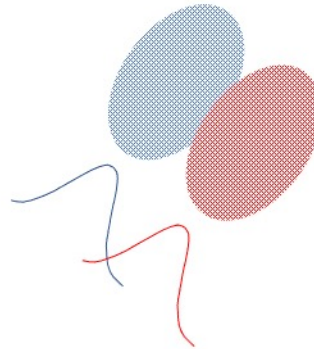


### Linear discriminant analysis (Ronald A. Fisher 1936)

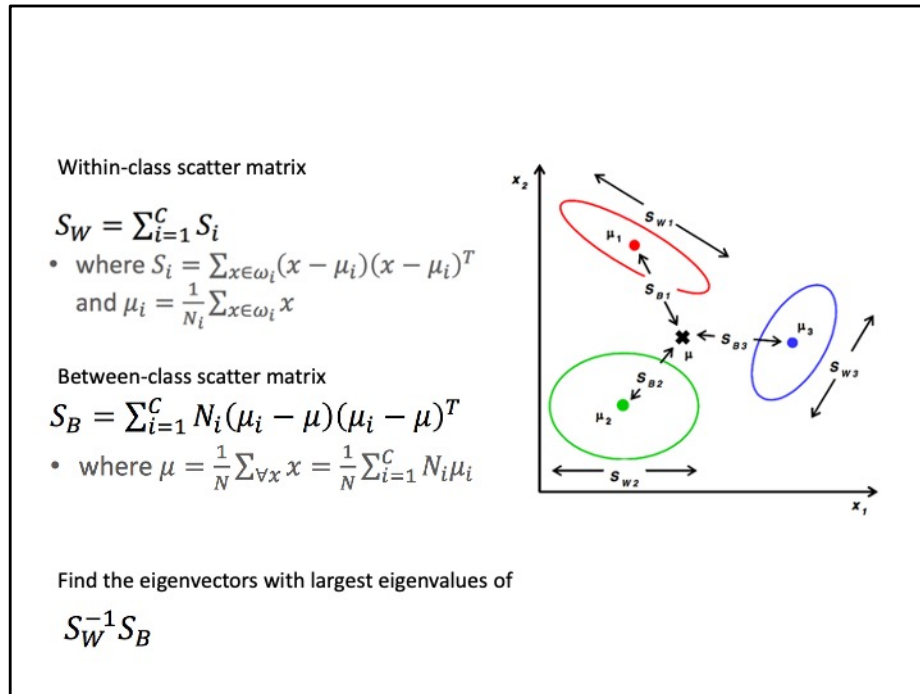
- find projections of data, which preserves most of the discriminability
- data vectors need to be labeled by classes
- yields matrix of discriminant vectors, ordered by their discrimination power
- discriminants are linear and therefore can be easily interpreted



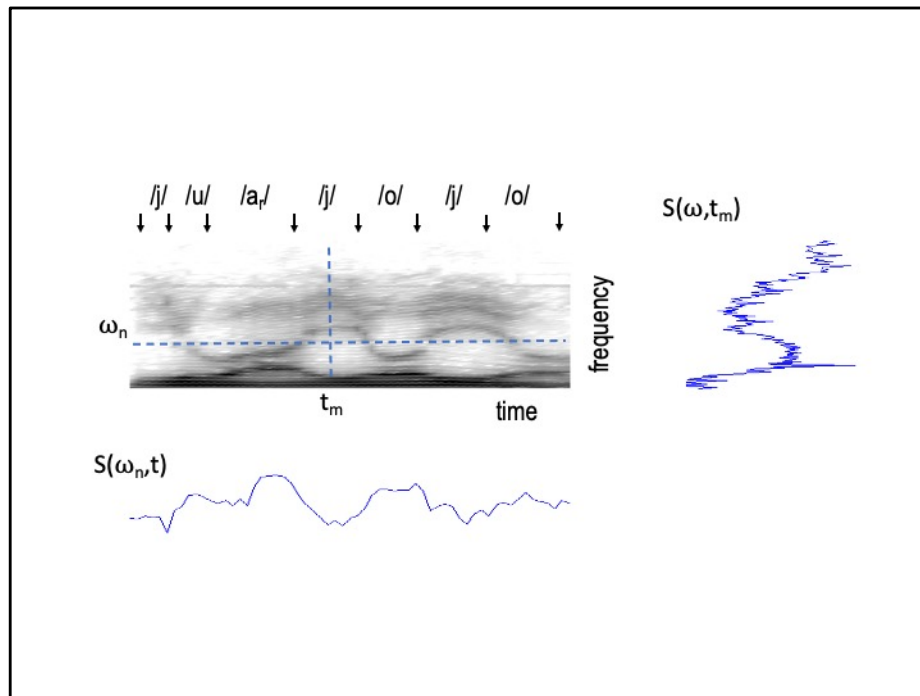
*The Acoustic Modeling Problem in Automatic Speech Recognition*  
Peter Brown, CMU CS Department, AFWAL-TR-87-1161

One well established technique to emphasize discriminability in a feature vector space is the linear discriminant analysis (LDA). LDA needs labeled vector space, where it knows which class each vector represents. LDA finds a projection of the space in directions of the class discriminability. To preserve all original discriminability, all dimensions of the new space are needed. However, most discriminability is in the first few dimensions of the new vector space. The amount of discriminability in each direction is indicated by the values of the eigenvectors of the discriminant matrix, used in the projection.

One of the first (if not The First) use of LDA in speech recognition was by Peter Brown at CMU in 1987.



