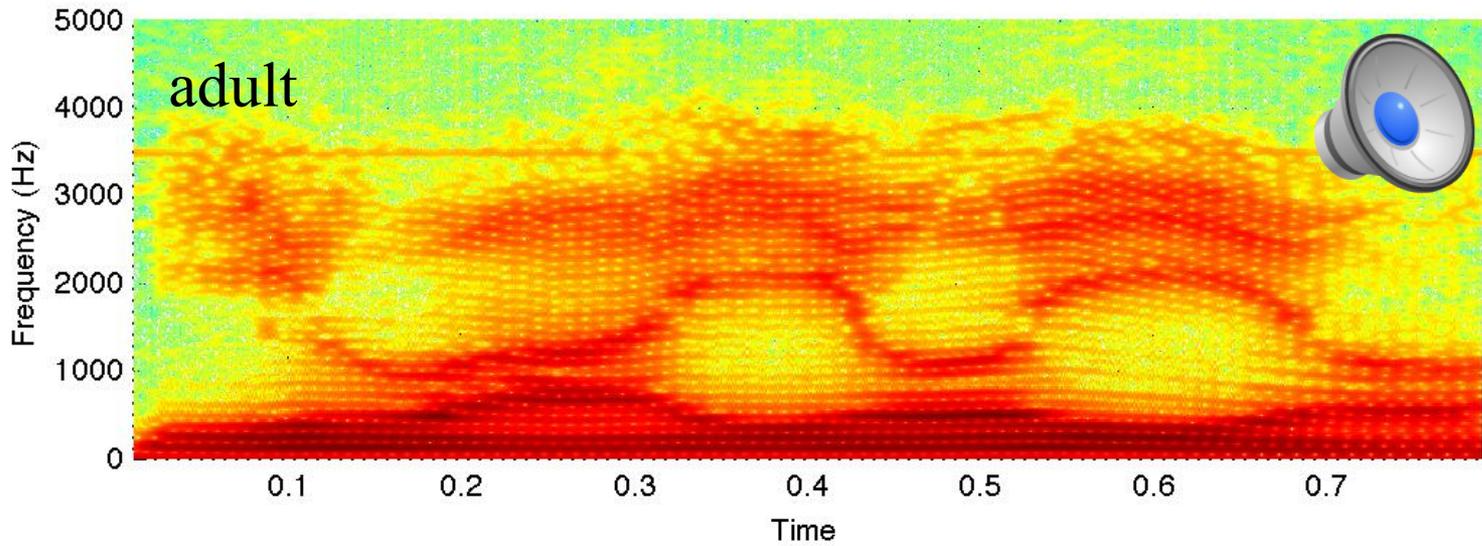
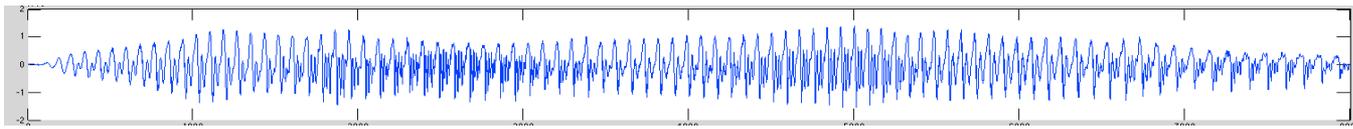
ki (*key*)

Potter, Kopp, and Green, Visible Speech 1947



filter

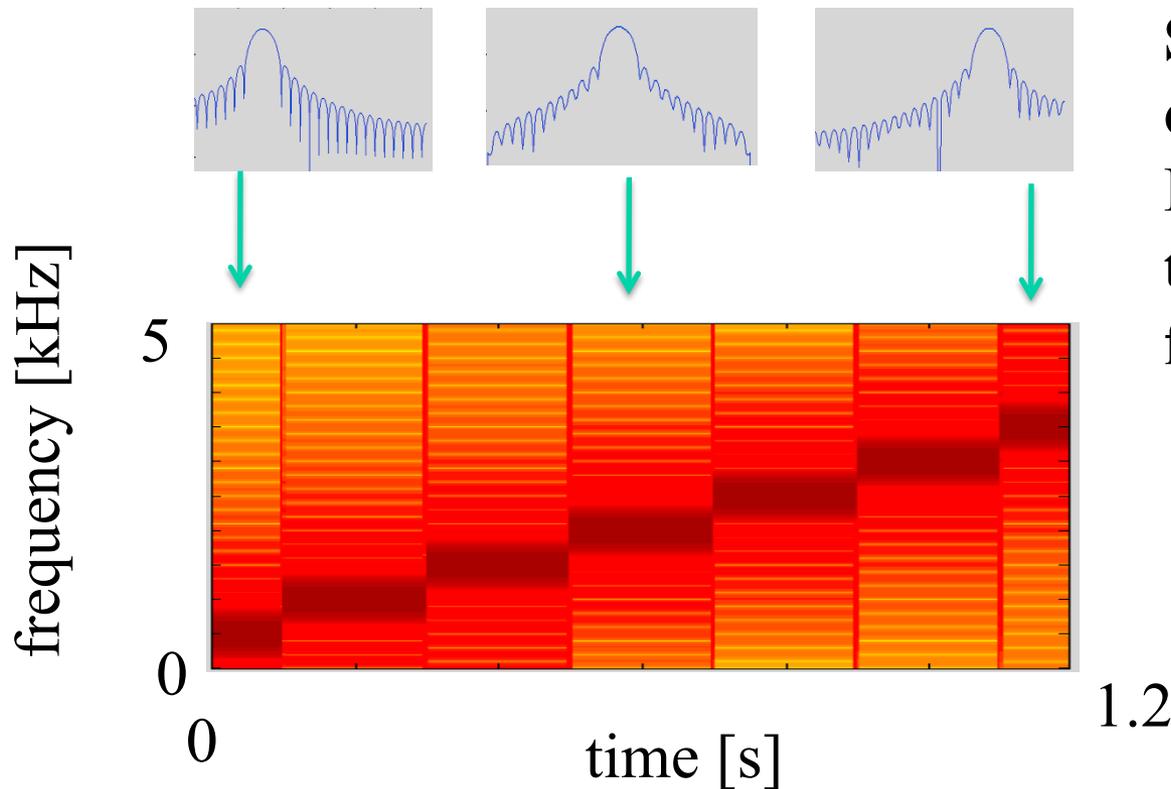






$$S_n(e^{j\omega}) = \sum_{m=-\infty}^{\infty} s(m) \cdot w(n-m)e^{-jm\omega}$$

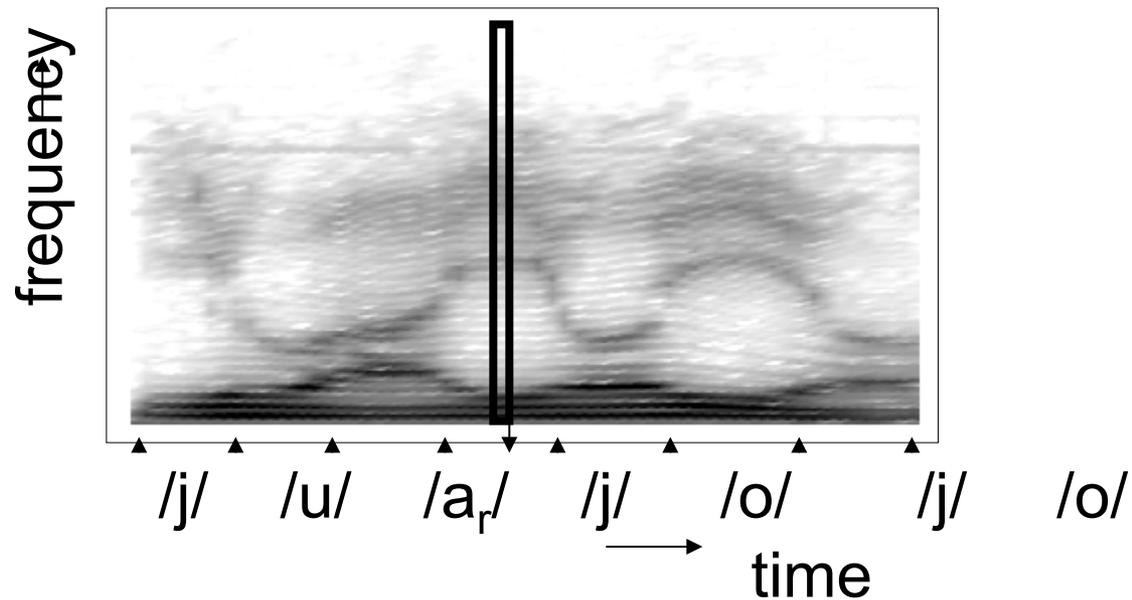
Fourier transform of the signal  $s(m)$  multiplied by the window  $w(n-m)$   
Spectrum is the line spectrum of the signal convolved with the spectrum of the window



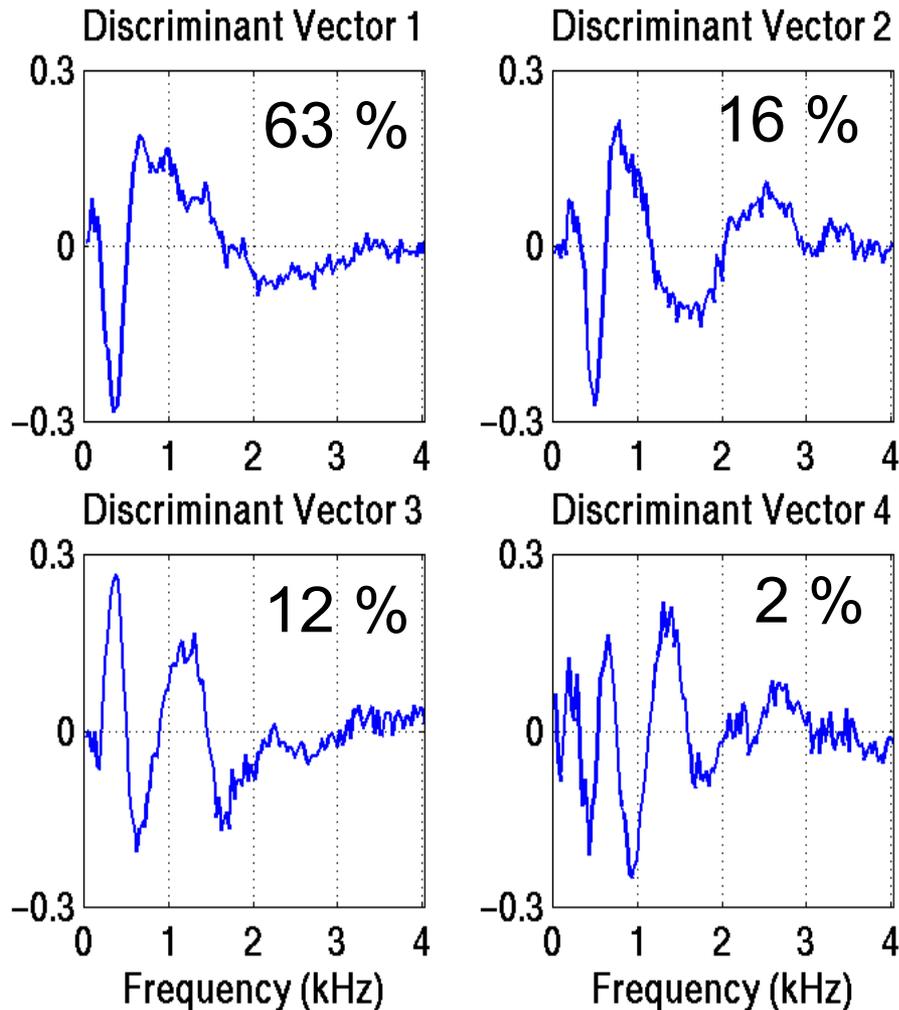
**Spectral resolution  
of the short-term  
Fourier analysis is  
the same at all  
frequencies.**

# Spectral Basis from LDA

LDA gives basis for projection of spectral space



# LDA vectors from Fourier Spectrum (OGI 3 hour stories hand-labeled database)



- Spectral resolution of LDA-derived spectral basis is higher at low frequencies

Psychophysics:

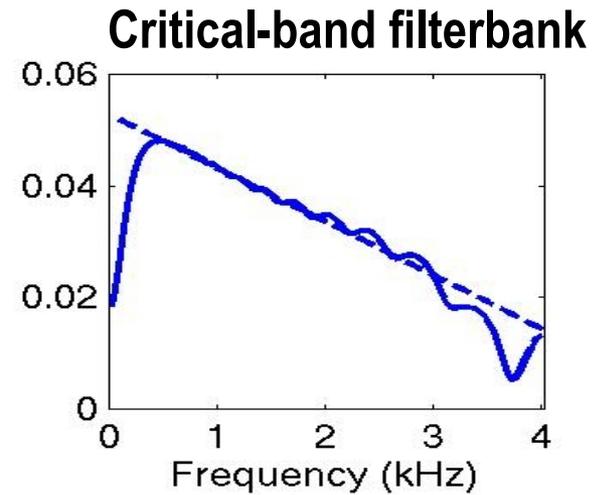
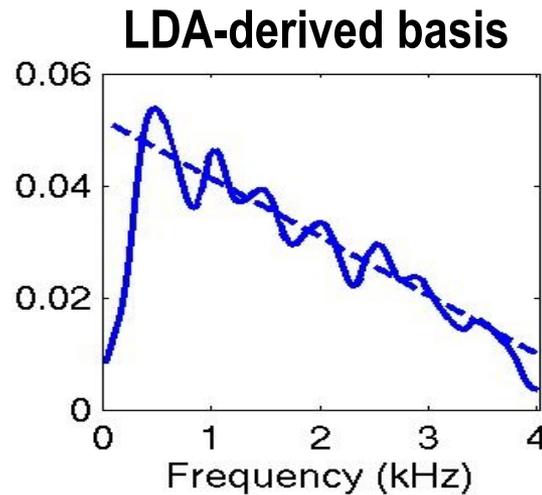
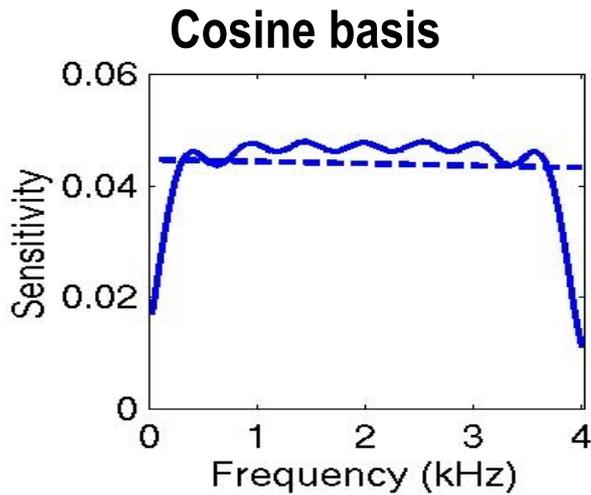
Critical bands of human hearing are broader at higher frequencies

Physiology:

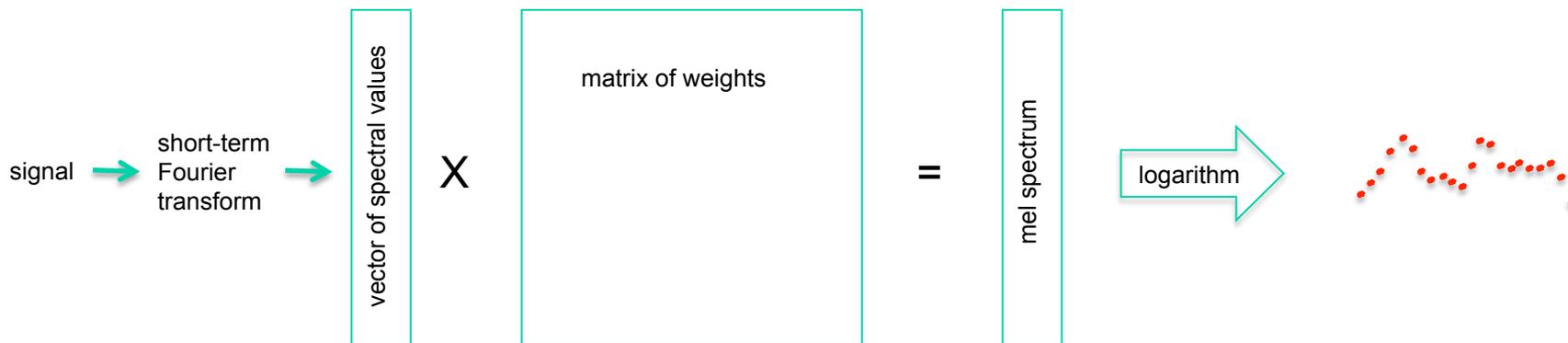
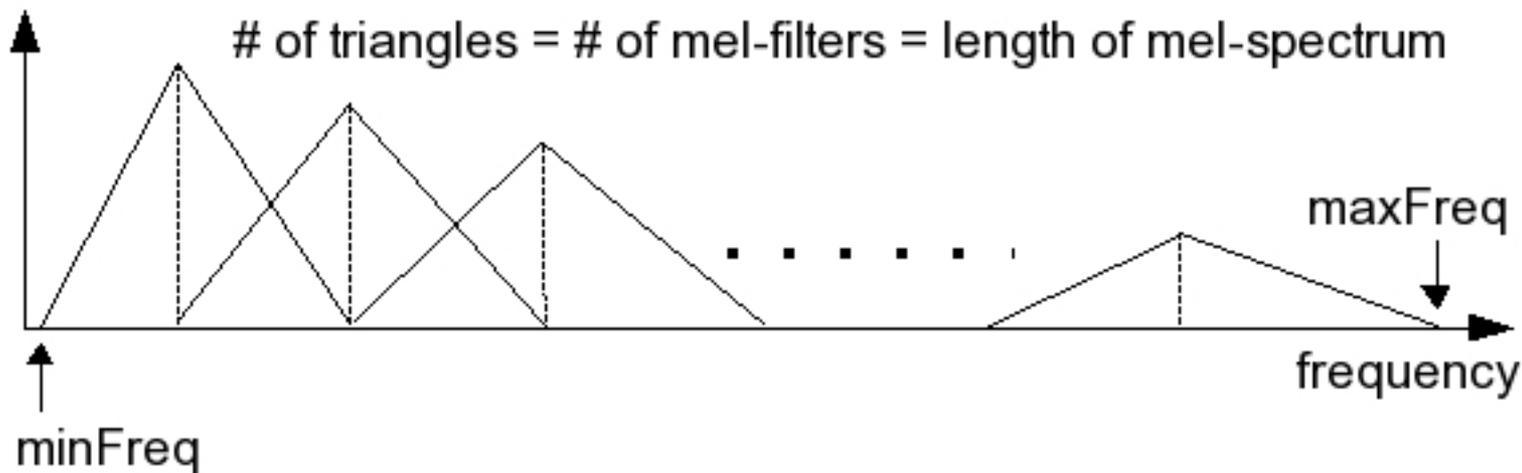
Position of maximum of traveling wave on basilar membrane is proportional to logarithm of frequency

