



Received 20 June 1969

Whither Speech Recognition?

J.R. PIERCE

Bell Telephone Laboratories, Inc., Murray Hill, New Jersey 07971

9.10, 9.1

Letter to Editor
J.Acoust.Soc.Am.

Research field of "mad inventors or untrustworthy engineers"

Funding artificial intelligence is real stupidity"

- supervised the Bell Labs team which built the first transistor
- President's Science Advisory Committee
- developed the concept of pulse code modulation
- designed and launched the first active communications satellite

.... should people continue work towards speech recognition by machine ? Perhaps it is for people in the field to decide.

To implement ASR, we need to apply ***intelligence and knowledge of language comparable to those of a native speaker !***

A short letter to editor of the Acoustical Society of America from the very influential researcher at Bell Labs almost stopped speech recognition research in USA. Read the letter by yourself, I believe that Dr. Pierce had some good advice, still valid even today.

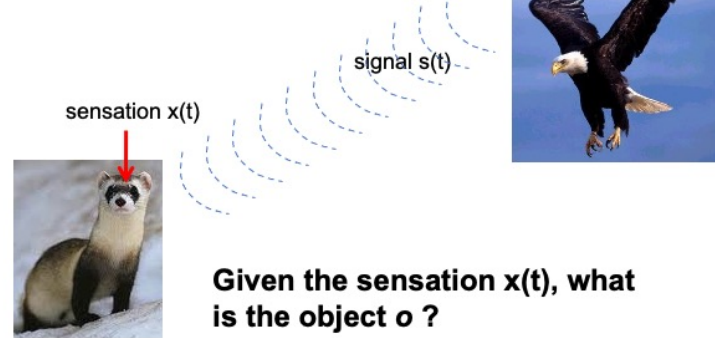
Since 1969

- Better speech features (linear prediction, cepstrum, auditory-like techniques...)
- Better pattern matching (dynamic time warping, Viterbi search)
- Stochastic models allowing for using huge amounts of speech data
- Iterative expectation-maximization (EM based training only from transcribed speech data (no need for data labeling))
- Explicit use of Bayes rule combining the evidence from the signal together with prior expectations from the language

2

The recognition field fortunately did not stop altogether and gradually it recovered to the point where we are today. Several reasons for its recovery are listed here.

How to survive in this hostile world?



This is a situation, which is often faced in nature. Ferret's perceptual system receives a signal $s(t)$. The signal may be multimodal (audio, visual, olfactory, tactile,..). The signal causes sensation $x(t)$ in the perceptual system of the animal. The task, which our ferret need to perform is to find out what the sensation $x(t)$ represents. In tis case, our ferret is concerned that the stimulus and the resulting sensation might have been produced by an eagle o .

Given the signal $x(t)$, what is the stimulus $s(t)$?

$s(t)$ – e.g. changes in acoustic pressure representing the sound

$x(t)$ – e.g. the activity in the auditory system

$s(t)$ comes from some probability distribution $P[s(t)]$

- different stimuli have different probabilities of occurrences

$x(t)$ occurs with the conditional probability $P[x(t)|s(t)]$

- $x(t)$ depends on $s(t)$ but the response is not unique (system noise?)

Need $P[s(t)|x(t)]$ (probability that $s(t)$ happened when the activity in the system is $x(t)$)

$P[x(t),s(t)]$ – likelihood that both $s(t)$ and $x(t)$ happen at the same time

- have we experienced $x(t)$ when $s(t)$ happened ?

$$P[x(t),s(t)] = P[x(t)|s(t)] P[s(t)] \quad \text{or} \quad P[x(t),s(t)] = P[s(t)|x(t)] P[x(t)]$$

$$P[s(t)|x(t)] = P[x(t)|s(t)]P[s(t)]/P[x(t)]$$

(Bayes Rule)