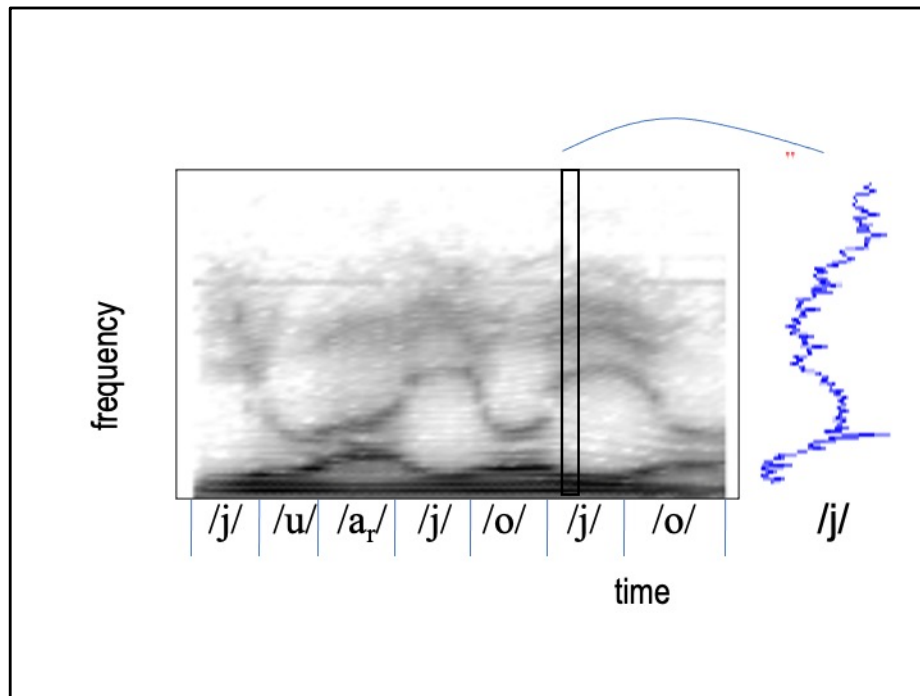
Here we have the math of LDA. Details can be found in any textbook on machine learning (e.g. **Pattern Classification and Scene Analysis, by Duda and Hart**). The bottom line is that we need to compute two correlation matrices, one rescribing correlations within classes and one correlation across classes. MATLAB loves to do LDA but one needs to be careful; not to ask for impossible, like using too small data sets. MATLAB always uses some tricks to give you the answers but the answers may be meaningless. Of course, this is the problem with most software packages and as a matter of fact with any "knowledge" obtained from the WEB.



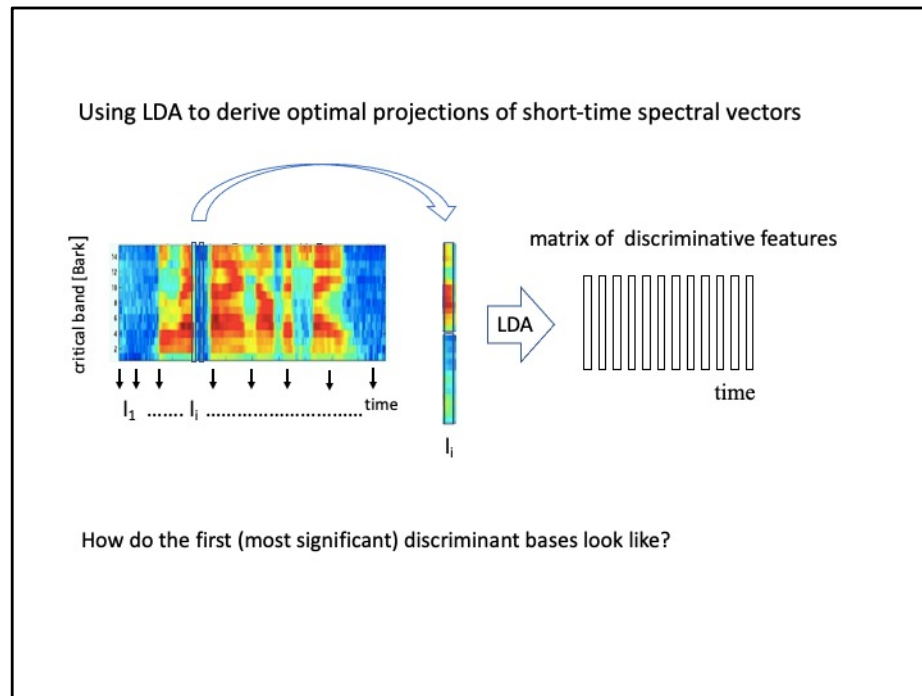
Spectrogram is 2D representation (matrix) of speech signal

Y-axis carries short-time spectra of speech at a given time in each column of the matrix.

X-axis carries evolutions of spectral energies at a given frequency over time in each row of the matrix



The spectral vector in blue is describing sound /j/ so it will carry the label /j/. Once we have the labeled vector space, we can do LDA. The result of the LDA will be the matrix of linear discriminants which can be used for projecting the old short-time spectrum vector space in the directions of the most discriminability on a new vector space.

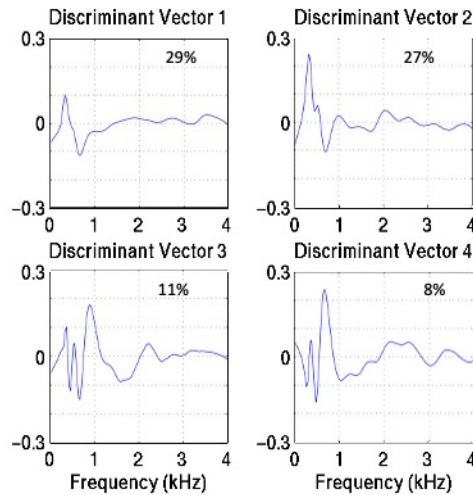


Once we understand that if we have labeled speech data, LDA can be used to find optimal projections of a vector space to discriminate among given classes, we can ask what would be the optimal projection of short-time spectral vectors to optimize their use in classification of speech sounds. When we have labeled data, we know from which sound every spectral vector came, so we have the vector space and the labels with each vector, so we can try LDA analysis.

### LDA-derived spectral bases

(30 hours of continuous telephone speech database – automatic labeling)

Valente and Hermansky 2006



These 4 discriminants represent about 75% of discriminating power. The spectral bases oscillate faster at low frequencies, this implies higher spectral resolution at low frequencies. Not much discriminability is carried beyond the 4-5 discriminants

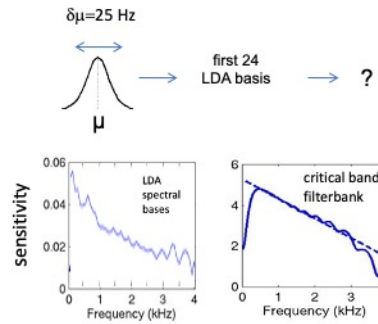
## Soectral nsitivity of LDA-derived spectral bases

(30 hours of continuous telephone speech database – automatic labeling)

Malayath and Hermansky 1998, Valente and Hermansky 2006

### Perturbation analysis

project Gaussian shape on 24 LDA-derived spectral basis and perturb its mean at different frequencies



Similar observations using different optimization techniques

Biem and Katagiri 1994, Cohen et al 1996, Kamm et al 1997, Palival et al 1997, Burget and Hermansky 2001

Spectral sensitivity of the derived LDA projections can be formally evaluated by the so called perturbation analysis. In this technique, a reasonable spectral object (here we used a Gaussian form resembling a resonant peak in spectral space with its mean at a given frequency) is projected on the new vector space. Its position is varied (perturbed) by a certain amount (here  $\Delta f = 30$  Hz constant at all frequencies) and variations in the projected output are evaluated. This yields the sensitivity of the projection at a given frequency. Evaluation of the variationa at different frequencies yields the sensitivity profile of the derived LDA projections.

Perturbation on linear frequency scale show higher sensitivity at lower frequencies.