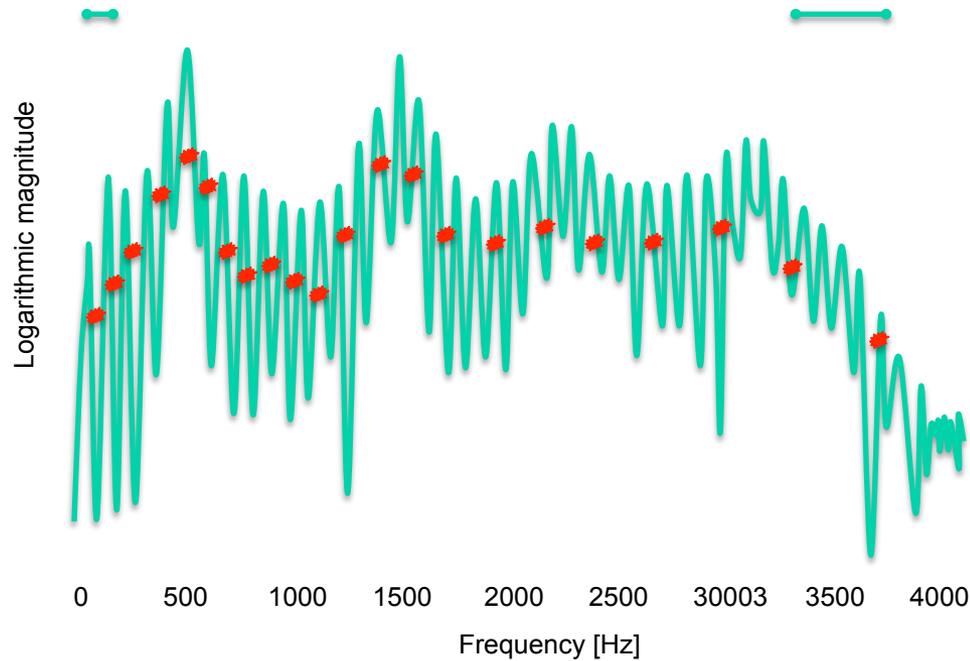
# Sensitivity to Spectral Change

(Malayath 1999)

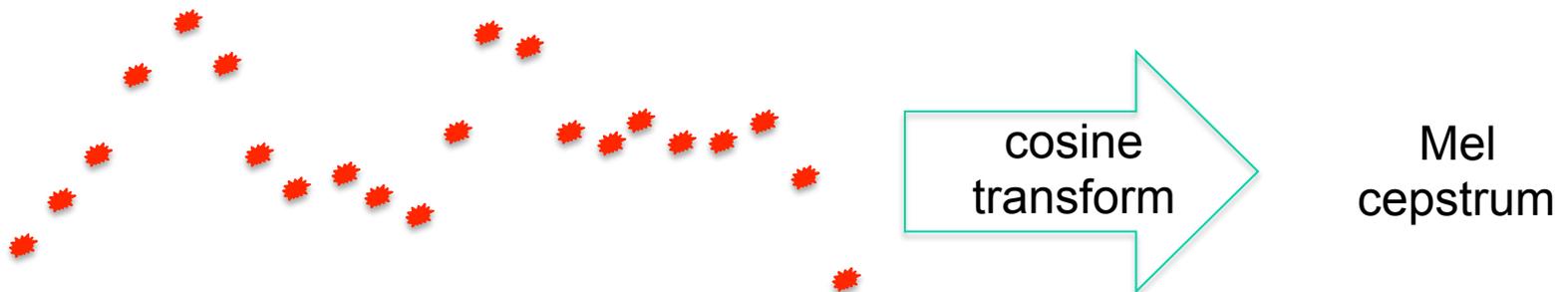


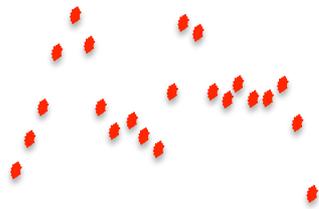
# Spectral weights



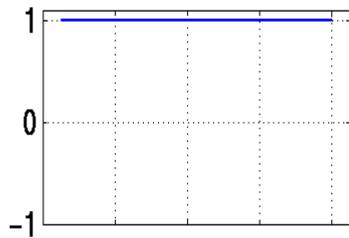


- Spectral resolution decreases with frequency.
- Temporal resolution stays the same (given by the length of the analysis window in computing spectrum)

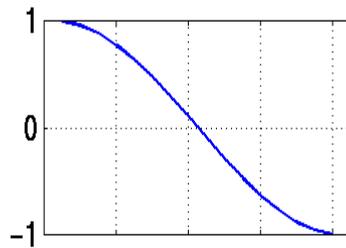




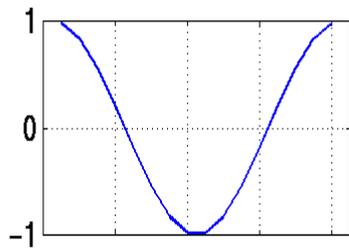
$C_0$



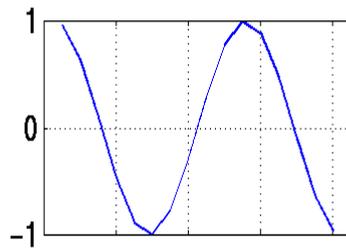
$C_1$



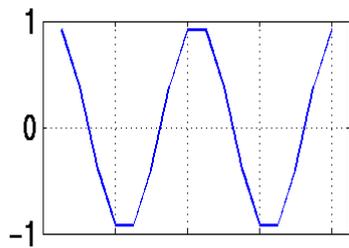
$C_2$



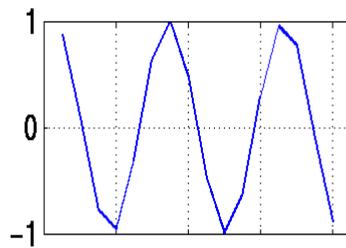
$C_3$



$C_4$



$C_5$

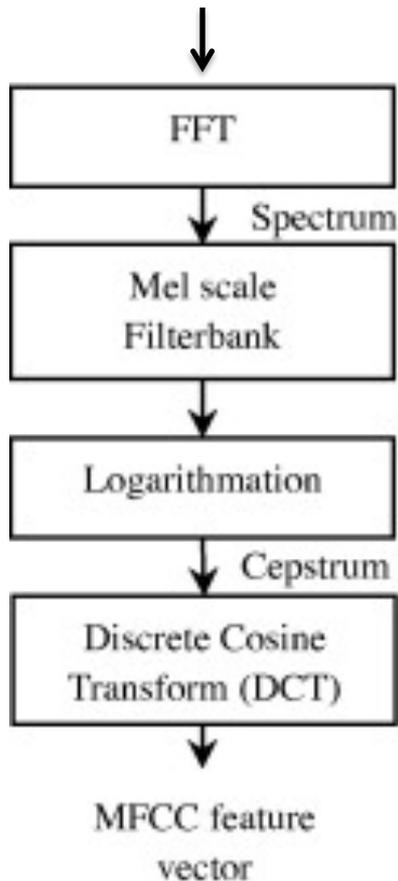


Critical Bands

Critical Bands

# Mel cepstrum

Segment of signal (~ 20 ms – windowed)



Short-term spectrum

- Frequency selectivity of hearing

Project on spectral weights

- Non-equal spectral resolution of hearing

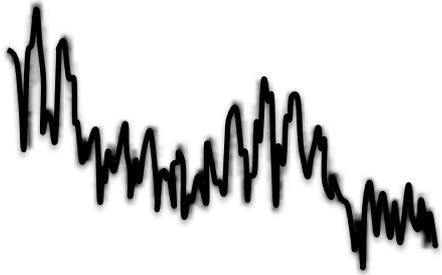
Take logarithm

- make distribution more Normal

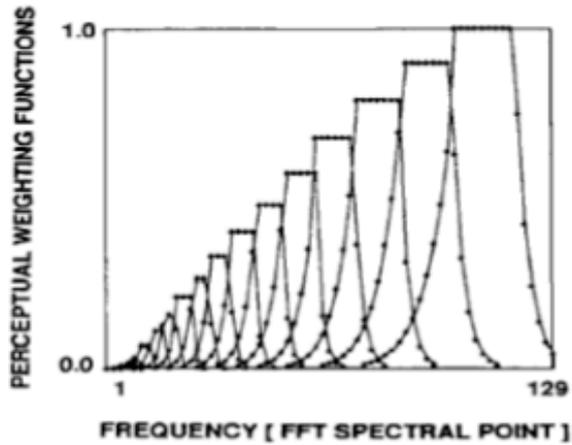
Cosine transform

- de-correlate

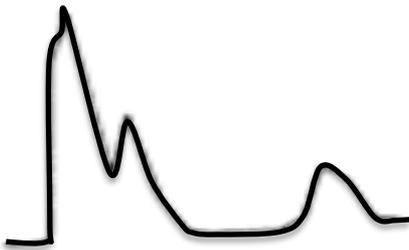
# Perceptual Linear Prediction



spectrum

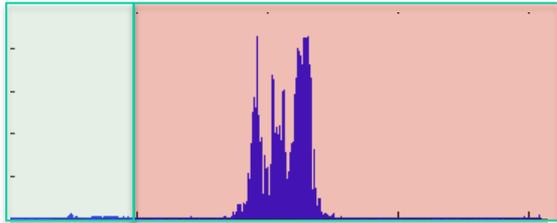


summation windows

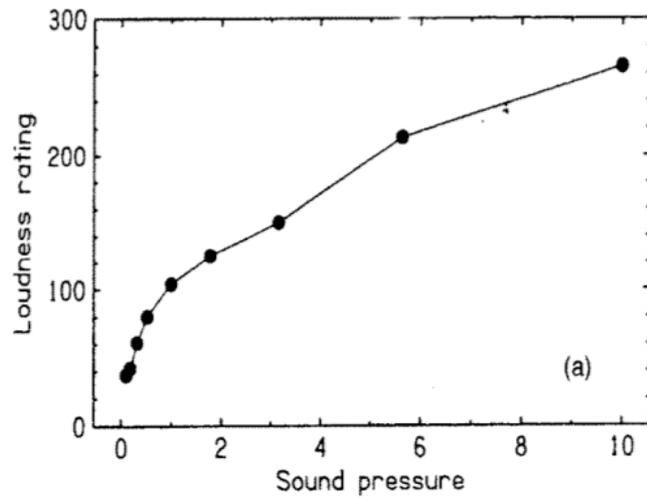
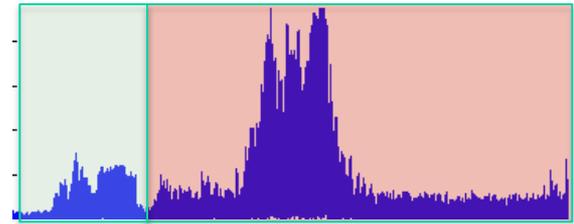


spectrum with auditory-like resolution

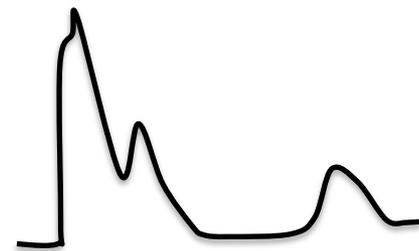
intensity  $\approx$  signal<sup>2</sup> [w/m<sup>2</sup>]



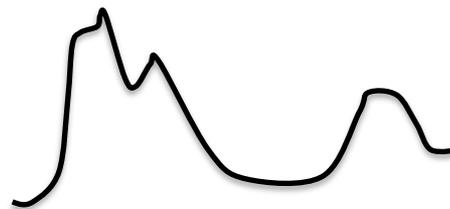
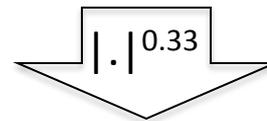
loudness [Sones]



loudness = intensity<sup>0.33</sup>



intensity  
(power spectrum)

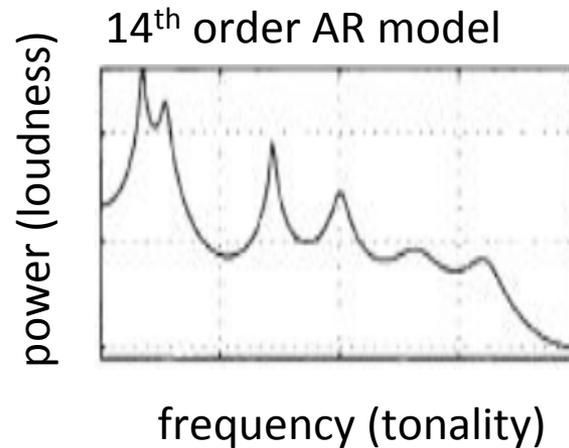
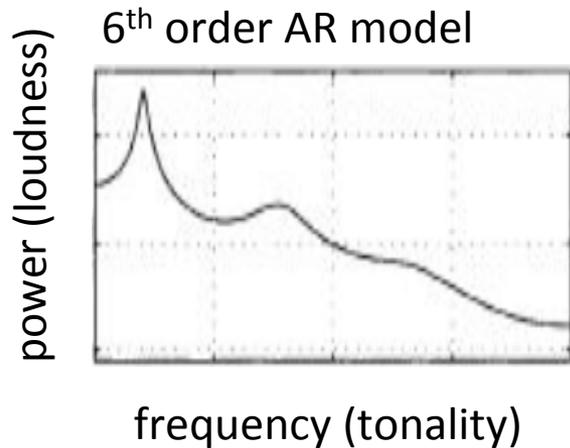


loudness

# Not all spectral details are important

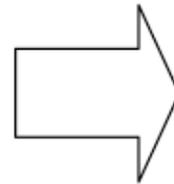
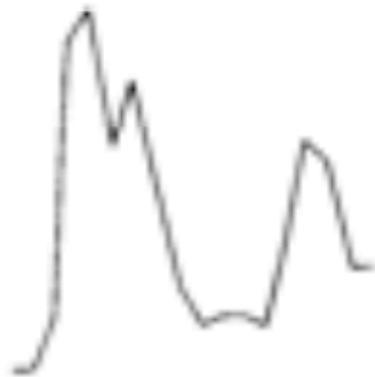
a) compute Fourier transform of the auditory spectrum and truncate it (cepstrum)

b) approximate the auditory spectrum by an autoregressive model

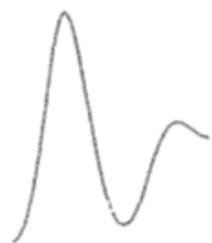
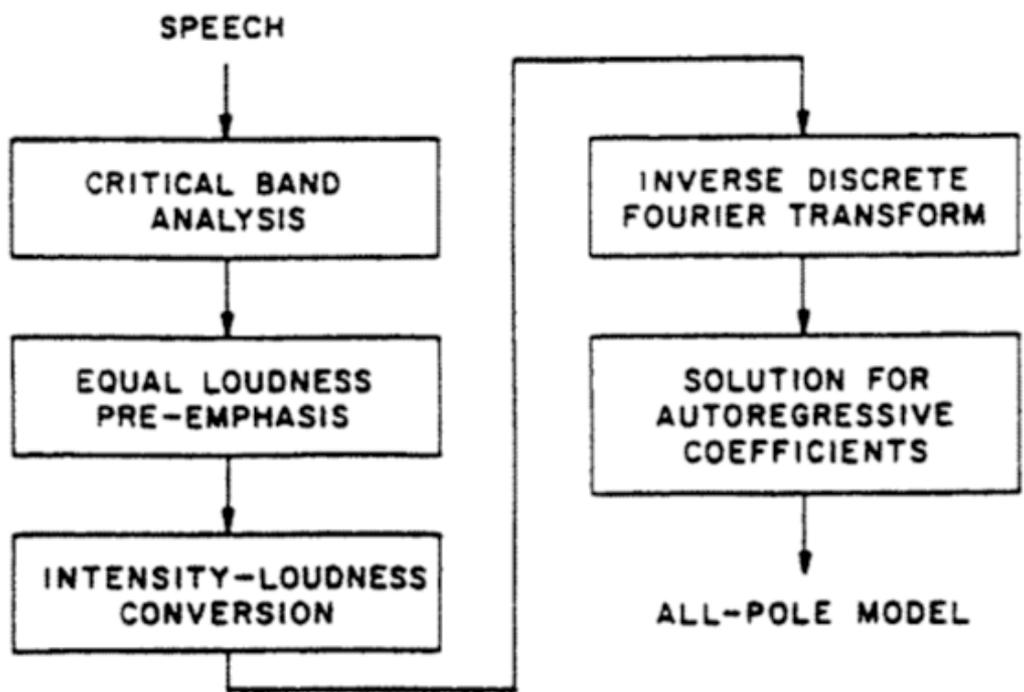
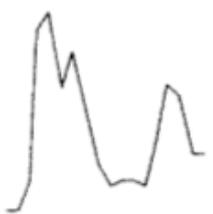
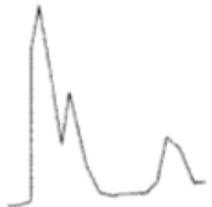