Bayes rule

$$P[s(t)|x(t)] = P[x(t)|s(t)] P(s(t)) / P[x(t)]$$

$s(t)$ - incoming stimulus that describes the object

$x(t)$ - activity in the system resulting from the stimulus

our ferret needs

$P[s(t)|x(t)]$ - probability of the stimulus given the data

our ferret learns

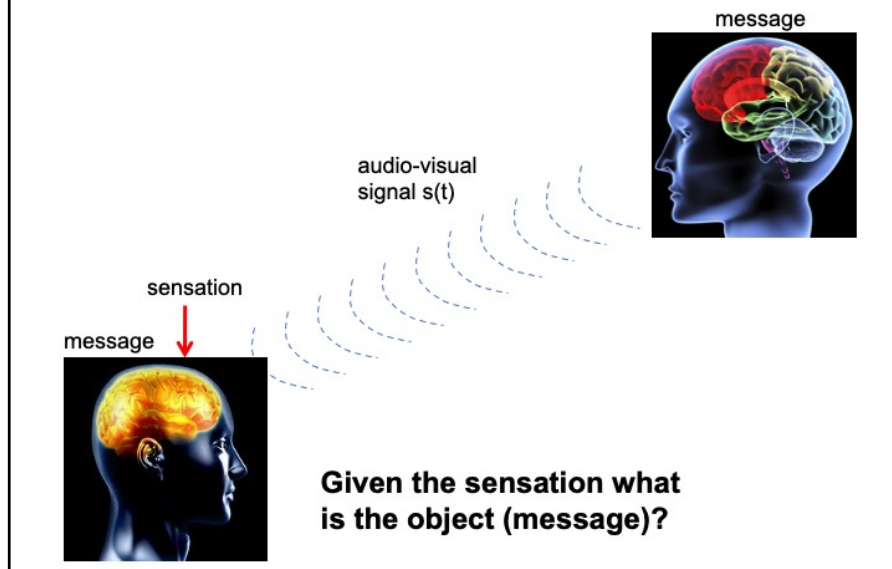
$P[x(t)|s(t)]$ - probability of the data given the stimulus

$P[s(t)]$ - probability of the stimulus

$P[x(t)]$ - probability of the data

This is just to summarize what we have talked about before. What our ferret needs is to estimate what is the probability $P[s(t)|x(t)]$ of the eagle when hearing the eagle's cry $x(t)$. The ferret knows from its prior experience what is the probability of eagle cry given the eagle is in the sky $P[x(t)|s(t)]$, probability of hearing eagle in the given context $P[s(t)]$, and also what is the probability of hearing the eagle cry $x(t)$.

How to survive in this hostile world?



Humans developed different acoustic means to survive, the communication by speech. Using speech, besides an obvious means for issuing warnings (many animals can do the same using sounds), an amazing amount of information can be exchanged. In particular, we believe that speech is unique among communication means since it is able to be used to communicate abstract concepts (truth, justice,...) or to refer to past or future. It is tempting to say that speech is what differentiates humans from other species.

Stochastic machine recognition of speech

$$P(M, x) = P(M|x)P(x) = P(x|M)P(M)$$

Joint probability that message M and data x occur together is given by the probability of the message M given the data x multiplied by the probability of the data x , or by probability of the data x given the message M multiplied by the probability of the message M .

Bayes rule

$$P(M|x) = P(x|M)P(M) / P(x)$$

To find the maximum of the probability, the probability of the data is not need

$$M = \operatorname{argmax} P(x|M_i)P(M_i)$$

How to get good model M ?

How to find the optimal M ?

What is the data x ?

Formerly, we express the joint probability of the particular sequence of models (the message M) and the observed data. Joint probability that message M and data x occur together is given by the likelihood of the message M given the data x multiplied by the probability of the data x , or by probability of the data x given the message M multiplied by the probability of the message M . This yields so called Bayes rule, which expresses the likelihood of a given message M given the data x . Since the probability of data is just a scaling factor (does not depend of the message M) it can be ignored. So to find the message M which most likely generated the data x , we need to search all possible messages which could be generated by the model and keep the one which yields the highest likelihood.

Easier said than done, since the model M may not be correct, the acoustic model likelihoods (learned from the acoustic training data) may be inaccurate,, the prior probabilities of particular messages (the language model trained on text data) may be also incorrect,

the number of possible messages is huge and searching through all of them is not easy, and the features x derived from speech signal may not carry the required information and may carry information about a number of irrelevant information sources.

How to get good models ?

The model architecture and **DATA !**

$$M = \operatorname{argmax} P(x|M_i)P(M_i)$$

There is no data like more data
attributed to Bob Mercer

More data is always better than
more thinking
attributed to Eric Brill



No knowledge is better than wrong knowledge

To get good model of speech, we need good model architecture and good data for the model training. The architecture largely stabilized as a sequence of speech sounds. What represents the is dependent on amount of training data which is available for the model training. More detailed models of sounds require more data. Similarly, more complex language models may require more data. So to get better models needs more data,

How to find the best $M(w)$?