One can do the same perturbation analysis on the spectrum modified by the auditory-like warping at in the mel spectrum or PLP. Looking at the PLP weigthed spectrm, we observe similar trends as on the LDA-based projections.

Just a few known facts – most of you know by now after out classe on human perception better than me ☺:

Equal changes of pitch require larger frequency changes at higher frequencies

Critical bands of hearing are broader at higher frequencies

Frequency selectivity of cochlea is approximately logarithmic.

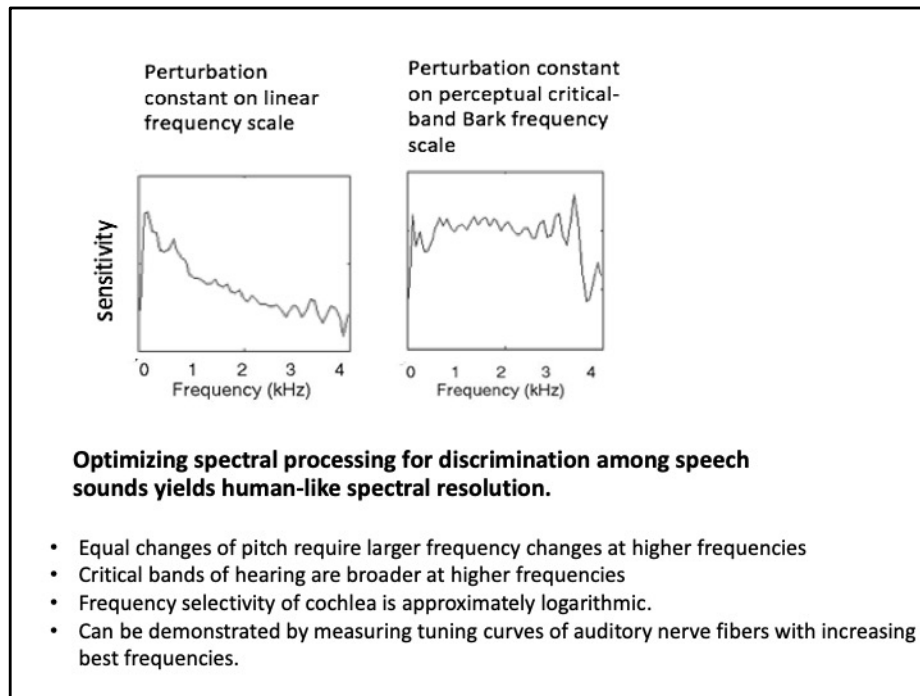
Can be demonstrated by measuring tuning curves of auditory nerve fibers with

increasing best frequencies.

No knowledge of human hearing used, just asking for optimal classification of speech sounds!

Human-like spectral sensitivity is optimal for classification of speech sounds. Known since early 20<sup>th</sup> century!





To double check. When the perturbations are equal not on the linear scale (as in the previous experiment where they were always  $\Delta f = 25$  Hz) but constant on the perceptual Bark frequency scale ( $\Delta f = 0.8$  Bark), the output from the LDA projection is constant, this verifying the lower spectral sensitivity at higher frequency, consistently with the properties of human hearing.

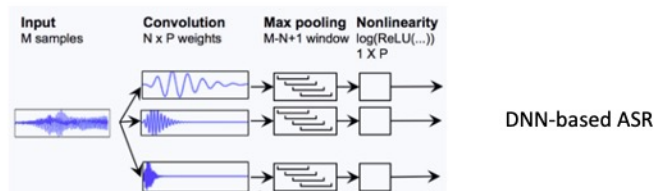
Just a few known facts about human hearing that we mentioned in this class:

- Equal changes of pitch require larger frequency changes at higher frequencies
- Critical bands of hearing are broader at higher frequencies
- Frequency selectivity of cochlea is approximately logarithmic.
- Can be demonstrated by measuring tuning curves of auditory nerve fibers with increasing best frequencies.

Such human-like spectral resolution emerging from just requiring optimal classification of speech sounds.