# Perceptual Linear Prediction (PLP) Autoregressive fit to the auditory-like spectrum

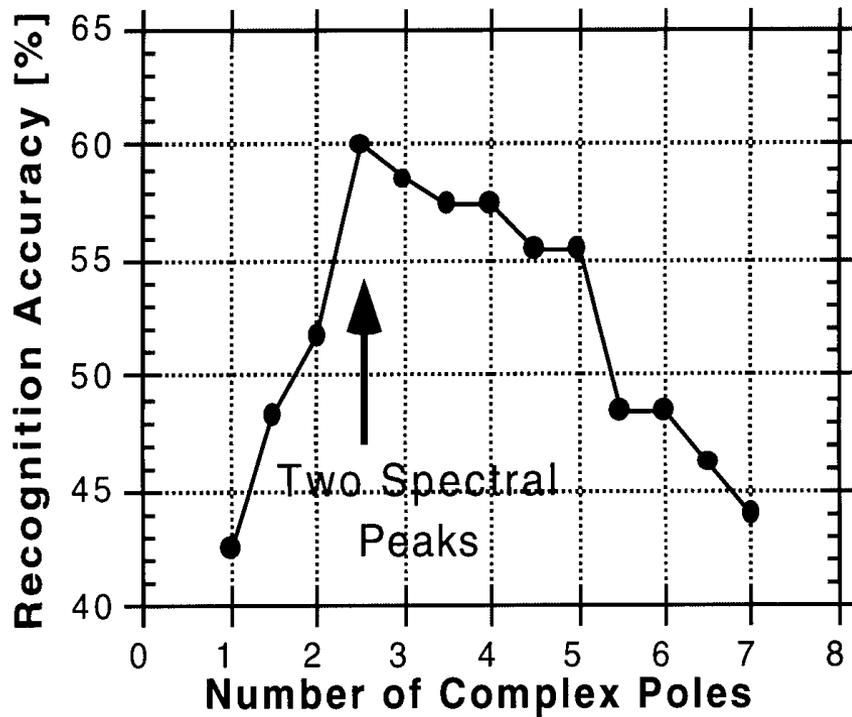
power  
(loudness)



frequency (tonality)



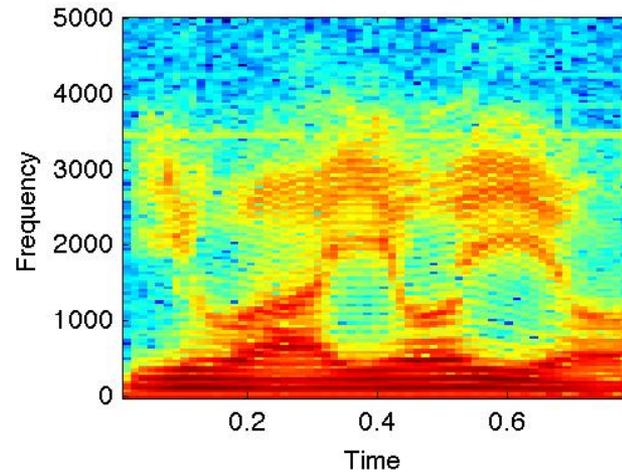
# Optimal Amount of Spectral Smoothing (order of PLP autoregressive model)



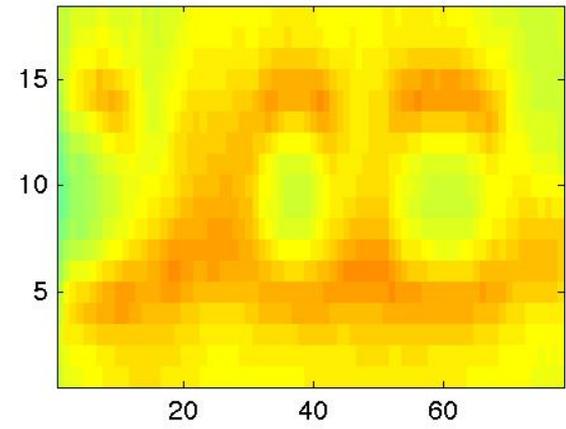
- cross-speaker ASR (trained on one speaker and tested on another)
- all speaker-dependent information harmful

adult male

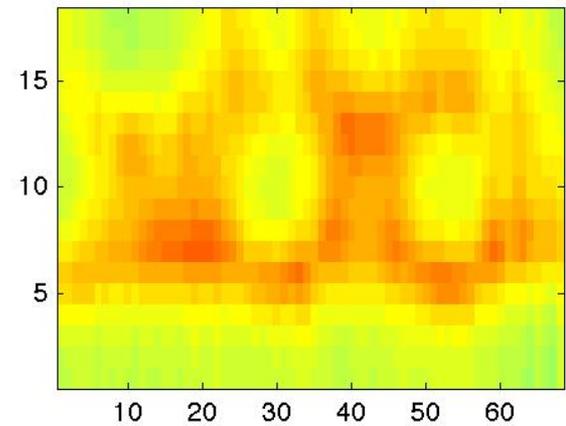
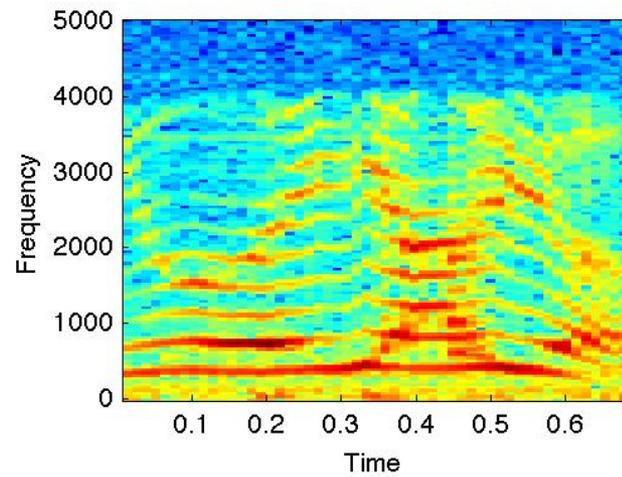
short-term spectrum

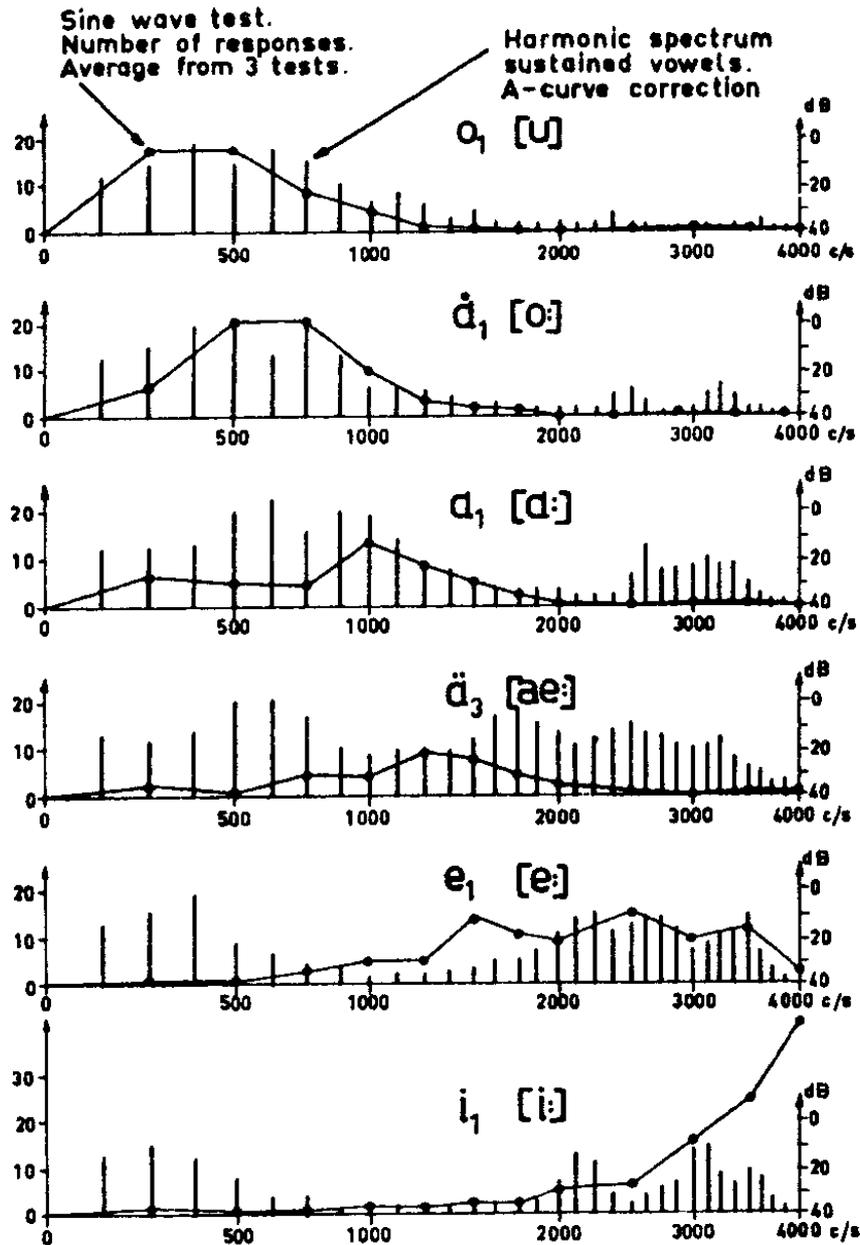


5<sup>th</sup> order PLP spectrum



4 year old child

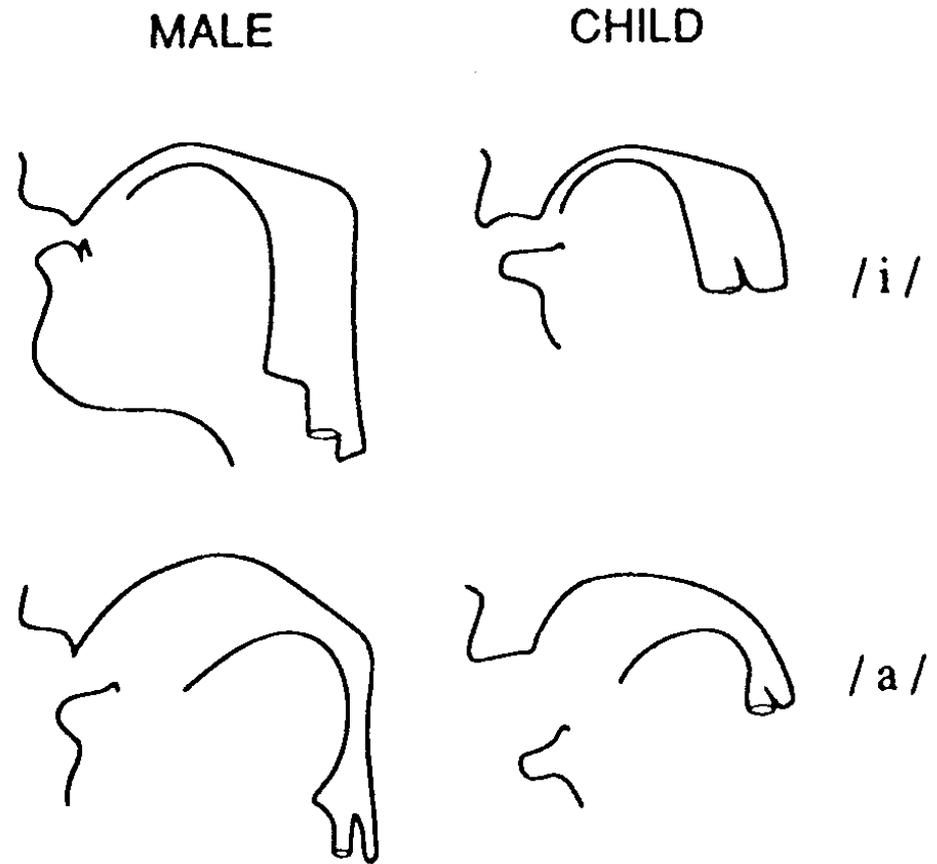




- affiliate vowel with sine wave tone (forced judgment)
- peak of histograms would correspond to resonance frequency of uncoupled front cavity in production of a given vowel
  - Fant 1947
- Perceptual F2'
  - position of second peak in two-peak simulation of vowels

# X-rays of Male and Child Vocal Tract in Production of Vowels

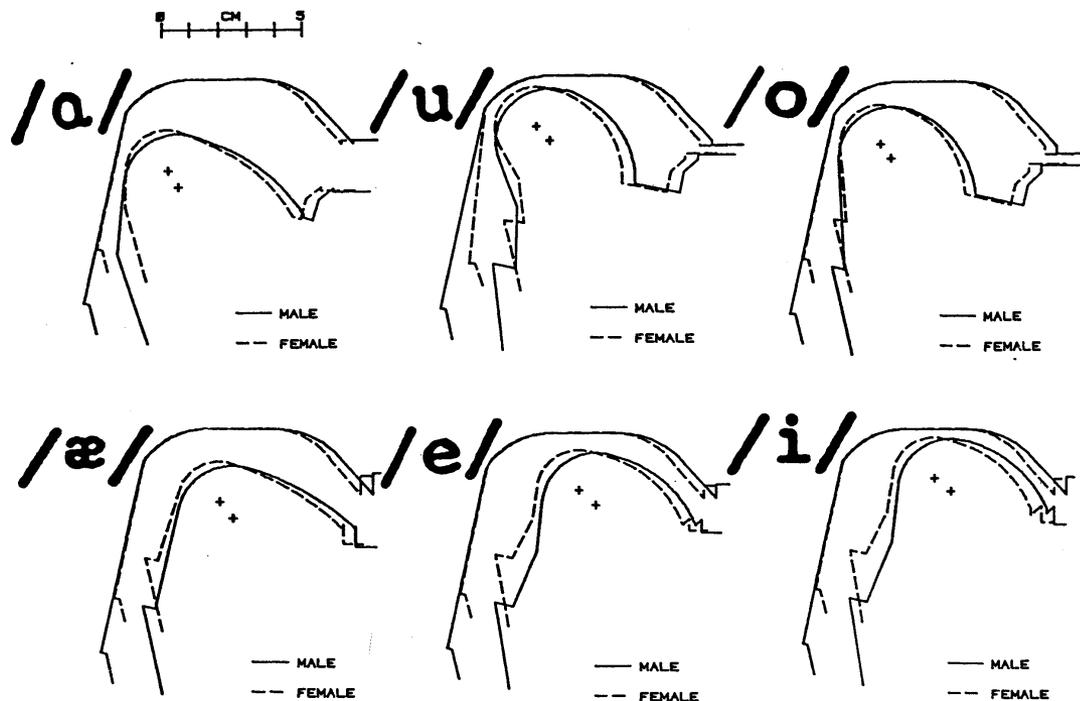
- In production of vowels, the front part of the vocal tract appears to be less speaker dependent than its back part
  - Hermansky and Broad 1990



# Female vocal tract from male

Ursula Goldstein, MIT PhD. Thesis 1980

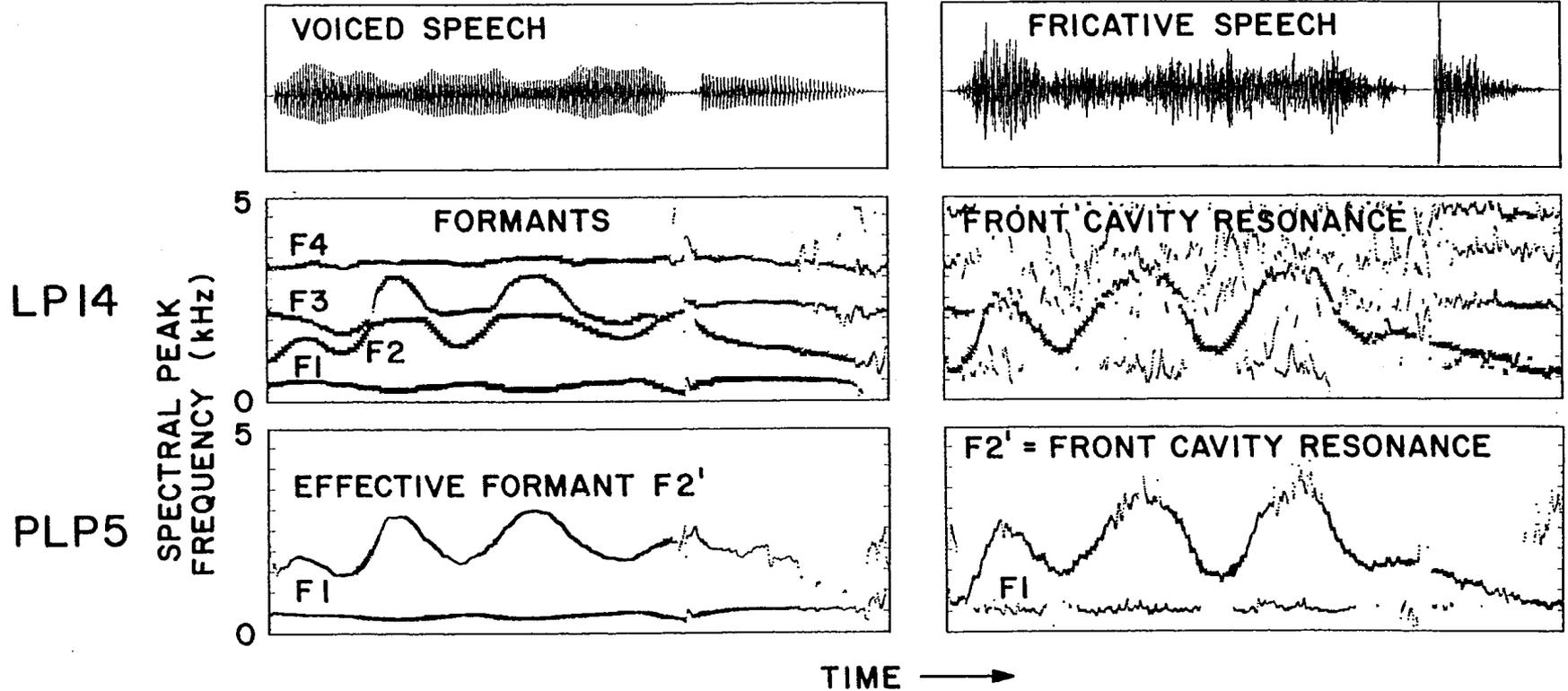
- Start with male vocal tract x-ray
- implement male-female anatomical differences
- change “resting dimensions” to “female”