Search through all possible models $M(w_i)$ to find which most likely produced x

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

Y – “fudge factor” (makes all statisticians uneasy)

The search is a form of dynamic time warping where templates are represented by stochastic models. Where distances are replaced by likelihoods of $P(x|M)$ combined with priors $P(M)$

$P(x|M)$ from acoustic data (labeled by speech sounds)

- trivial when sound labeled training data
- not labeled but only transcribed
 - know the sound sequences, find boundaries through iterative expectation-maximization techniques (remember the line of boys and girls)

$P(M)$ – sequences of words with probabilities derived most often from texts

- Problems for low or zero probability words!

The search for the best message involves two stochastic models, the acoustic model $P(x|M)$ which is trained using labeled speech acoustic data, and the language model $P(M)$, typically trained using the text data. Relative contributions of these two models is controlled using the “fudge factor”

Y , which is chosen experimentally for the best performance on some development data. Notice that messages for which the $P(M)$ is zero cannot ever be chosen. Thus, the words which are the most likely in the language are also the most likely to be selected. Remembering that from the information point of view, the less likely items carry more information, we see some inconsistency.

Almighty language model

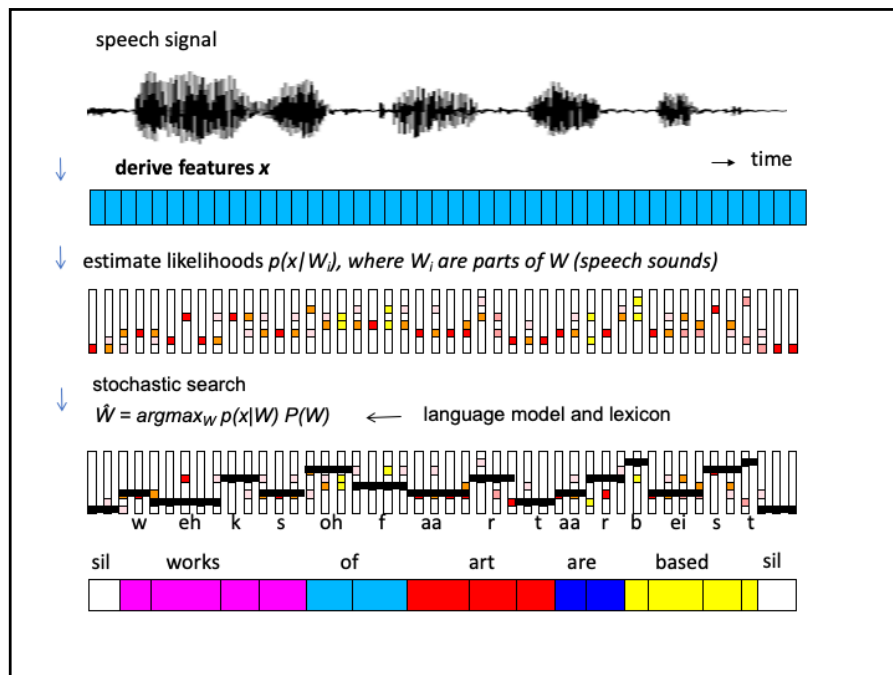


Koupil jsem si nový **computer**, který nefunguje.



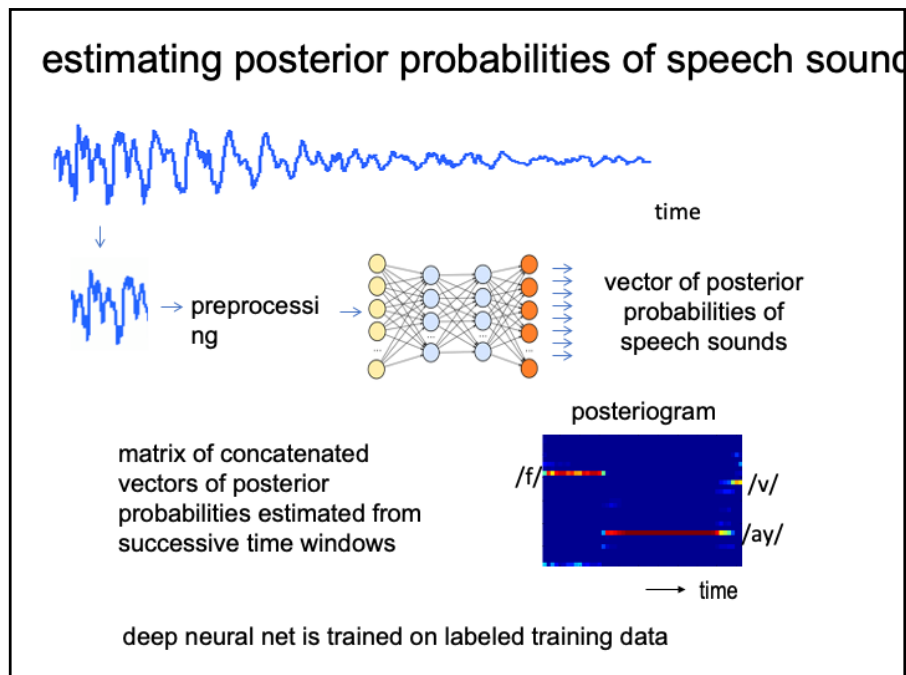
although some sort of the **computer** can either way
hopefully cin-cin o-bi **computer** connected with