## Using deep neural net classifiers, filters directly from speech signal

Sainath et al ASRU 2013

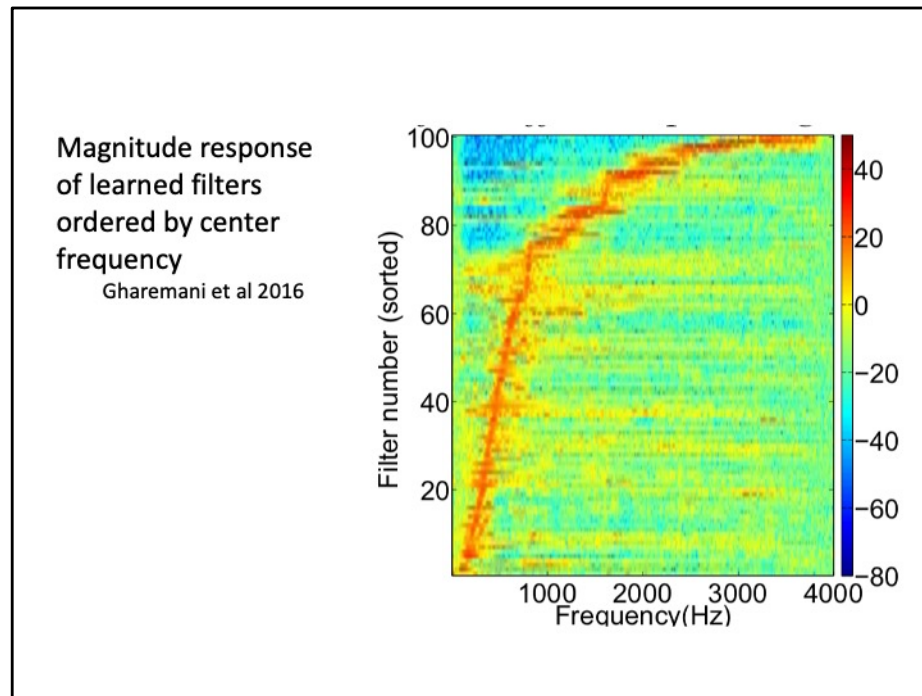


**Q: what are the learned weights in the convolution input layer?**

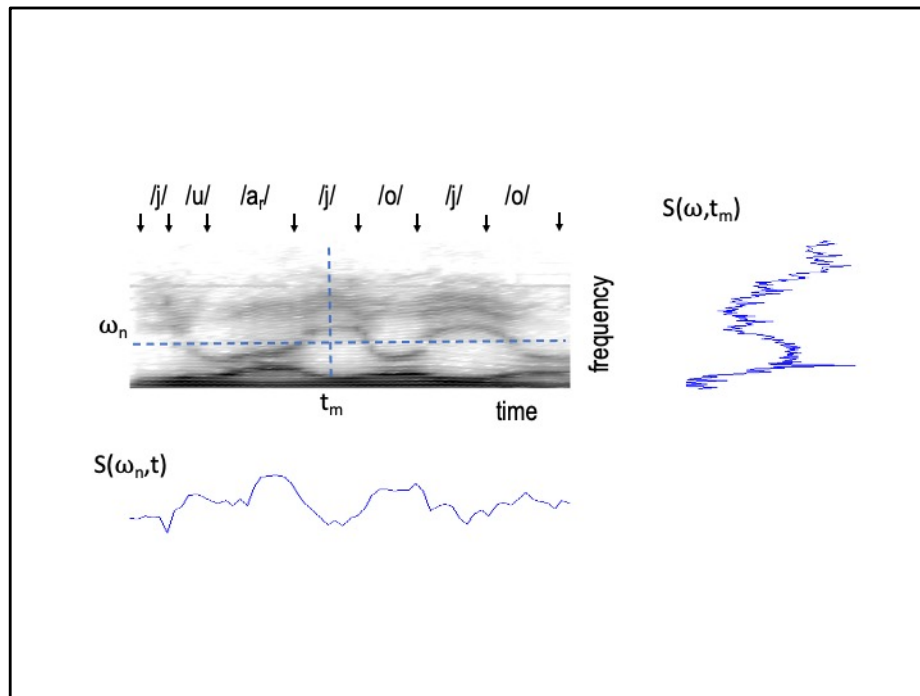
**A: impulse responses of filters consistent with critical bands of hearing**

also Palatz et al 2013, Tieske et al 2014, Golik et al 2015, Gharemani et al 2016, Luo and Mesgarani 2018, ...

DNN architecture can be set up in such a way that on the DNN input nodes are weightings of the data which can be interpreted as the finite impulse response (FIR) filters (signal convolutions for the FIR filters). After the DNN system training, the FIR filters can be examined.



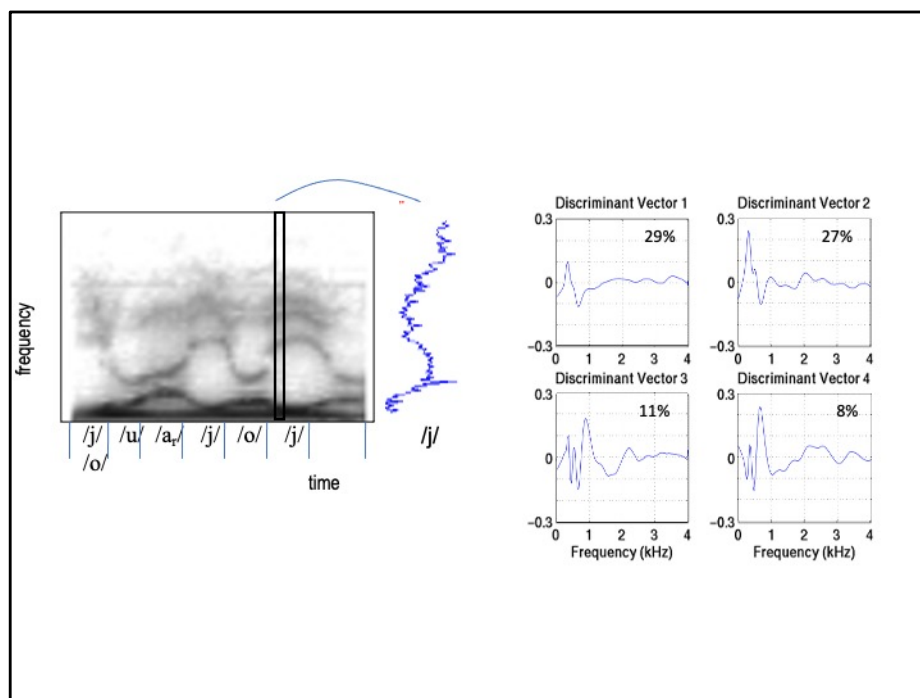
The derived bank of FIR convolutive filters typically show human-like spectral resolution.



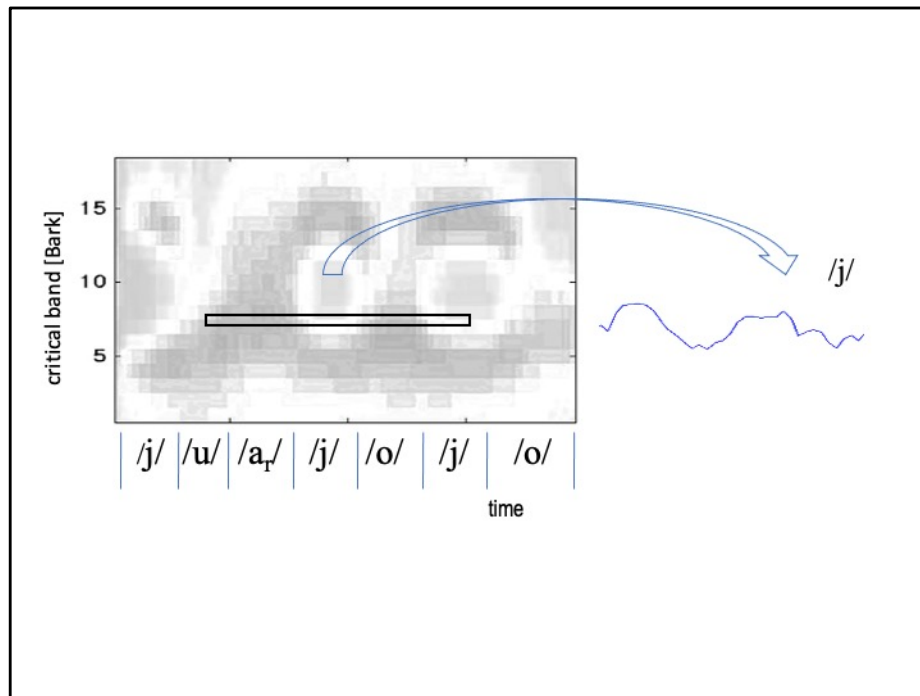
Spectrogram is 2D representation of speech signal

Y-axis carries short-time spectra of speech at a given time

X-axis carries evolutions of spectral energies at a given frequency over time



LDA on short-time spectrogram suggested the use of human-like spectral resolution (critical-band like)



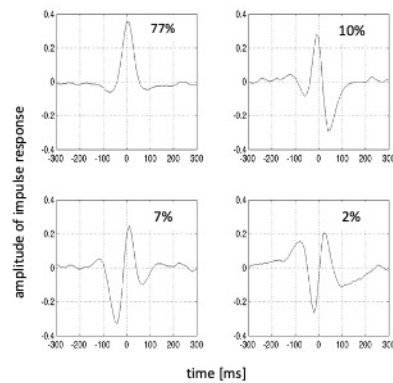
Now we know that we should use critical-band-like spectrum. Using such spectrum, we can re-use the same technique for deriving optimal modulation frequency filters. This time the derived discriminant matrix will correct optimal weighting of neighbouring spectral values to optimize the speech sound classification. For those skilled in DSP – this represents impulse responses of optimal finite impulse response (FIR) filters.