# Front Cavity - F2' Hypothesis

- F2' correlates with resonance frequency of decoupled front cavity of vocal tract in production of vowels
  - Fant 1960
- Front part of the vocal tract
  - grows relatively little during lifetime
  - is easy to manipulate without special training
  - for many consonants, the front part dominance is well accepted

# Voiced and fricative speech



# PLP-estimated $F2'$ and Front Cavity Resonance Frequency

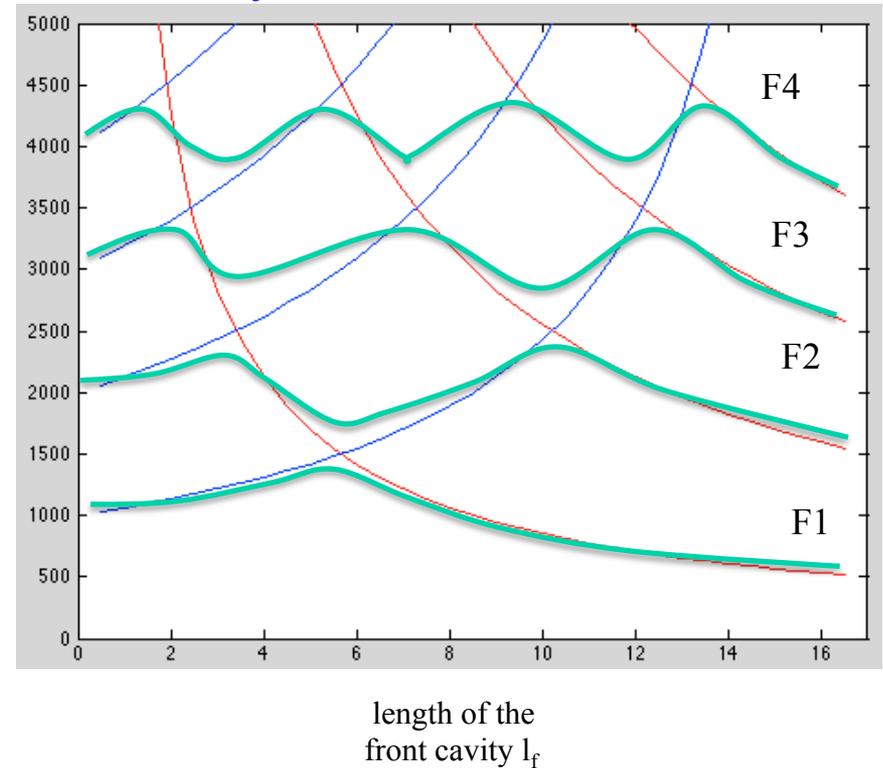
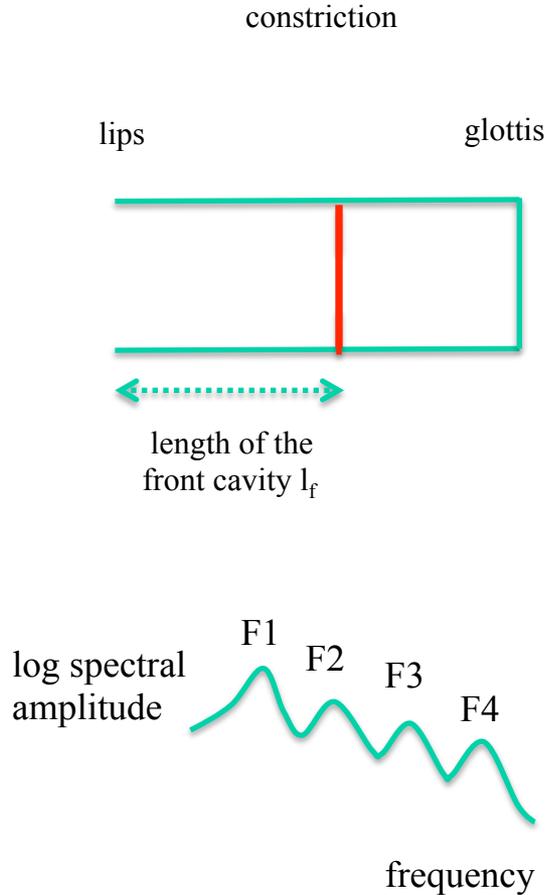
- Articulatory Synthesis
  - formants known
  - resonance frequency of decoupled front cavity can be computed
  - synthetic speech is available for analysis by PLP ( $F2'$  can be estimated)

“quarter wave resonator”

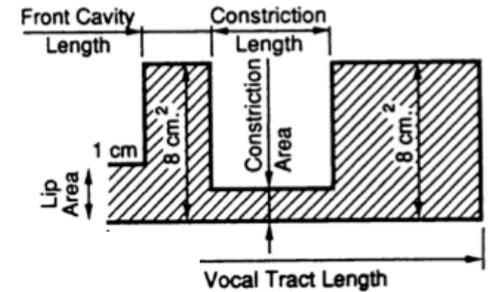
$F1 = 500 \text{ Hz}$ ,  $F2 = 1500 \text{ Hz}$ ,  $f3 = 2500 \text{ Hz}$ , ...

if the length = 17 cm and  $c = \text{m/s}$

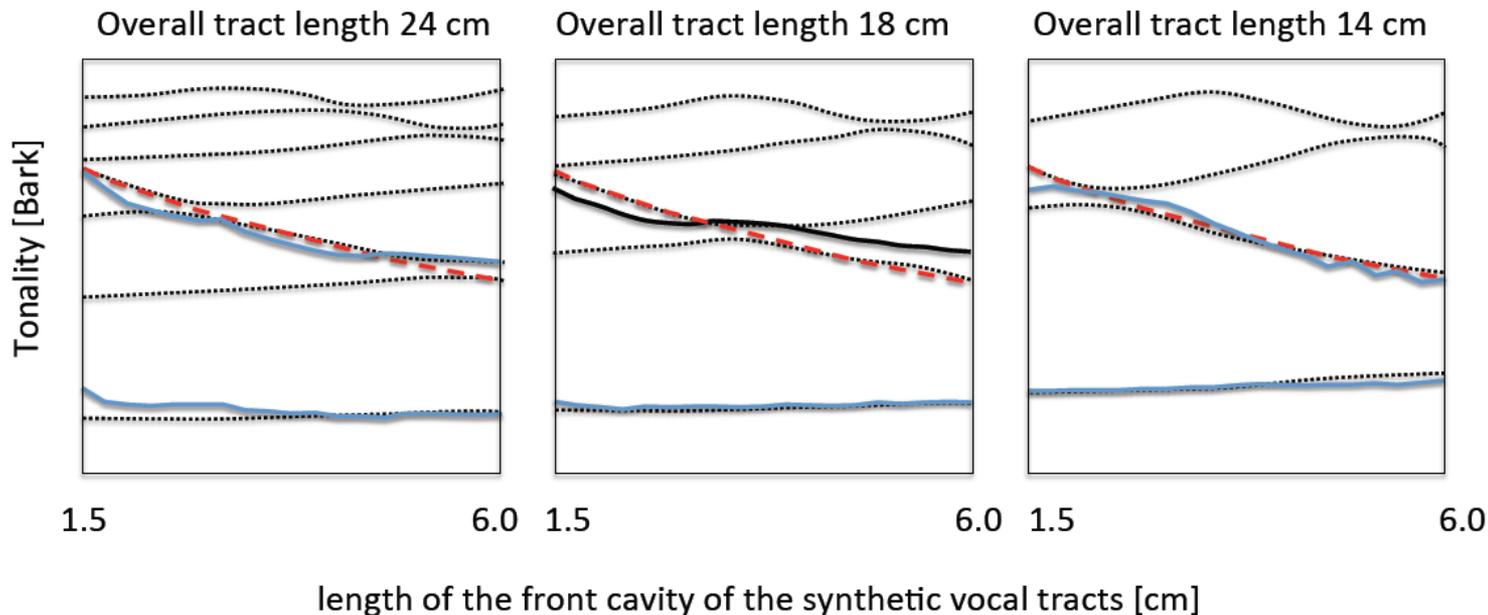
front cavity resonance modes  
back cavity resonance modes



# Front Cavity Resonance Experiment Using Articulatory Synthesis



- ..... resonance frequencies of synthetic vocal tracts (formants)
- - - first resonance of the front cavities of synthetic vocal tracts
- frequencies of peaks of the 5<sup>th</sup> order PLP autoregressive models

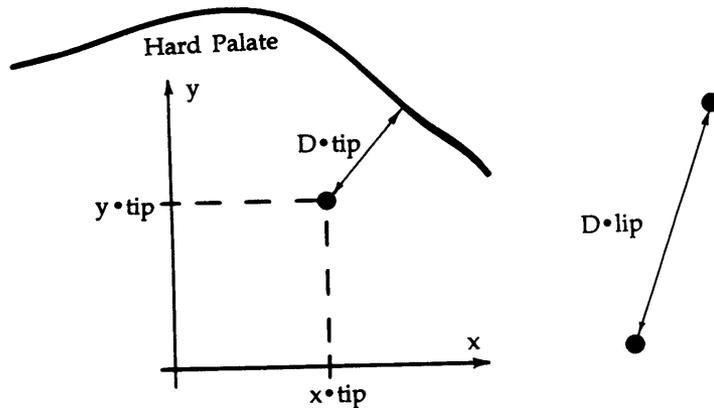


# Result of Experiment with Synthetic Vowels

- correlations on about 11 000 synthetic front vowels
  - (back vowels for which PLP formed only one peak were excluded)
  - tract length varied between 14 and 24 cm

	tract length	front cavity resonance
Second peak of PLP model	<b>-0.18</b>	<b>0.9</b>
formants (averaged)	-0.71	0.22

# X-ray Microbeam Experiment (Broad and Hermansky 1989)



- Shape approximated by cosine with period of  $2L$  and amplitude  $\Phi$
- Resonance frequency given by  $L$  and  $\Phi$  (Schroeder, Mermelstein)

$$L = k_1 - \alpha x$$

$$(a) \quad \mathbf{x} = x_{\text{tip}} \cos\theta + y_{\text{tip}} \sin\theta$$

$$\Phi = k_2 + b_1 \ln D_{\text{tip}} + b_2 \ln D_{\text{lip}}$$

(b)

$\frac{1}{F_2'}$	$=$	$\frac{4L}{c}$	$\frac{2}{2 + \Phi}$
------------------	-----	----------------	----------------------

(c)

PARAMETERS:

$k_1, k_2, \alpha, \theta, b_1, b_2$

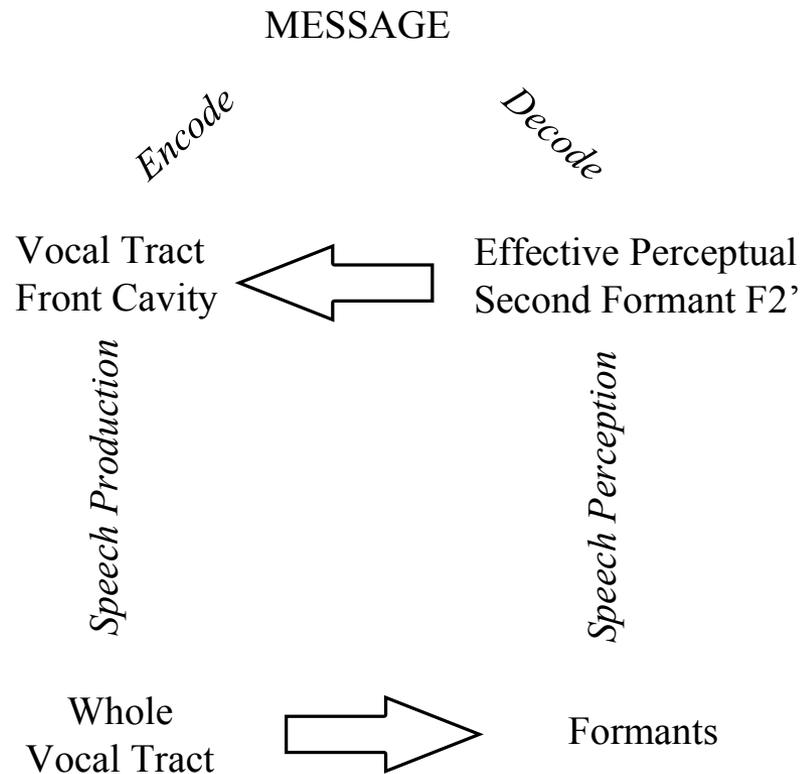
# Results of X-Ray Microbeam Experiment

- two male speakers
  - “where were you a year” three times each
- front cavity resonance from articulations
- PLP-estimated F2’ from acoustic data

## CORRELATION BETWEEN RESONANCE FREQUENCY OF FRONT CAVITY AND PLP-DERIVED F2’

speaker 1	correlation 0.95
speaker 2	correlation 0.96

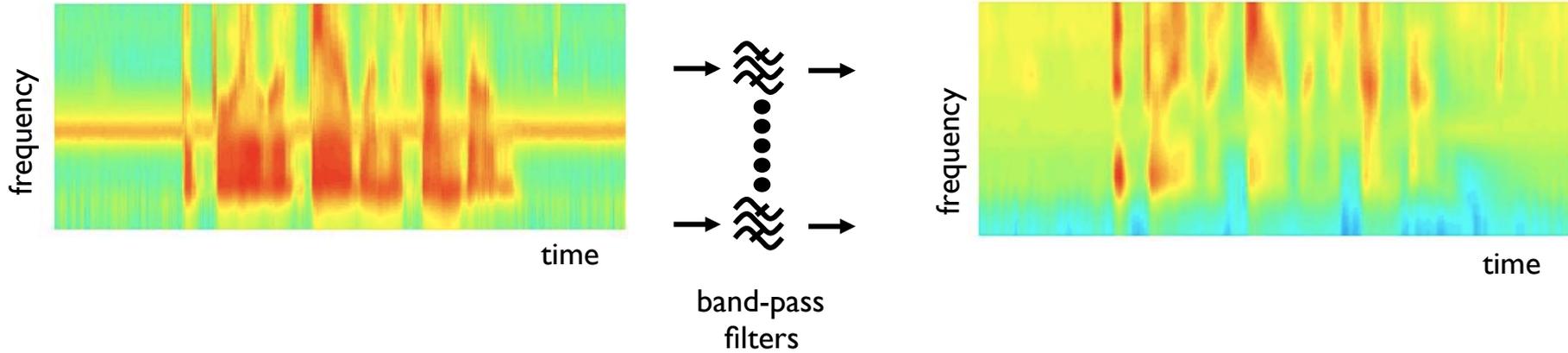
# Front Cavity - F2' Hypothesis



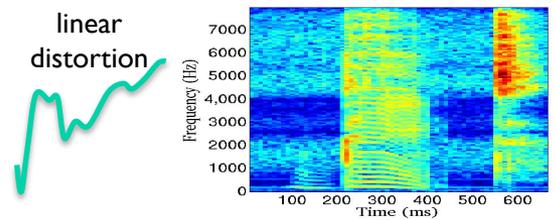
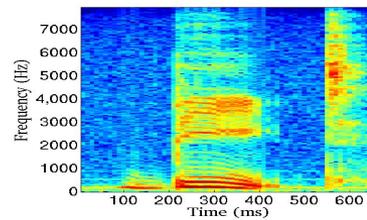
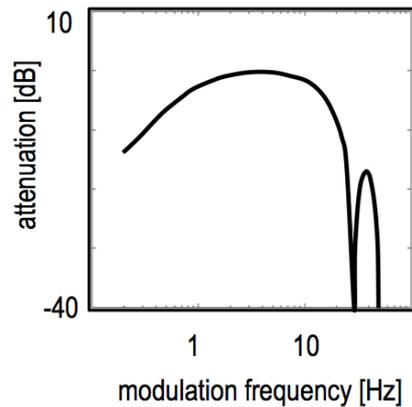
- **Our limited experimental data do not contradict the hypothesis**

# RASTA processing

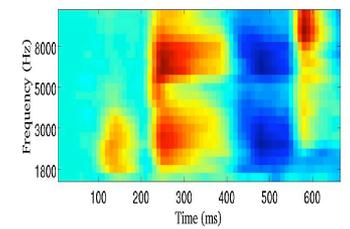
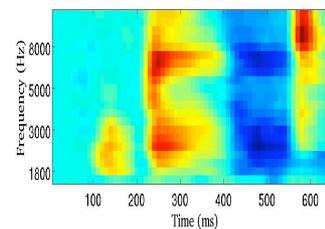
Hermansky and Morgan 1990



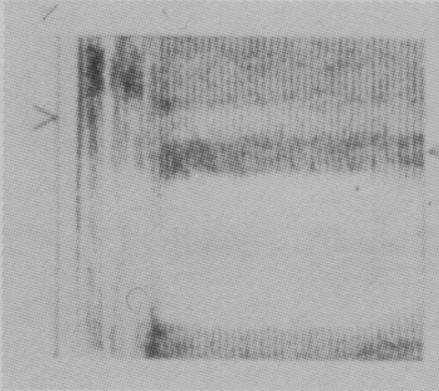
optimized RASTA filter



spectrum from RASTA-PLP



/k/



ki (*key*)

Potter, Kopp, and Green, Visible Speech 1947

need to know the following vowel before  
identifying the consonant ?

recognize whole syllables ?