Language model can severely limit what can be recognized. This is obvious when you try to recognize a sentence in a language which your recognizer does not expect. Two outputs from two different state-of-the-art systems available today in response to Czech language are shown here.

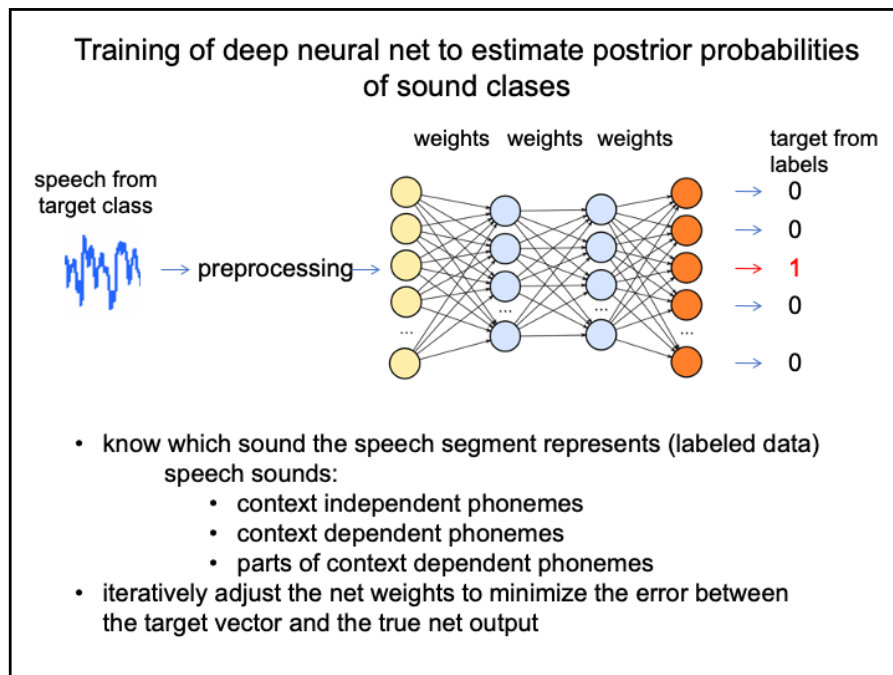


This scheme shows schemitacally the HMM recognition process. Short speech segments (what is “short” may be discussed later) are described by the segment features x . Likelihoods of speech sounds given the features are derived using trained acoustic models. Search for the best path though the sequence of likeligood vectors with further help of the language model, which contains preferences for a particular word combinations, yields the sequence of speech sounds (phonemes). Having this sequence, pronunciation rules allow for transcribing the unknown utterance/