## LDA-derived FIR filters

(30 hours of continuous telephone speech database – automatic labeling)

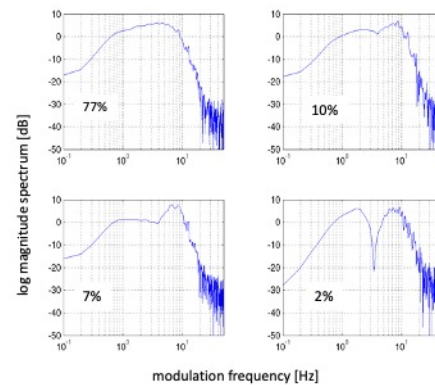
### impulse responses

active parts of impulse responses > 200 ms



### frequency responses

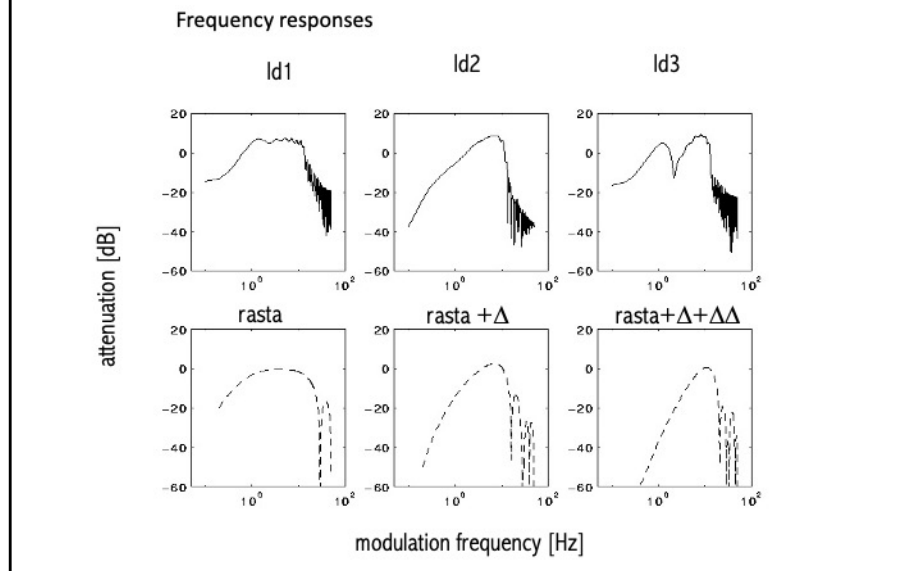
band-pass roughly 1-10 Hz



van Vuuren and Hermansky 1997, Valente and Hermansky 2006

Impulse responses are rather long – more than 200 ms. Frequency responses show that filters suppress very slow modulations as well as modulations faster than around 10 Hz. This is similar to RASTA filters. However, the impulse responses are symmetric, implying the zero-phase linear FIR filters.

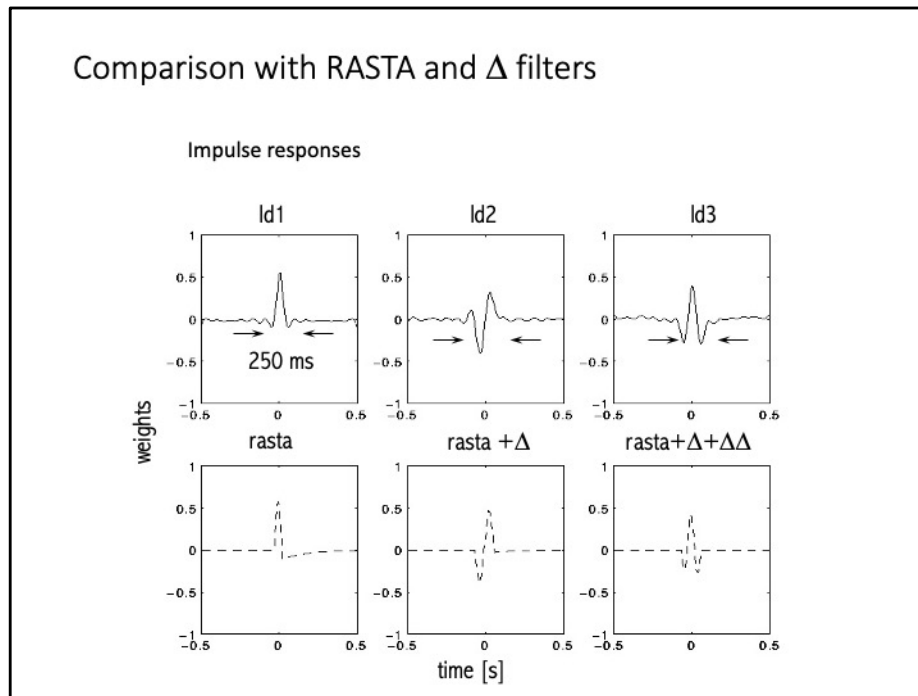
## Comparison with RASTA and $\Delta$ filters



RASTA features were typically more effective when used with the dynamic (differential and double-differential) features. It is interesting that the second and the third discriminant from the temporal LDA filter design resemble the combination of the RASTA and dynamic RASTA features.

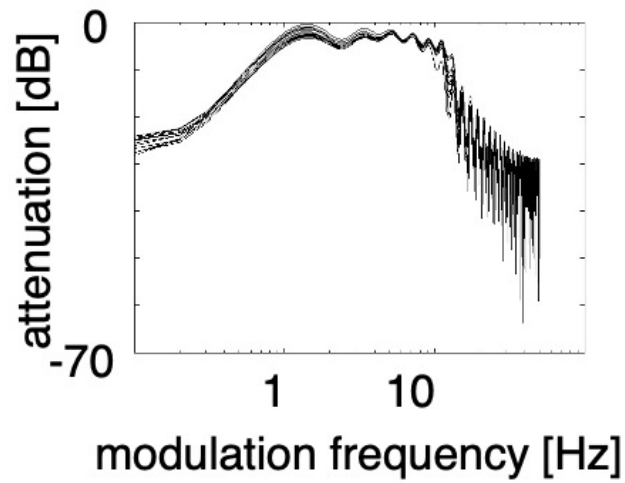


## Comparison with RASTA and $\Delta$ filters



RASTA features were typically more effective when used with the dynamic (differential and double-differential) features. It is interesting that the second and the third discriminant from the temporal LDA filter design resemble the combination of the RASTA and dynamic RASTA features.