**recognize phonemes but use information from  
syllable-length segments of the signal !**

- V. A. Kozhevnikov and L. A. Chistovich, Speech: Articulation and perception. 1967

h e l l o w o r l d

about 70 ms



about 200 ms

time

> 200 ms

classifier



FDPLP

spectro  
temporal  
features

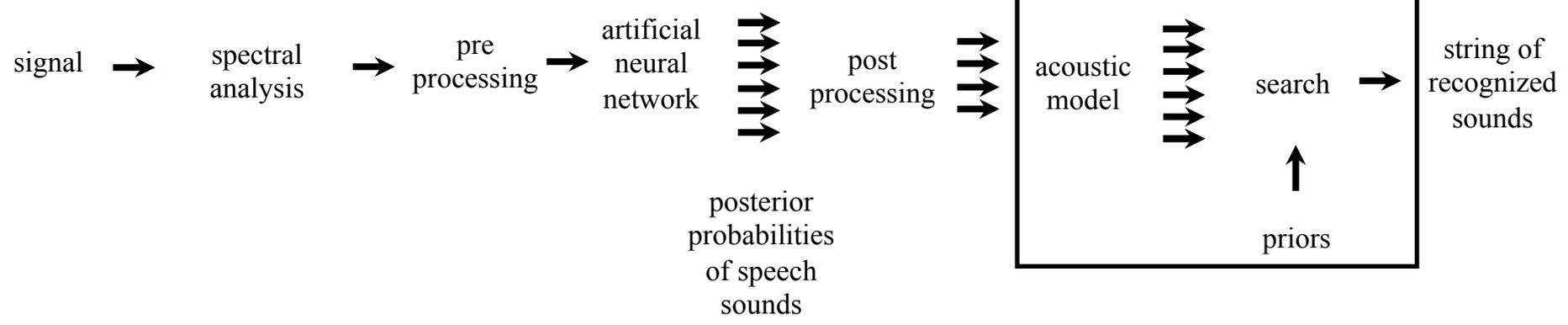
serial

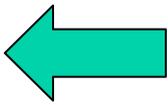
TANDEM

MRASTA

parallel

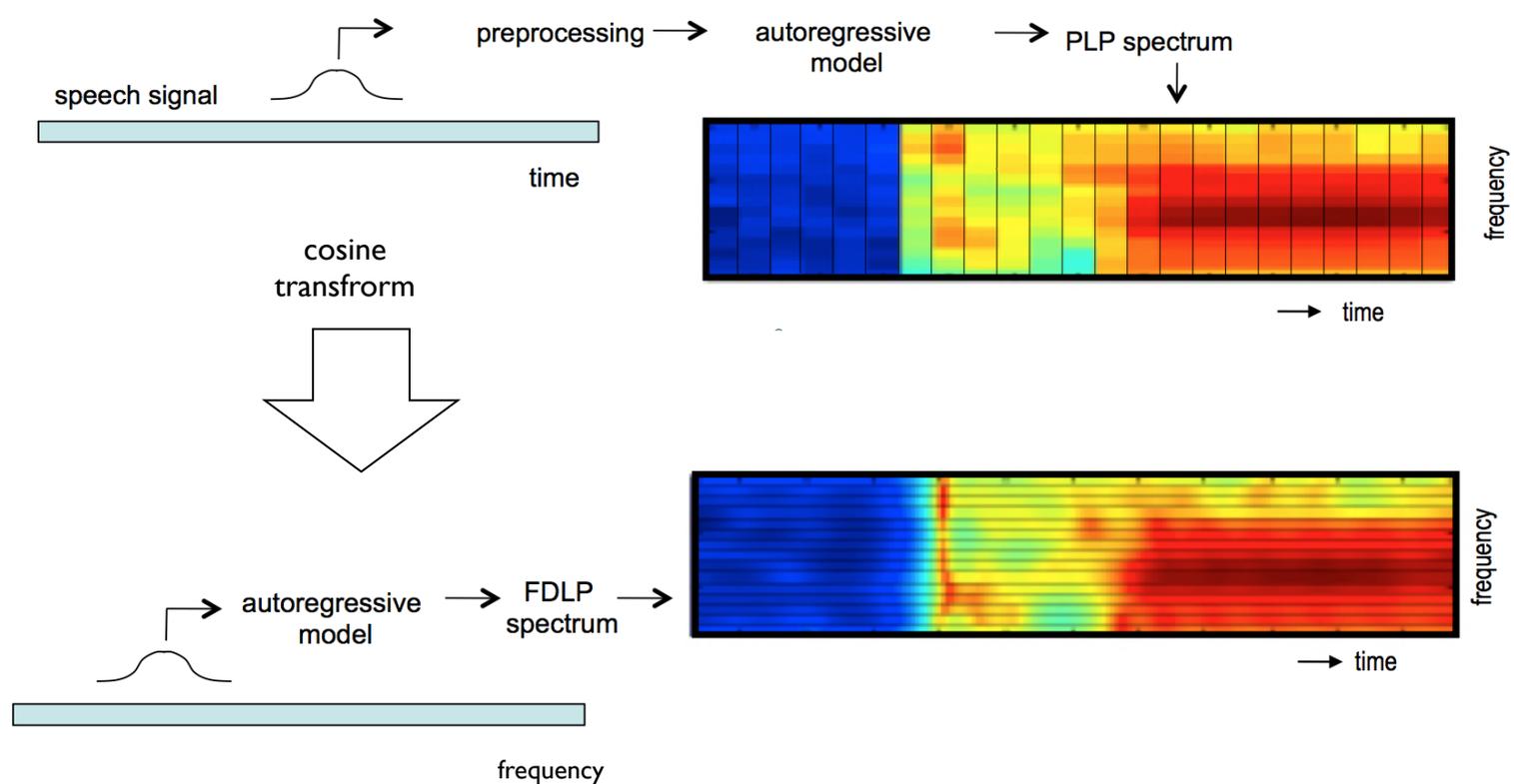
HATS  
(bottleneck)





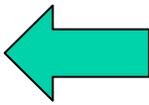
# Frequency Domain Perceptual Linear Prediction (FDPLP)

-with Marios Athineos, Dan Ellis, Sriram Ganapathy and Samuel Thomas

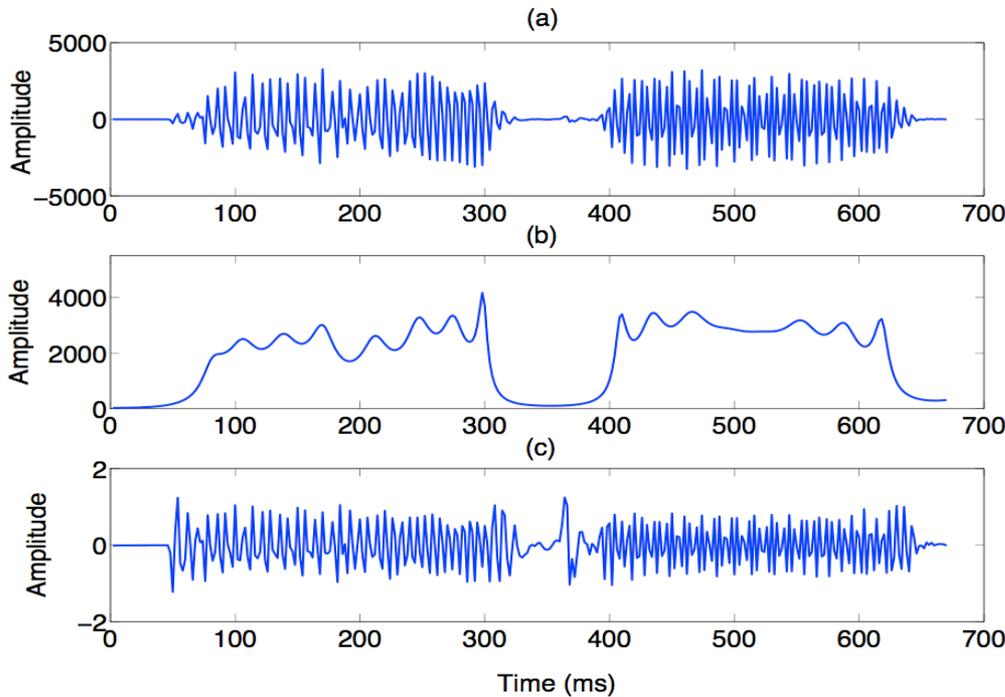


Decomposition into AM and FM components.

Straightforward alleviation of effects of linear distortions and , reverberations .



# FDLP decomposition of the signal

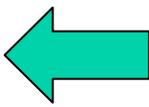


signal

AM component  
(temporal envelope)

FM component  
(carrier)



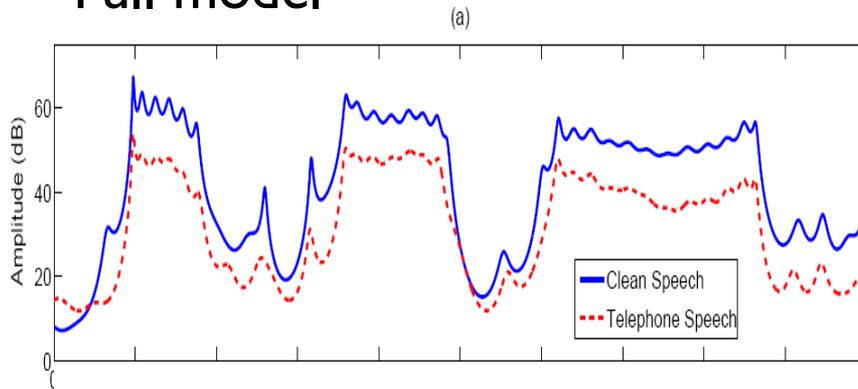


# Varying communication channels

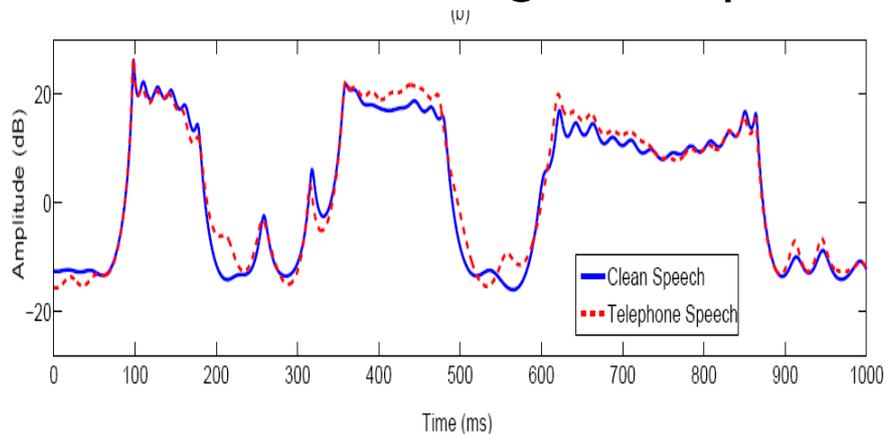
(convolution with a short impulse response of a channel)

Convolution turns into addition in log spectral domain

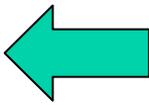
## Full model



## Model without its gain component



Ignoring FDPLP model gain makes the representation invariant to linear distortions introduced by the communication channel.



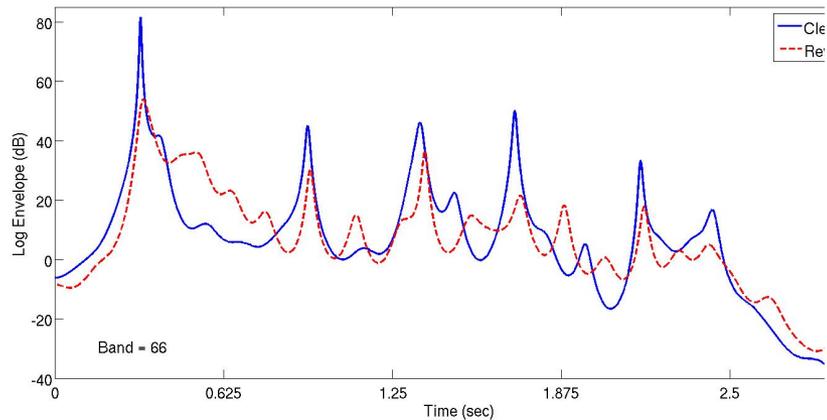
# Reverberant speech

(convolution with a long impulse response of the room)

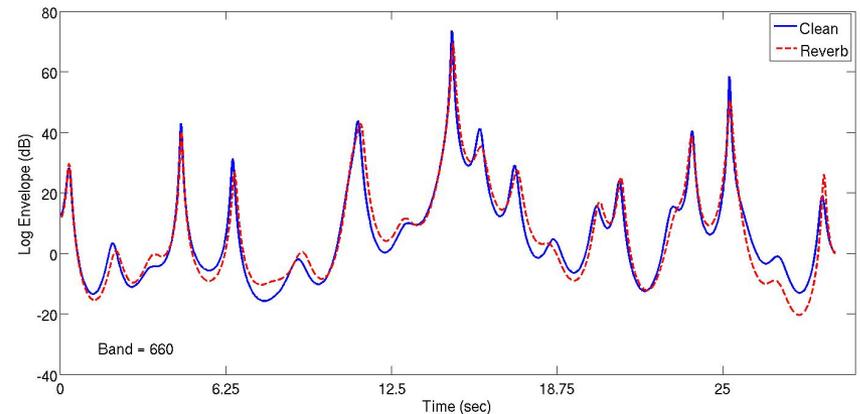
Convolution turns into addition in log spectral domain, **as long as the most of the room impulse response fits into the analysis window!**

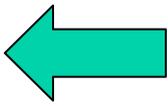
Ignoring FDLP model gain makes the representation invariant to reverbs.