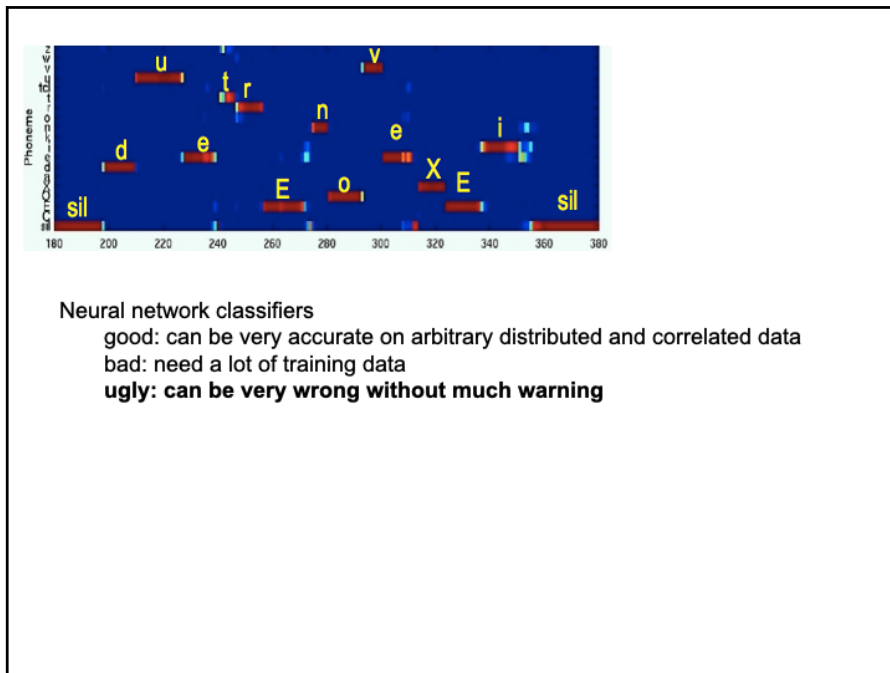
Deep Neural Nets in Speech



DNNs can be used for deriving posterior probabilities of speech sounds. Input is a set of some features describing the information in the short speech segment (think short-time spectrum of the windowed speech segment), on the output of the trained net we obtain estimate of posterior probabilities of all speech sounds in which the DNN was trained.



For their training we need again labeled speech data. The DNN is trained to deliver highest output for the sound which underlines the given segment of the speech signal. This is done using the “one-hot (that is one output is close to 1, others are close to 0) vectors as the target of the DNN during the training. The weights of the DNN are iteratively adjusted during the training to achieve this as well as possible.










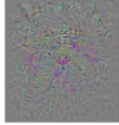




Well trained net on good clean data can estimate posterior probabilities of speech sounds amazingly well. As inputs DNNs do not require uncorrelated and Normally distributed data, the main requirement typically is that the data are “whitened”, i.e. mean and variance normalized. For a good performance, DNN need significant amounts of training data. This also implies that during the inference (during recognition) they can be very sensitive to any deviations from the statistics of the training data.



CLASSIFICATION

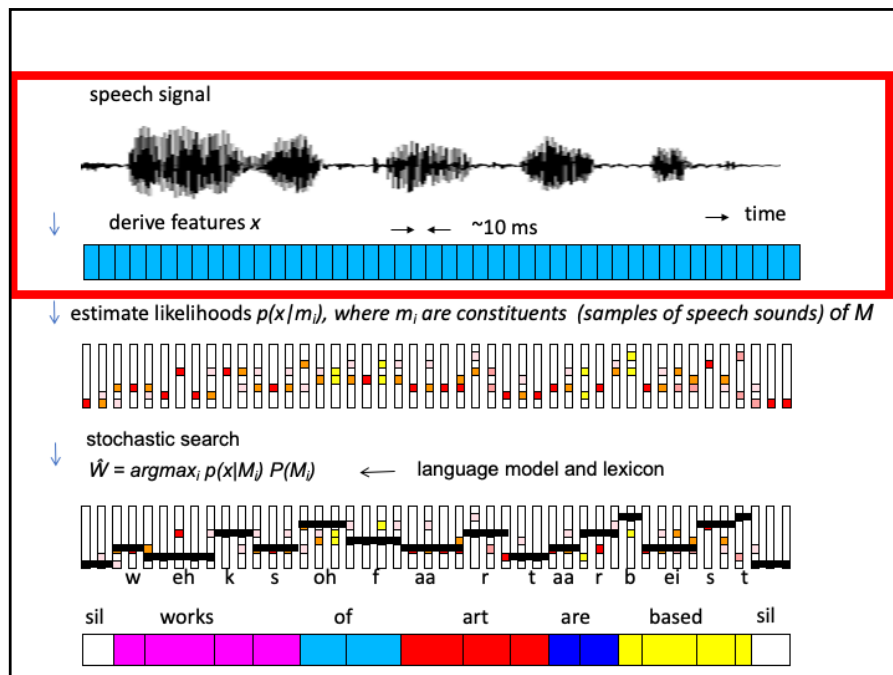
Szegedy et al 2014

	SOAP DISPENSER			OSTRICH	
	MANTIS			OSTRICH	
	DOG			OSTRICH	

He who knows not and knows not that he knows not is a fool; avoid him.
He who knows not and knows that he knows not is a student; teach him.
 He who knows and knows not that he knows is asleep; wake him.
 He who knows and knows that he knows is a wise man; follow him.