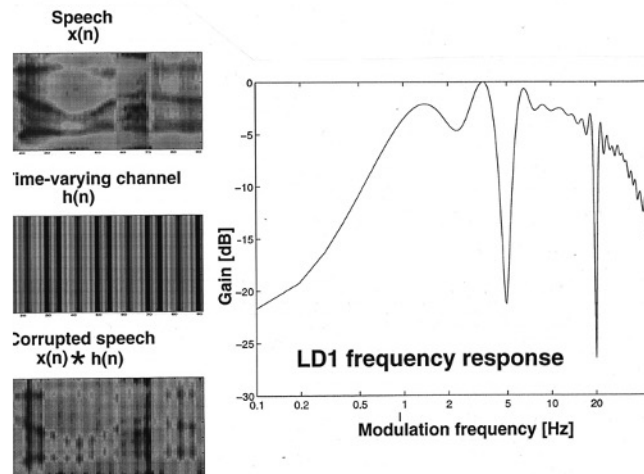
## Dependency on Carrier Frequency



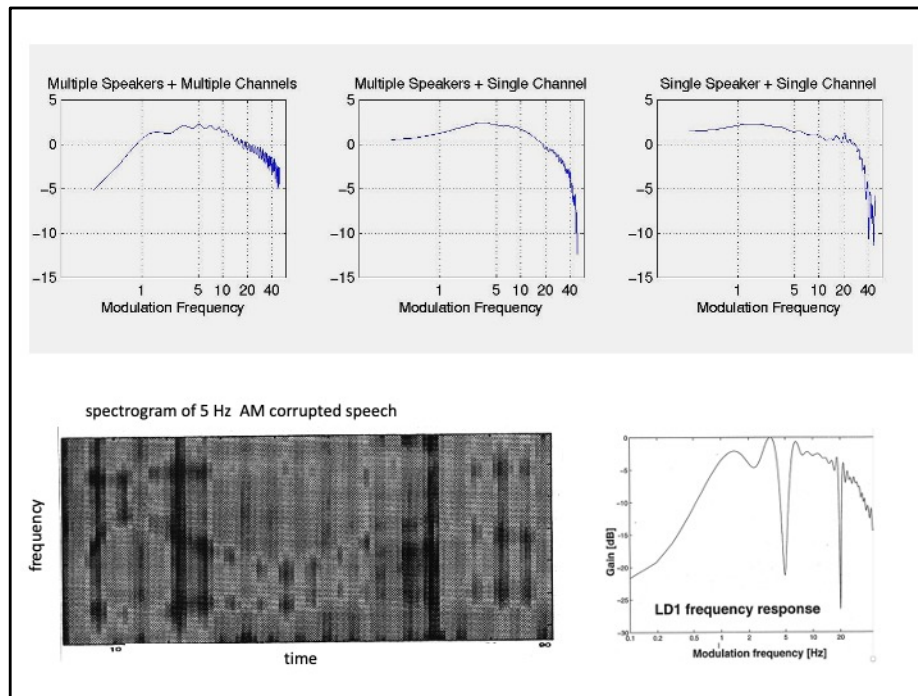
- filters at all carrier frequencies practically identical

An interesting question is if the filter properties are dependent on frequency trajectories (carrier frequencies) to which they are applied. The answer seems to be "no".

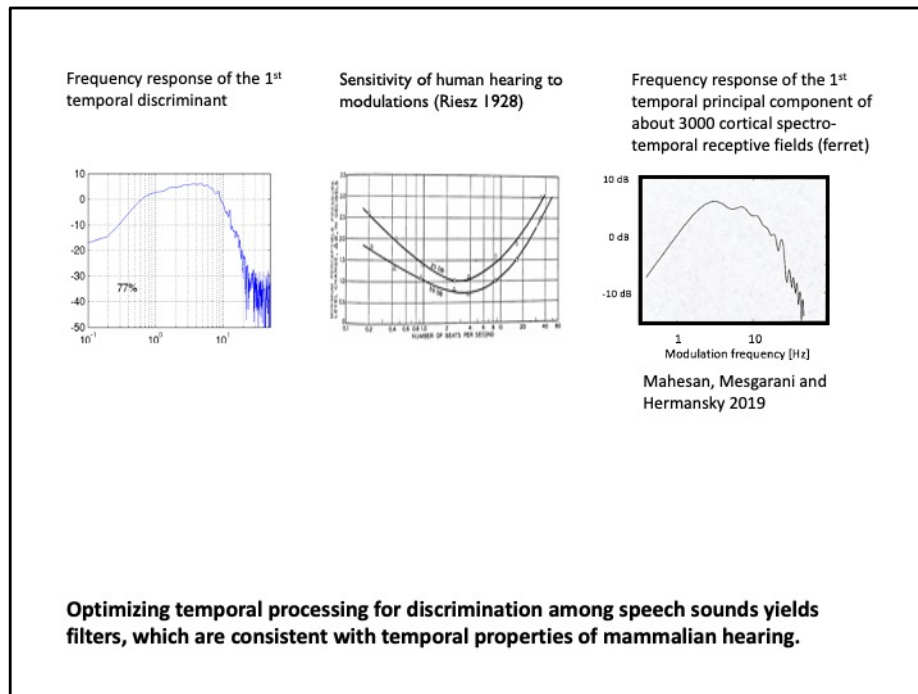
## Filter depends on training data



As expected, the filters reflect the data on which they were derived. When the data were artificially corrupted by amplitude modulating the speech signal at 5 Hz, the filter emerged with the notch at 5 Hz.



However, they are dependent on the type of speech material on which they were derived. For the general case where different speech utterances come from different speakers and different environments, the filters are noticeably suppressing modulation frequencies below 1 Hz. When the channel variability is not present (all speech samples coming from the same acoustic environment) the low modulation frequency suppression is much weaker. When the speaker variability is not present (all files from the same speaker), low modulation frequencies are preserved.



Just a few facts from human hearing:

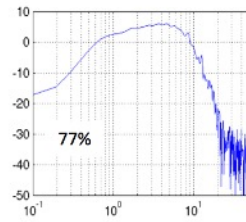
Sensitivity of human hearing to modulations is highest around 5 Hz

-known to speech engineers since 1928

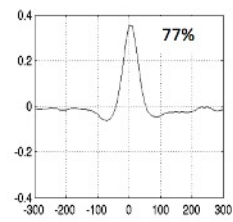
Auditory cortex seems to be the most sensitive to slow modulations.

No knowledge of human hearing was used to derive temporal features directly from speech data, just asking for optimal classification of speech sounds!