## 3 s window



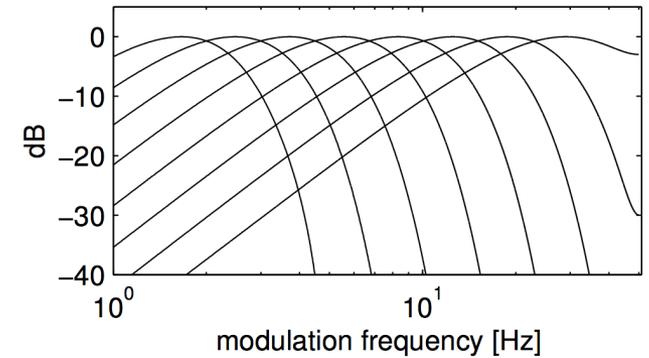
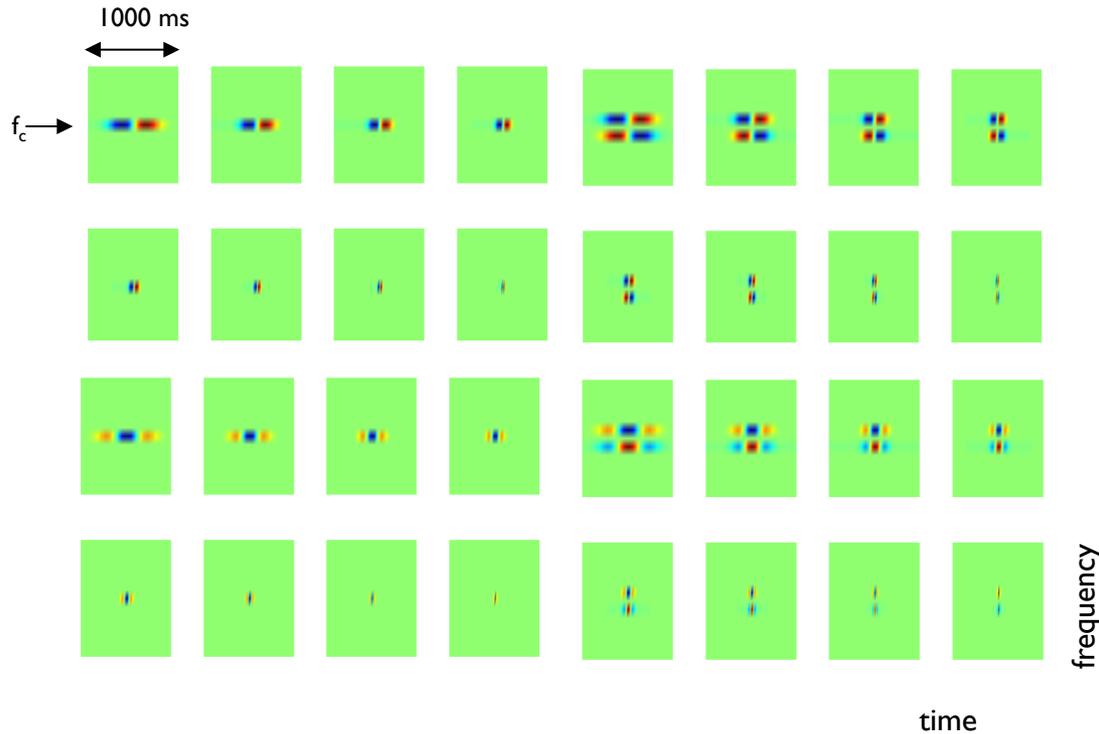
## 30 s window





# MRASTA

Hermansky and Fousek 2005



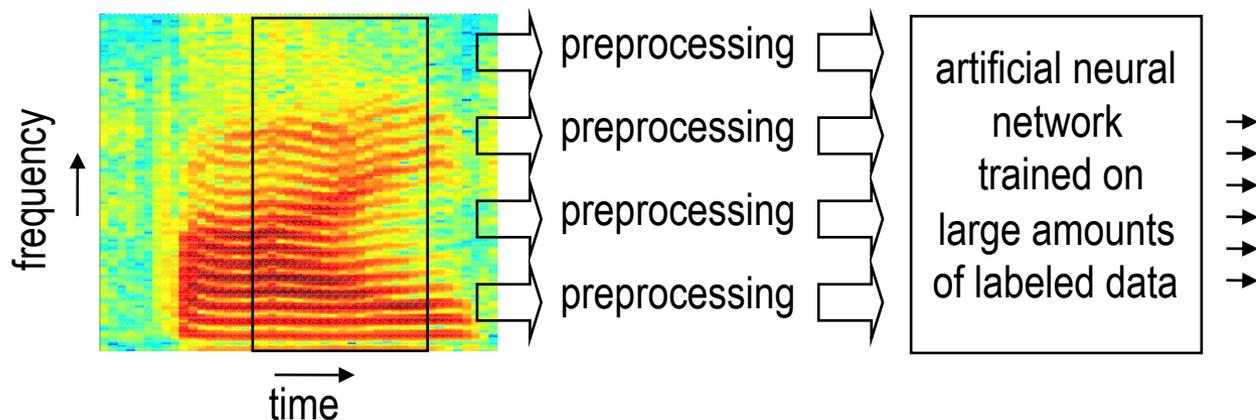
32 features at each of 14 frequencies

**448 dimensional vector of features every 10 ms**

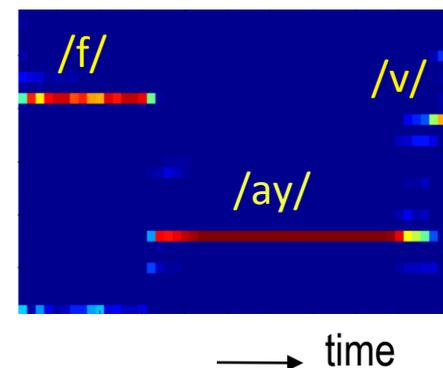
multi-resolution band-pass filtering of modulation spectrum

Optimal (lowest dimensionality) features are posterior probabilities of classes

## Spectrogram



## Posterlogram



Training of the artificial neural net

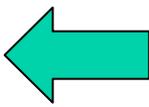
data representing  
phoneme 4



artificial  
neural  
net



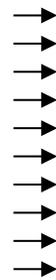
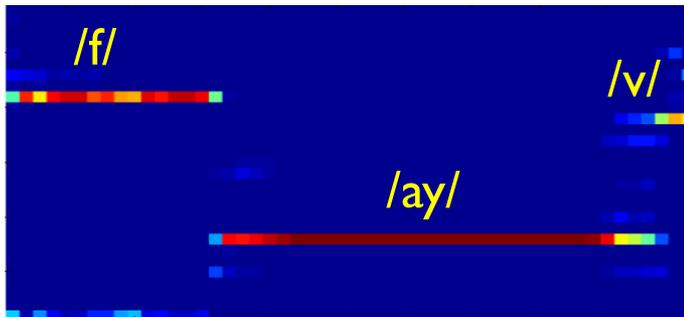
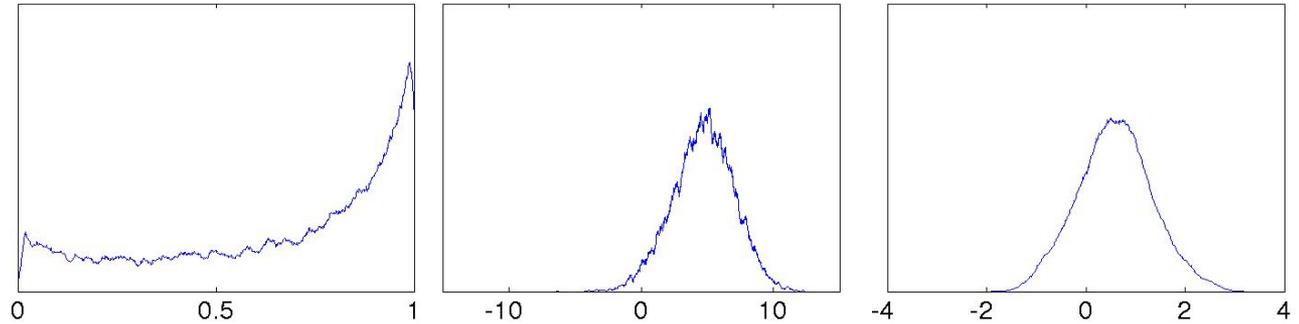
# TANDEM Features for HMM/GMM System



Hermansky, Ellis, and Sharma 2000

good attributes for state-of-the-art ASR systems  
should be Normally distributed and uncorrelated

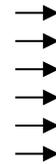
histogram of  
one feature



pre-  
softmax  
outputs

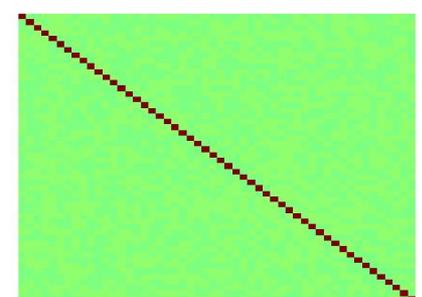
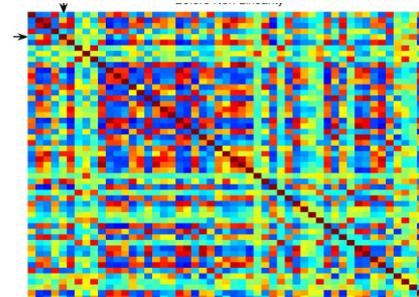
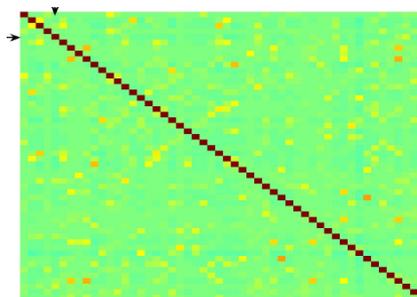


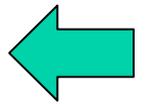
principal  
component  
projection



to HMM

correlation  
matrix  
of features

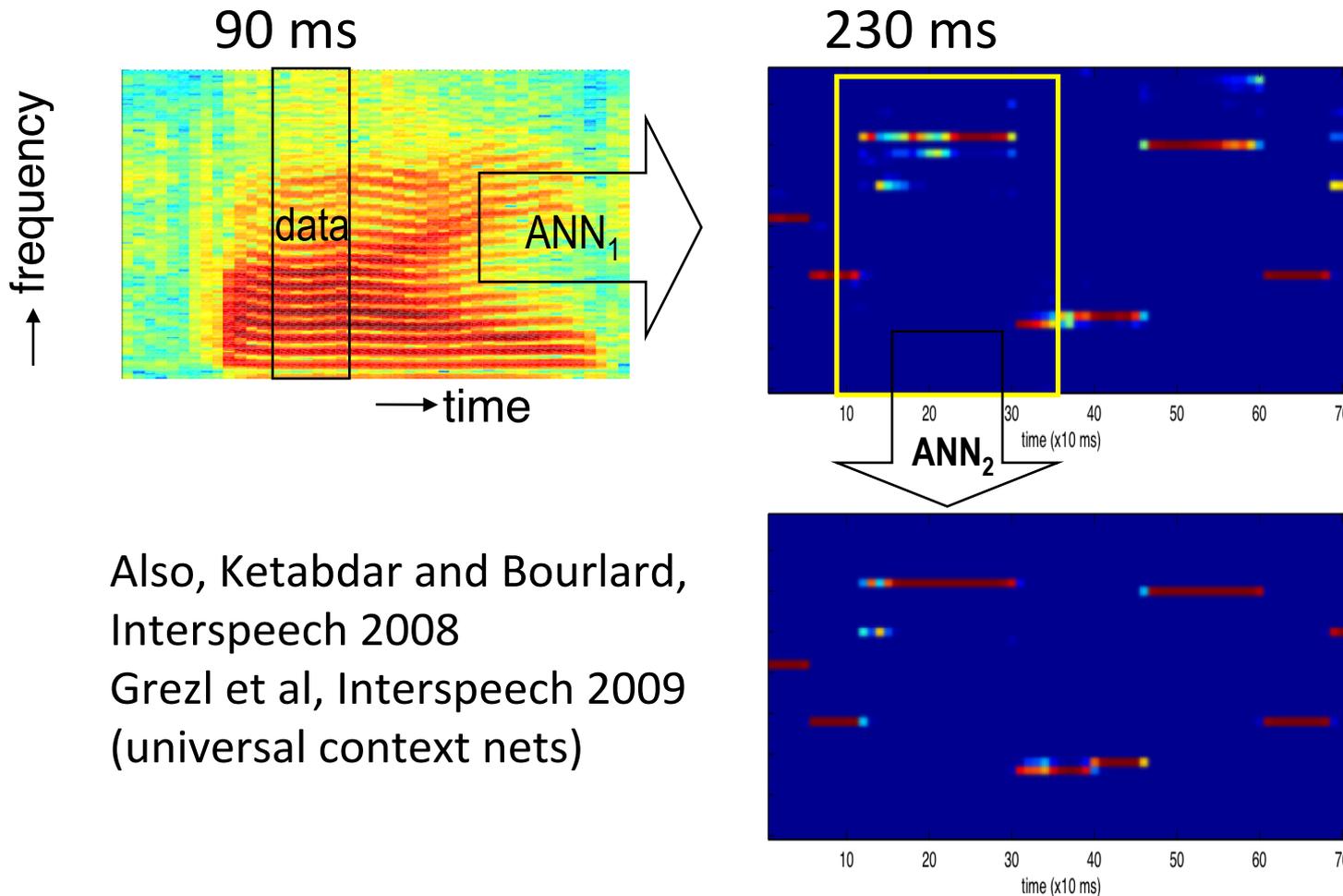




# Serial hierarchical estimation

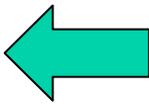
(Pinto et al, Interspeech

2008)



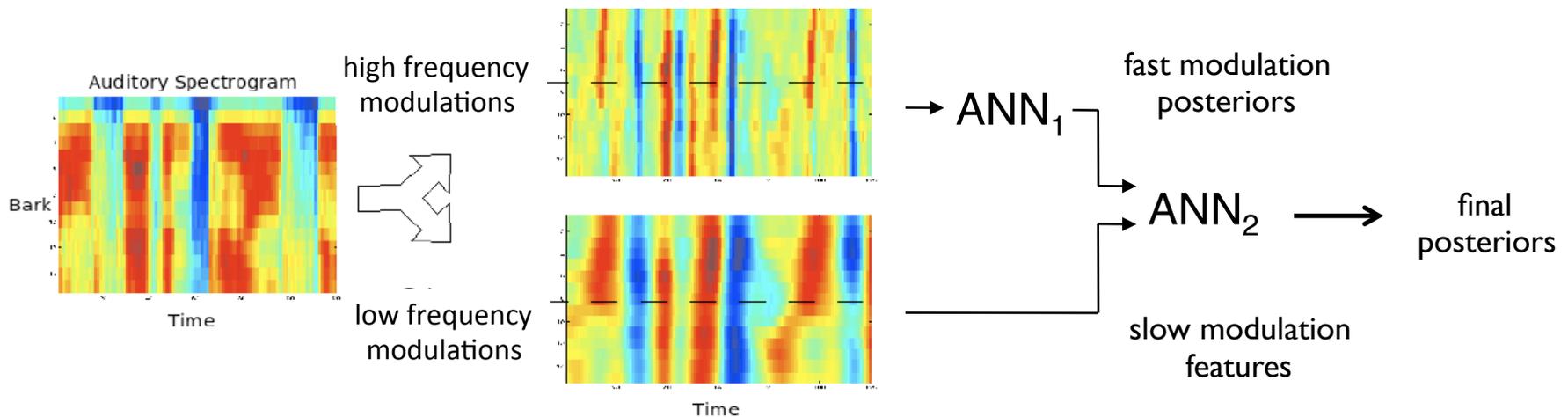
Also, Ketabdard and Boulard,  
Interspeech 2008  
Grezl et al, Interspeech 2009  
(universal context nets)

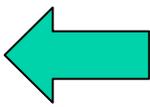
# Parallel hierarchical estimation



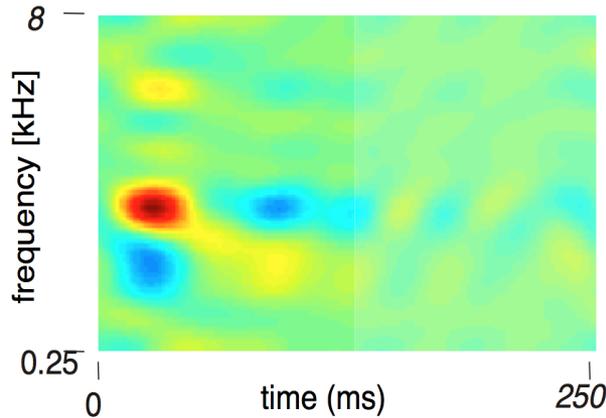
Valente and Hermansky, ICASSP 08, Interspeech 2008

- one-stage processing on coarse (slow modulations) representation
- two-stage processing of finer (faster modulations) representation





# Auditory cortical spectro-temporal receptive fields (STRFs)



indicate “optimal” stimulus that excites a given cortical neuron

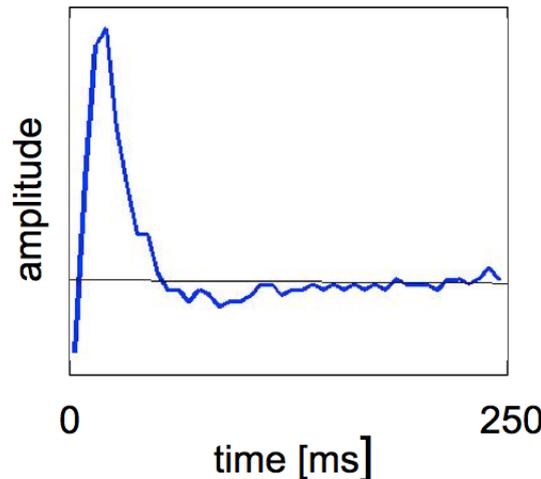
from Shihab Shamma

compute principal components along temporal axis of about 300 STRFs

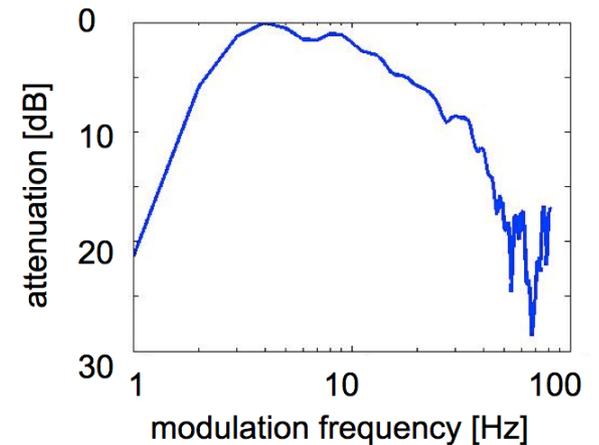
Nima Mesgarani (in preparation)

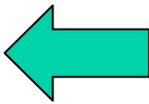
higher components similar but shifted in time

1<sup>st</sup> principal component (41% of variance)

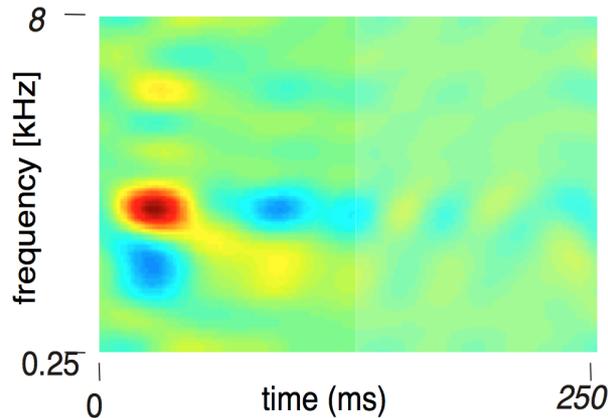


1<sup>st</sup> component frequency response





# Auditory cortical spectro-temporal receptive fields (STRFs)



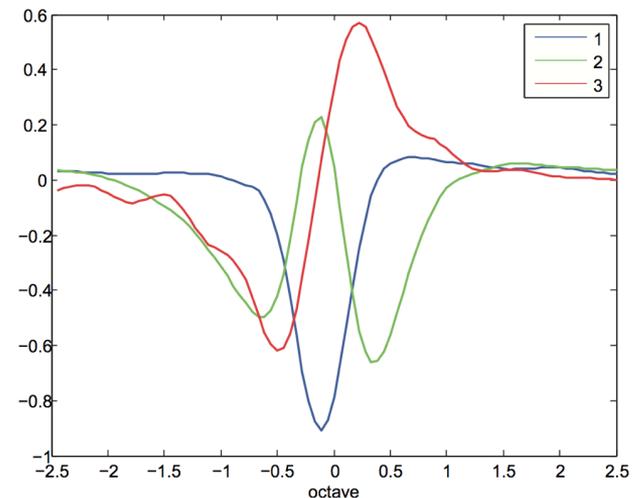
from Shihab Shamma

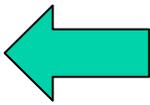
indicate “optimal” stimulus that excites a given cortical neuron

align maxima of STRFs in frequency and compute principal components along frequency axis of about 300 STRFs

Nima Mesgarani (in preparation)

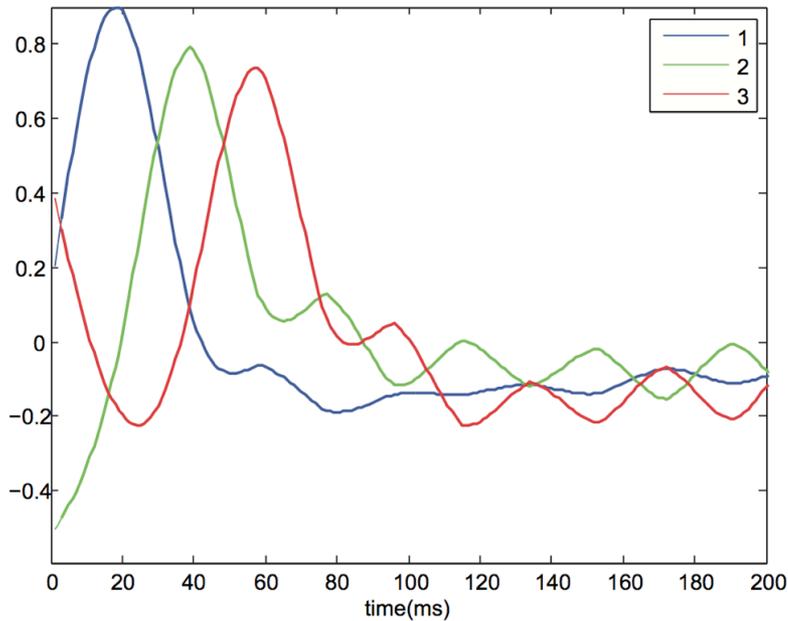
Principal components of spectral axis



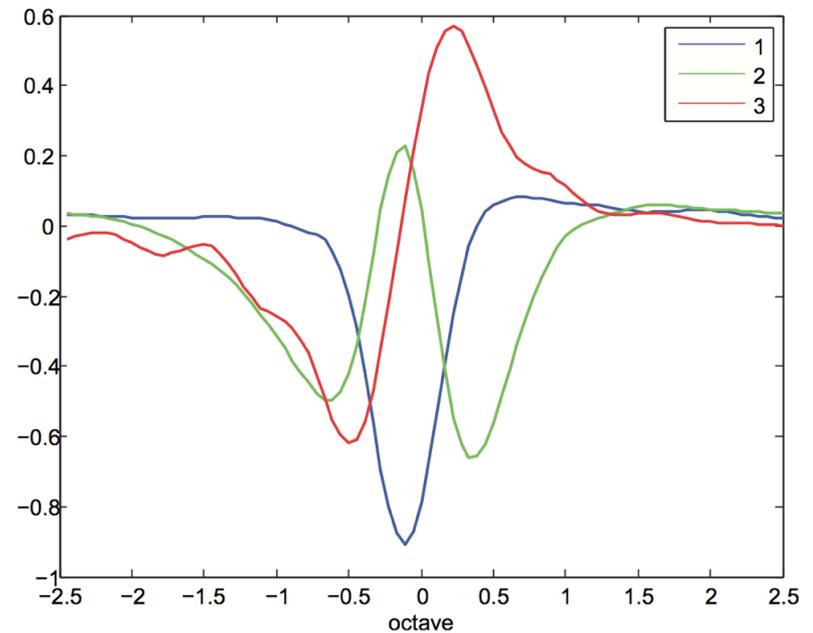


- principal components of about 300 STRFs
  - Nima Mesgarani (in preparation)

Principal components of temporal axis

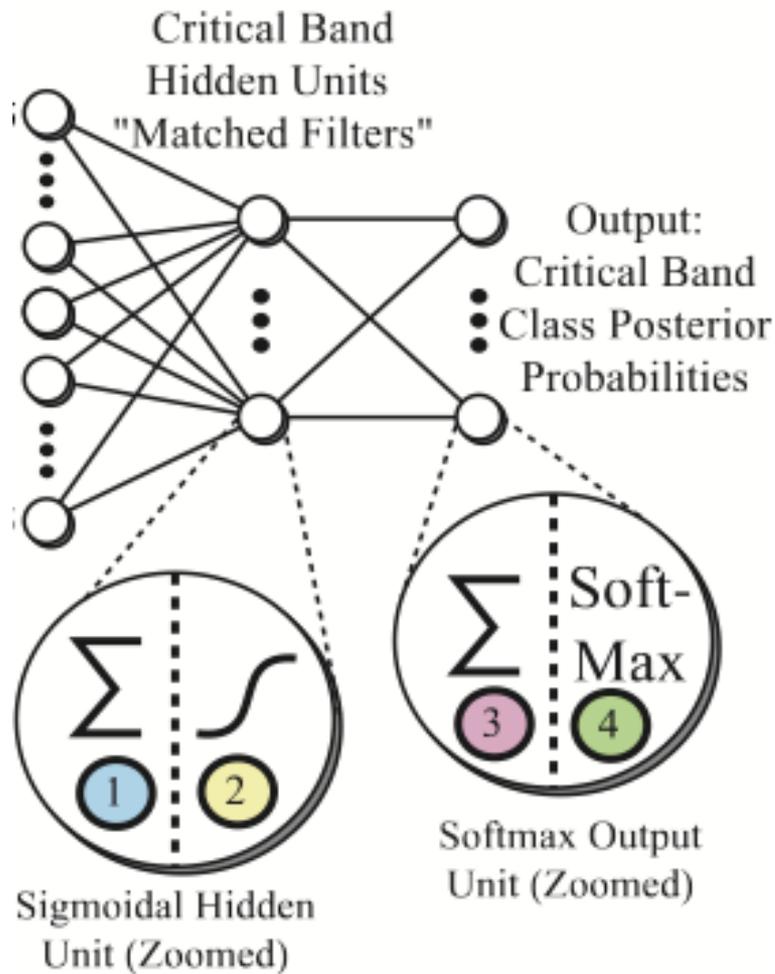
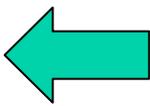


Principal components of spectral axis

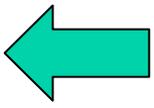


# HATS (bottleneck features)

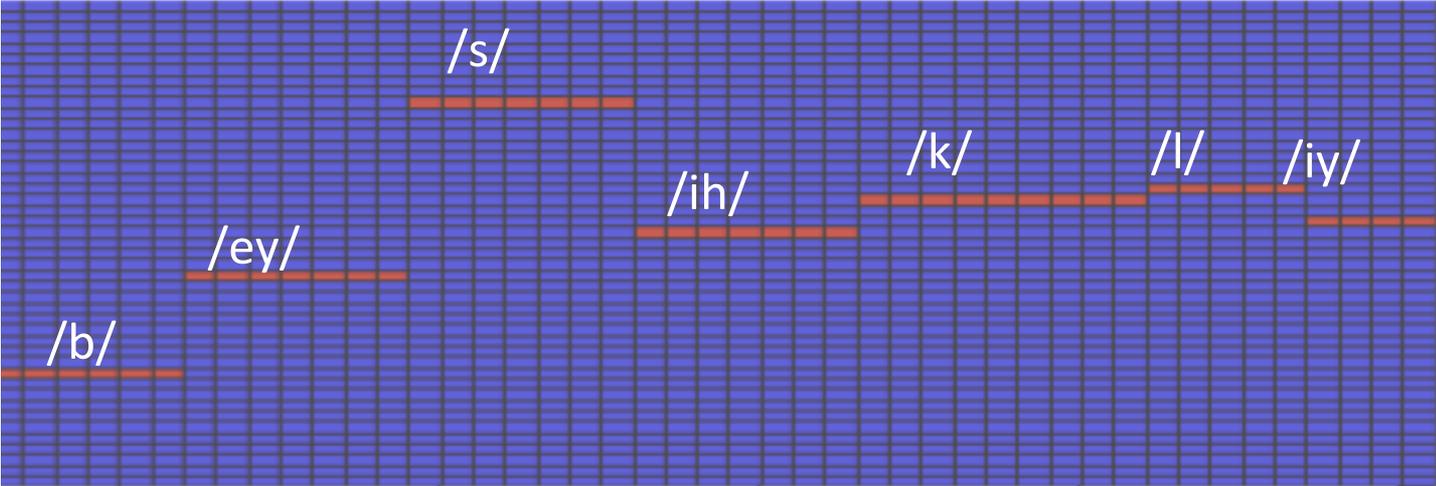
Chen et al 2005, Grezl et al 2007



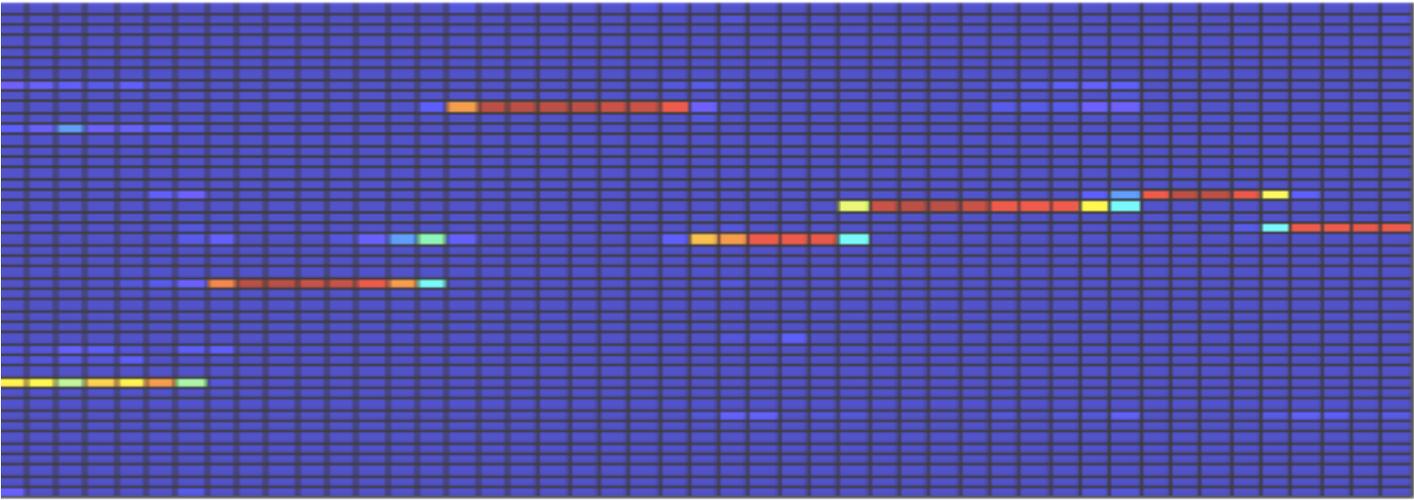
4 - posteriors  
3 - for TANDEM  
**2 - HATS**

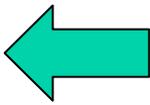


truth

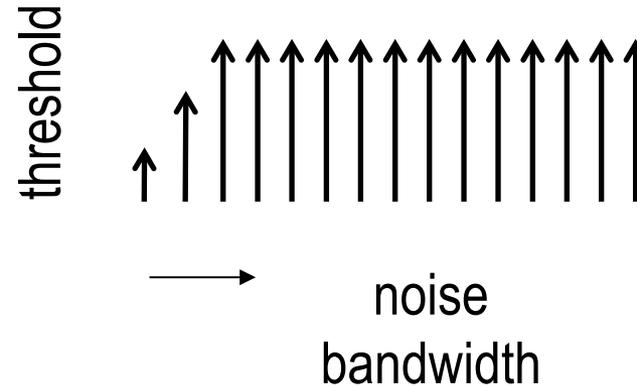
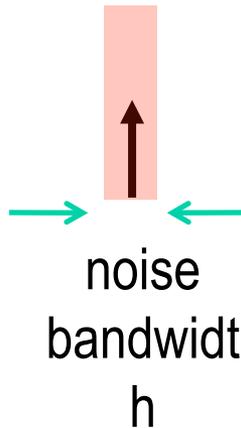
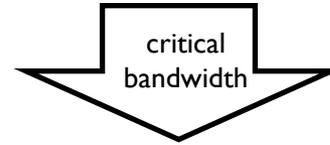


estimate



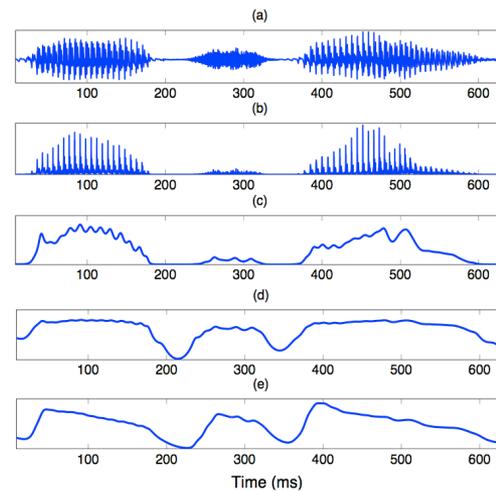
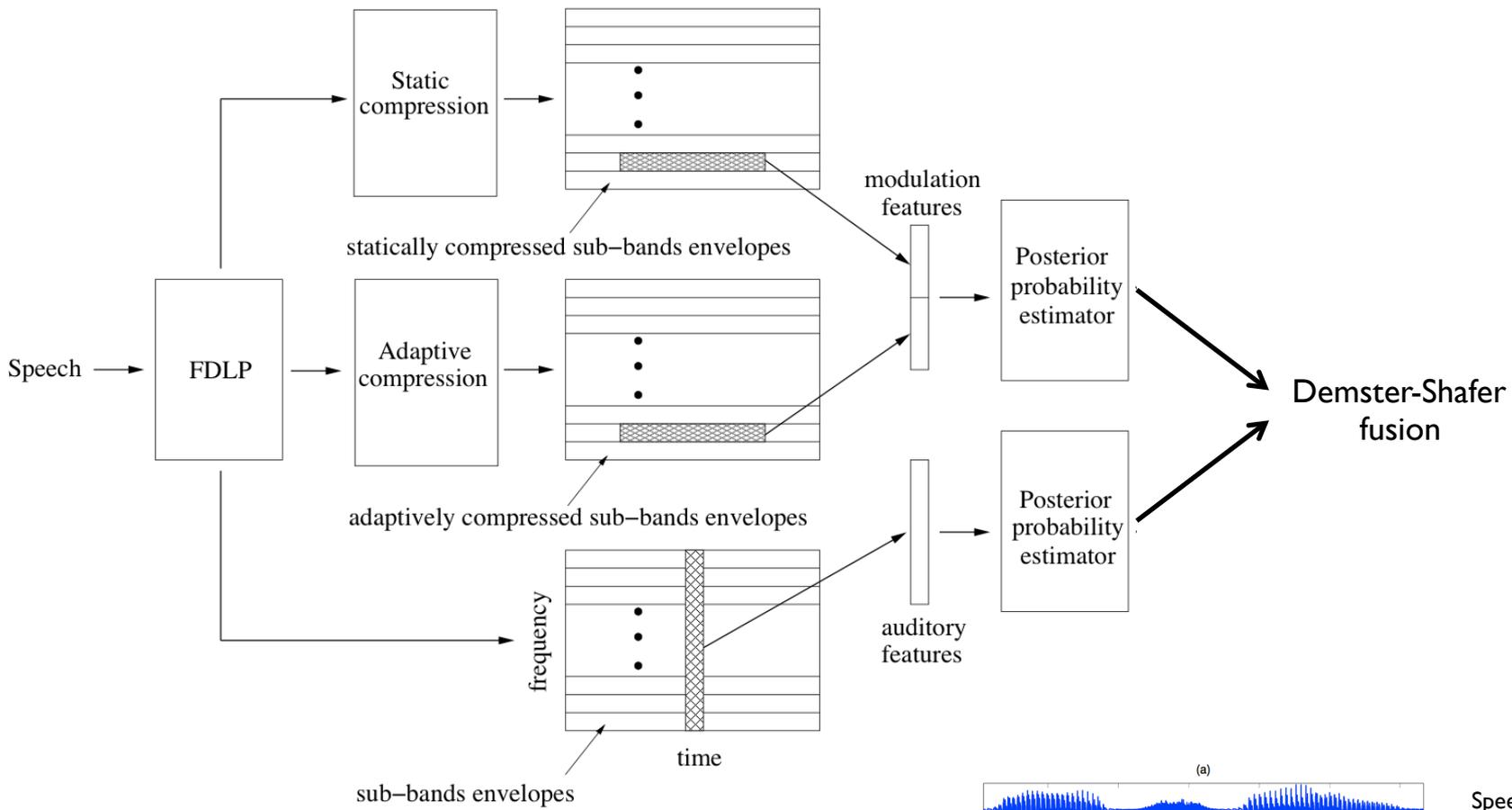


simultaneous  
masking



spectral components inside the (critical) band interact differently with components inside the band than they do with components, which are outside the band