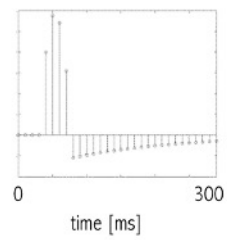
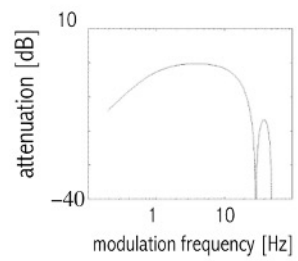
Frequency response  
(peak around 5 Hz)



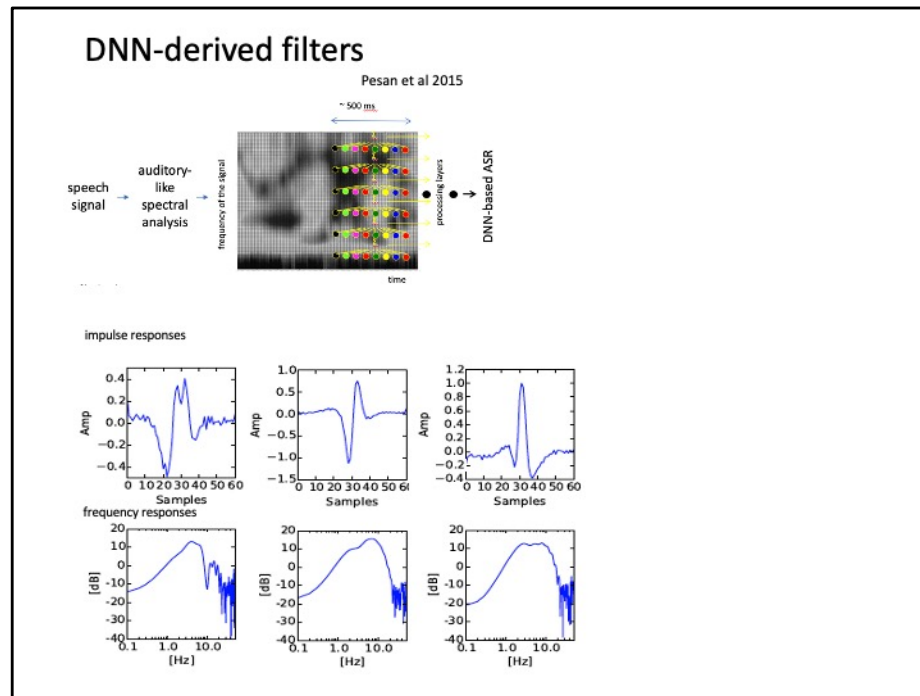
Impulse response  
(effective length around 200 ms)



1<sup>st</sup> LDA filter



Original  
RASTA filter

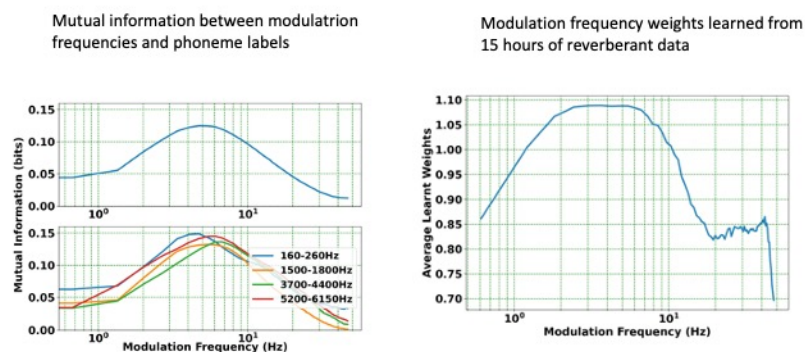


**Impulse responses are rather long – more than 200 ms.**

**Frequency responses show that filters suppress very slow modulations, as well as modulations faster than around 10 Hz.**

## Data-driven design of modulation weights

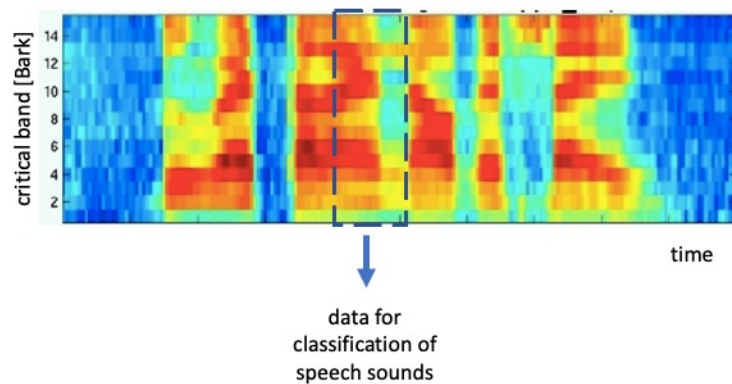
Sadhu and Hermansky, submitted to Interspeech 2022



Mutual information studies and data derived modulation frequency processing can be also studied, conforming relations between linguistic information carried in modulation frequency and optimal modulation frequency weighting derived by optimizing speech recognition system during the machine training.

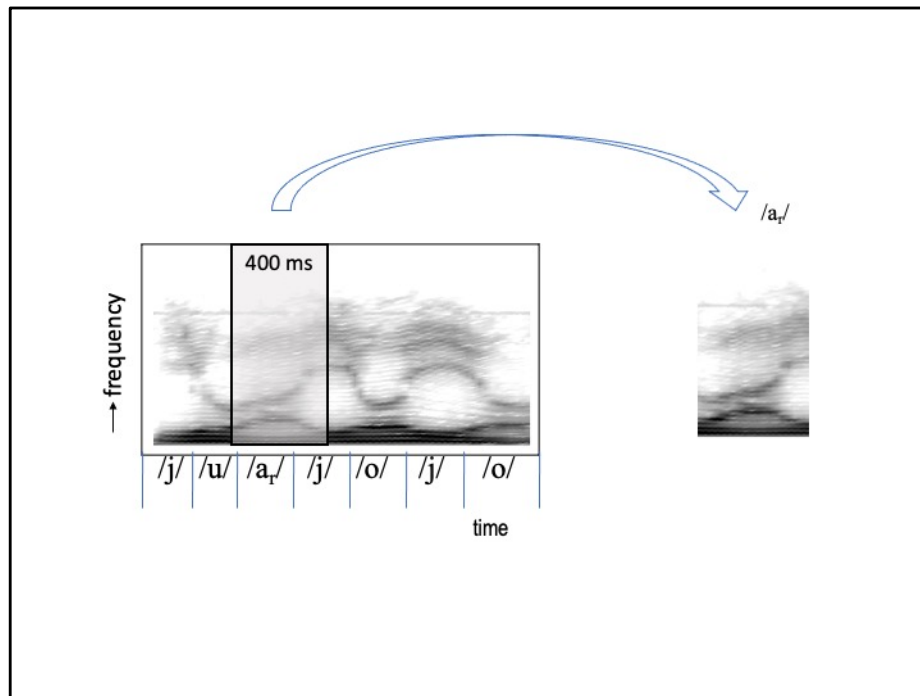


Optimizing for classification of speech sounds suggest critical-band-like spectral resolution and processing within at least 200 ms temporal intervals



Does all that mean that we should use critical-band spectral resolution (of course we all know that) and temporal spans larger than 200 ms (that is what most advanced ASR systems started to use now).