**hearing periphery does spectral analysis to allow for separation of corrupted signal elements at higher levels of auditory processing**



Speech signal

Hilbert transform of the signal

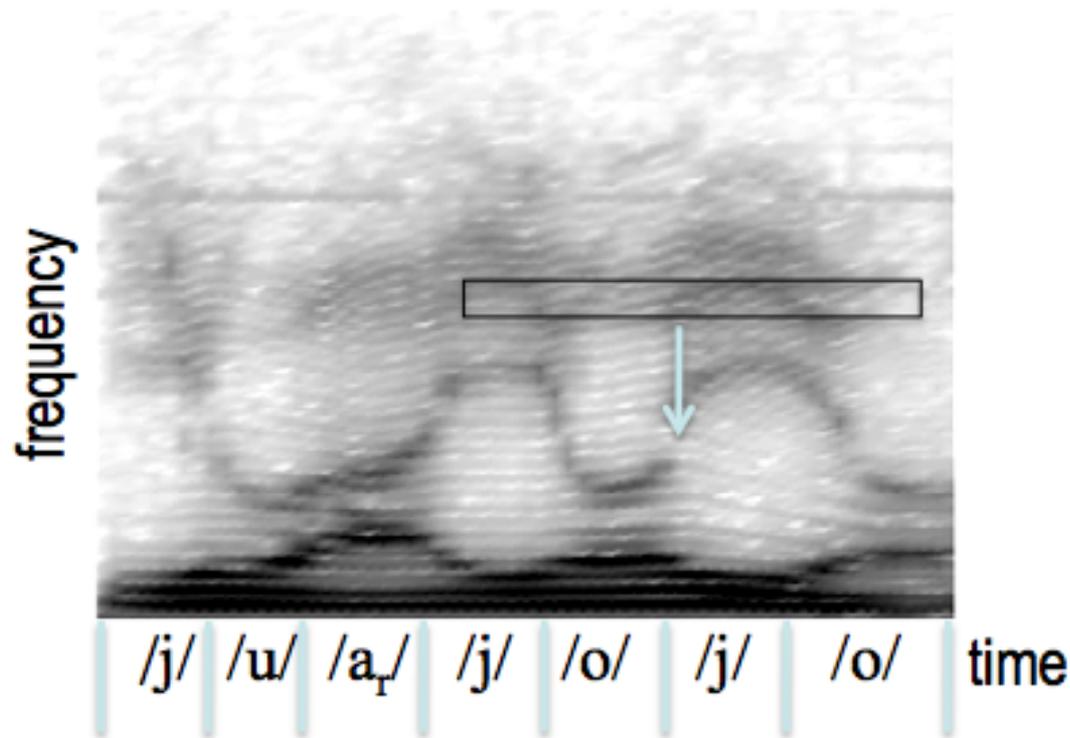
FDLP envelope

Static compression of the envelope

Dynamic compression of the envelope.

# Data-guided FIR RASTA filters

van Vuuren and Hermansky 1997, Valente and Hermansky 2006

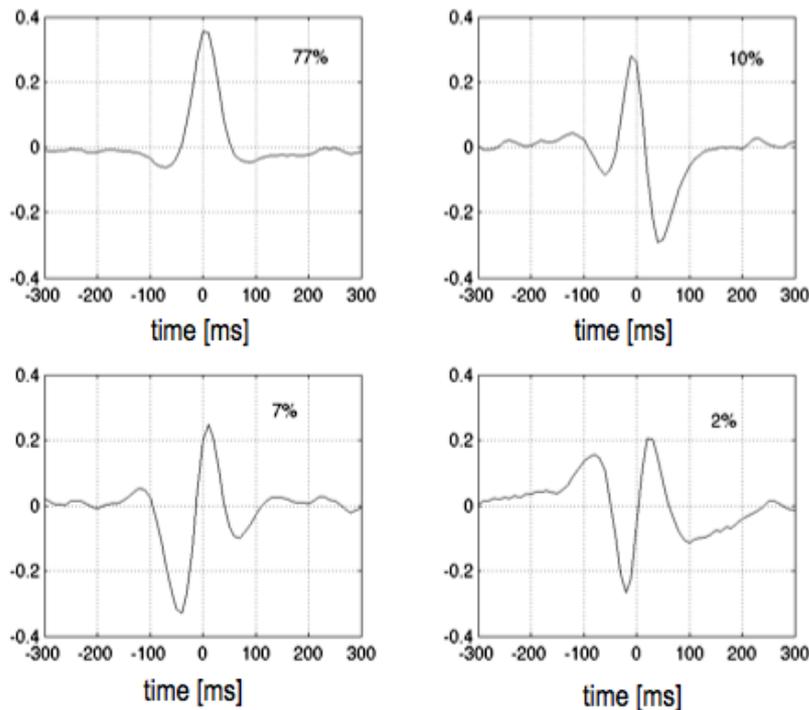


# Data-guided (LDA-based) FIR RASTA filters

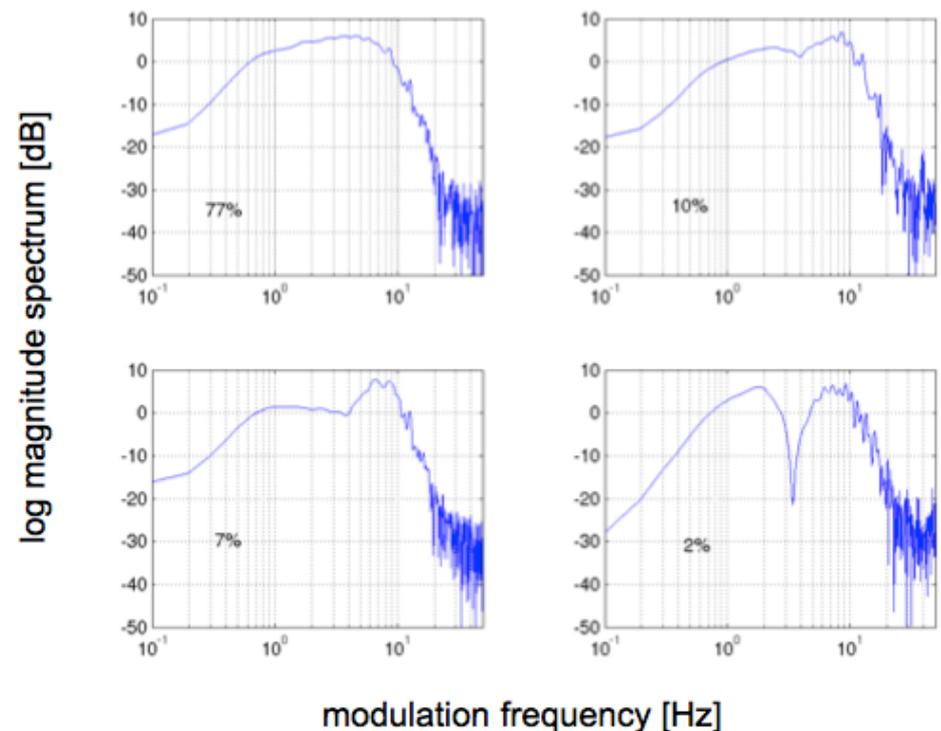
van Vuuren and Hermansky 1997, Valente and Hermansky 2006

first 4 temporal linear discriminants

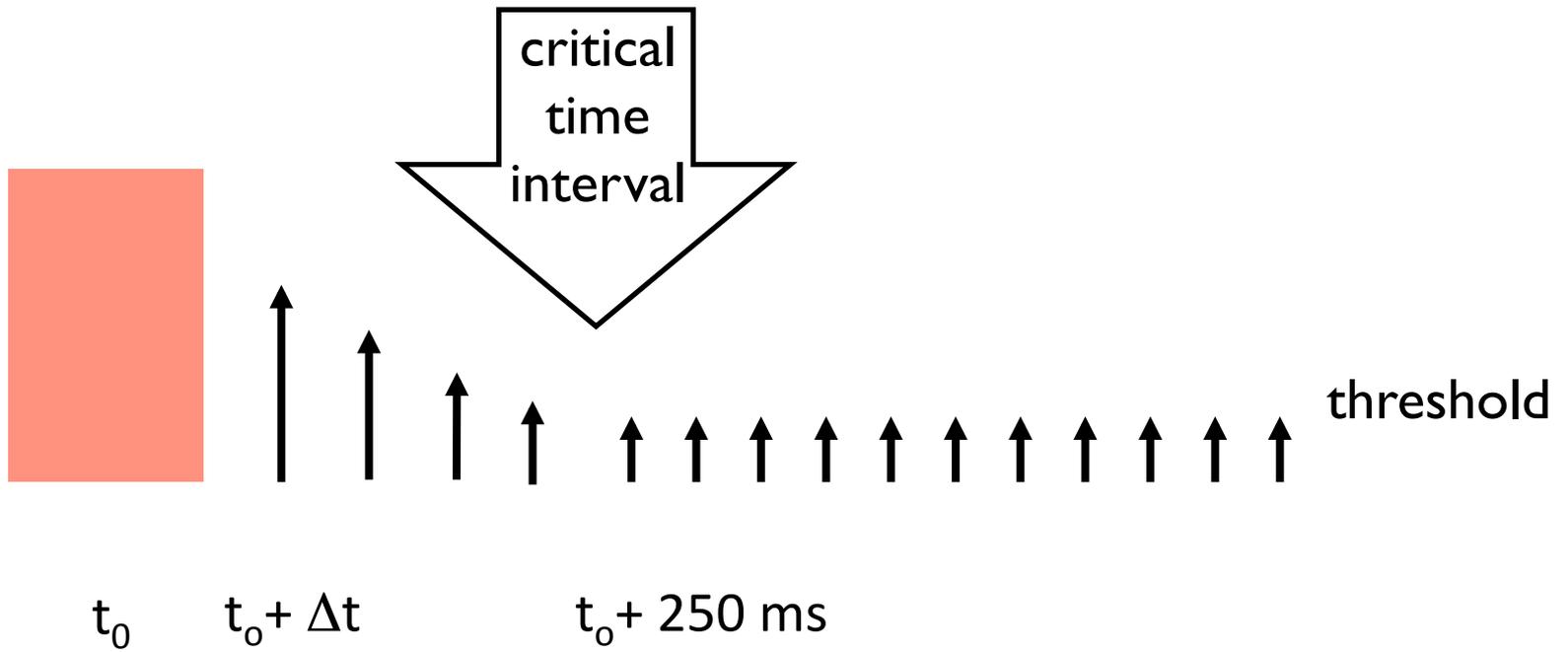
impulse responses



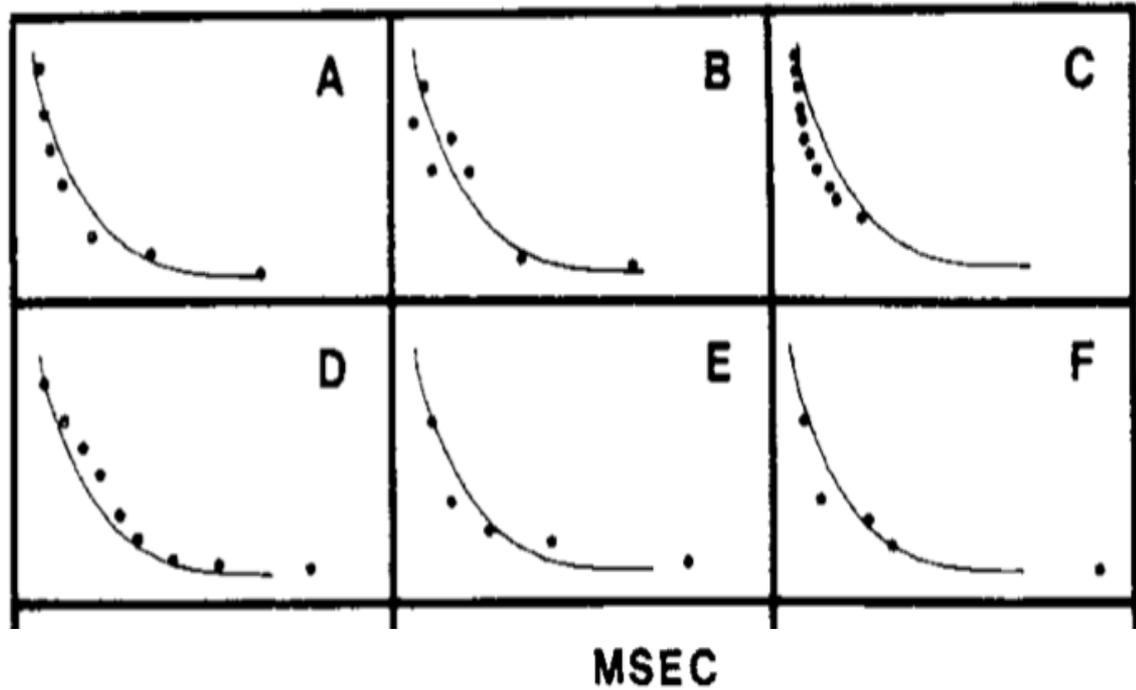
frequency responses



forward  
masking

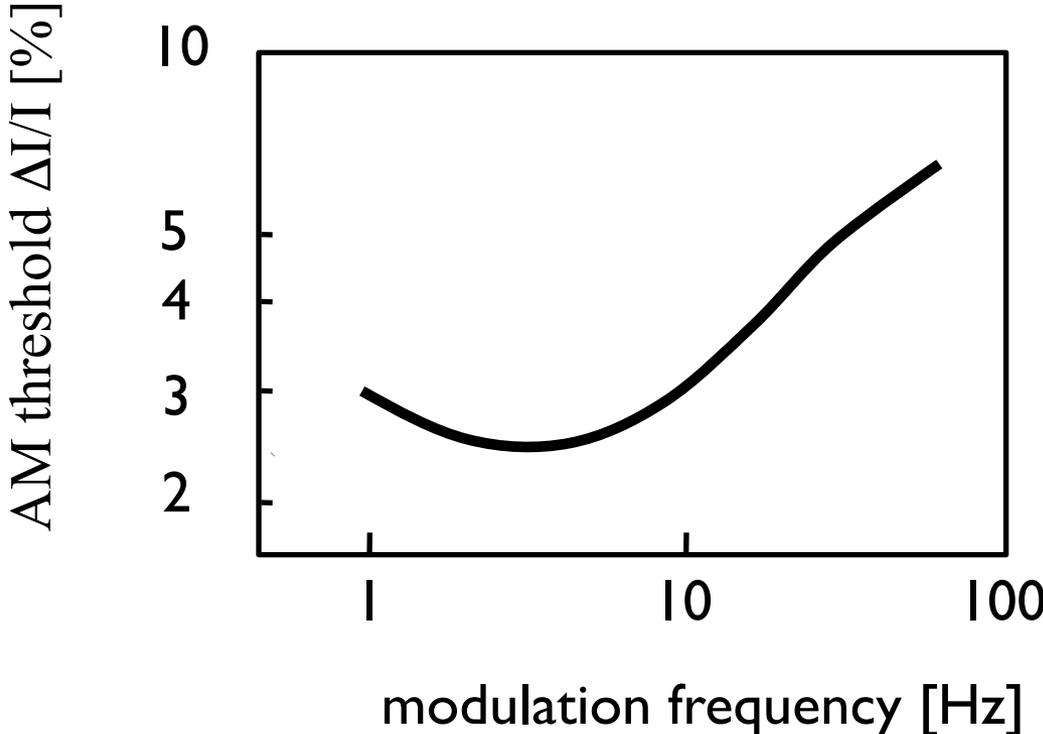


- A-forward masking
- B-backward masking
- C-gap detection
- D-overestimation of short burst duration
- E-loudness decrement
- F-JND in frequency



N. Cowan, On Short and Long Auditory Stores, Psychological Bulletin 1984

# Riezs 1928



Kozhevnikov and Chistovich  
(Speech: Articulation and  
Perception, 1965)

- reaction times for identifying consonants and vowels in CV syllables
- consonant always identified before a vowel

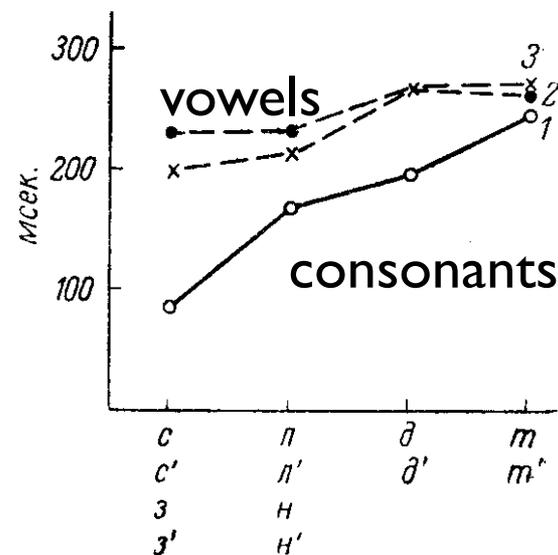


Рис. 6.6. Задержки буквенной записи согласных в зависимости от их качества (1) и задержки буквенной записи гласных в зависимости от качества предшествующих согласных (2 — твердый согласный, 3 — мягкий согласный).

To recognize phoneme one needs to collect information distributed over the whole syllable