Should we probe even further for appropriate signal pre-processing?

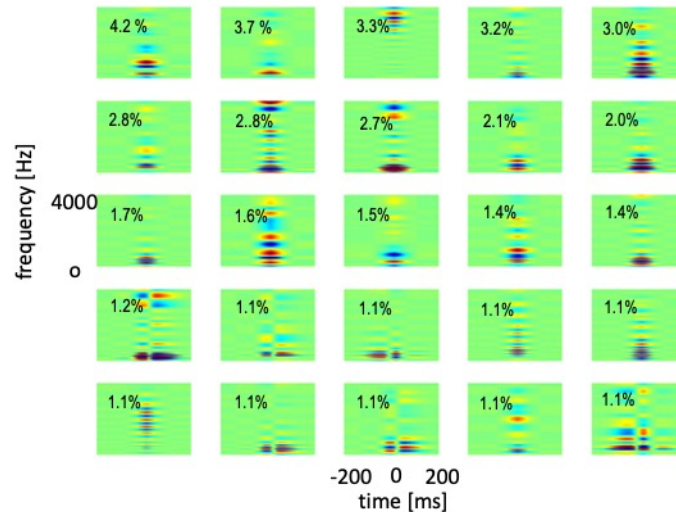


Using the same LDA technique on larger segments of LPC-smoothed spectra should give 2-D spectro-temporal (cortical-like??) projections for next stages of phoneme classification.

## 2D time-frequency discriminants

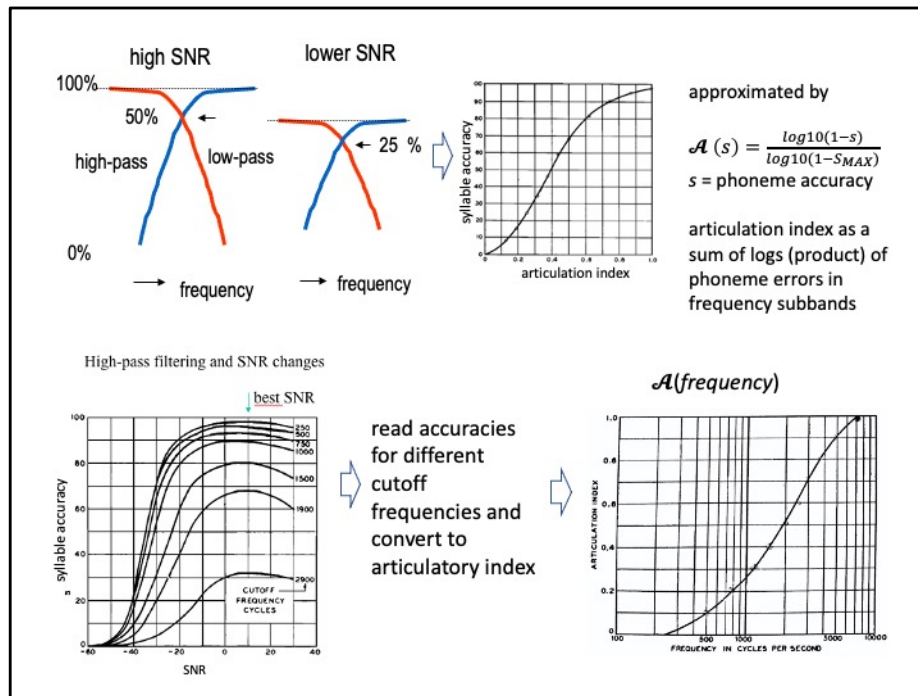
Valente and Hermansky 2006

Many 2D discriminants are frequency-selective, emphasizing particular parts of speech spectrum



Spectral resolution is typically coarser at higher frequencies. Temporal spans of filters typically cover more than 200 ms. We knew that from our earlier LDA spectral and temporal optimizations.

However, most interestingly, different filters often focus on different parts of speech spectrum !! This suggests that the optimal set of time-frequency derived features would be the features which look at different temporal spans of speech signal (up to 200 msec) and often do not cover all frequencies, each feature focusing only on some frequencies of the speech spectrum.

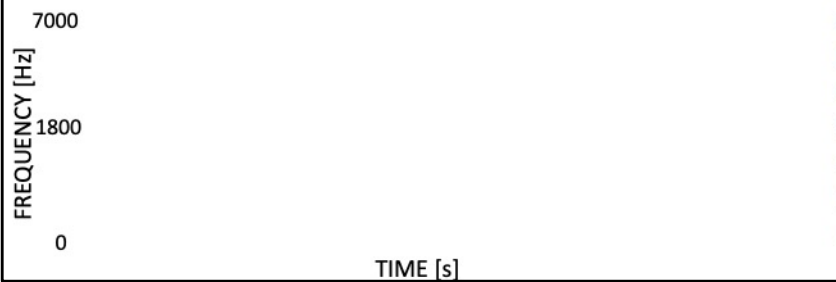
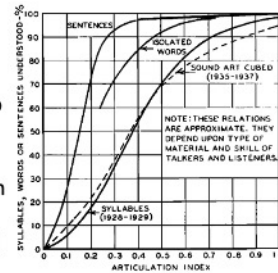


## Articulatory Bands

French and Steinberg 1949

250-375-505-654-795-995-1130-1315-1515-1720-1930-2140-2355-2600-2900-3255-3680-4200-4860-5720-7000 Hz

- 20 frequency bands in speech spectral region
- each band with SNR > 30 dB contributes equally to human speech recognition
- bands with SNR < 0 dB do not contribute at all
- any 10 bands sufficient for 70% correct recognition of nonsense syllables, better than 95% correct recognition of meaningful sentences [Fletcher and Steinberg 1929]



This turned out to be true not only for two bands, but also for more bands (up to 20 bands). Thus, the multichannel model predicts that the total error will be given by

$$e = \prod_{i=1}^K e_i$$

- SERIAL PROCESSING

- phonemes in a nonsense syllable are decoded independently of each other  $S=c.v.c$  (**probabilities of correct recognition multiply**)

- PARALLEL PROCESSING

- errors in phonetic judgment in nonsense syllables in individual sub-bands are independent (**probabilities of errors multiply**)