How bad can DNNs be in classification? This is nicely illustrated in examples of classification of still images. When the signal is clean, the classifier works perfectly, **soap dispenser**, **praying mantis** and **dog** are classified as such.

When a small amount of noise, not perceived by human observers at all, is added to the images, all images are classified as **ostrich**. No surprise, this bothers at least some researchers.



We have spent time discussing how to derive a good set of features x for ASR. As long as the likelihood estimates were derived by the quadratic diagonal-covariance Gaussian mixtures classifier, the features x needed to be approximately Normally distributed and uncorrelated. That limited the techniques which were used for their estimation. There is no such a constraint for the deep neural net (DNN) classifiers.

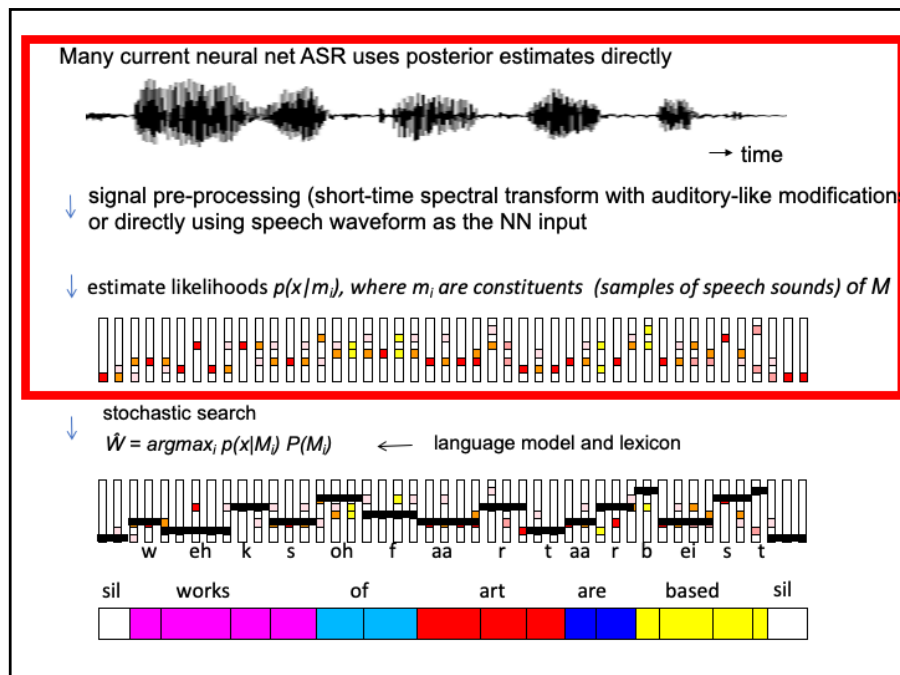
DNN-HMM hybrid

- Convert posterior probabilities to likelihoods (divide by training priors) to be used in Viterbi search for the best word sequence

Bourlard and Morgan, NIPS 1990

$$p(c_i | \mathbf{x}) = \frac{P(c_i | \mathbf{x})}{P(c_i)}$$

One way of using posterior probability estimates in ASR is to forget about the classical Gaussian classifiers used for the sound class likelihood estimates in HMM paradigm since the seventies of the last century and use the posterior probabilities from the DNN directly. To convert the posterior probabilities to likelihoods for the HMM search, we can divide each posterior estimate by the class prior (which is known from the training data). Today is sounds easy and typically works well but in the early days of DNNs, estimating probabilities of large sets of sound classes which were used in ASR, was not easy. Additionally, GMM-HMM recognizers were well developed and it was not easy to abandon them.



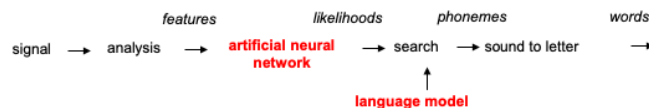
Here is the simplified ASR system, which directly uses estimates of likelihoods of speech sound from the DNN.

Speech recognition since 1990

Artificial neural nets (ANNs) to estimate probabilities for statistical pattern recognition

Morgan and Bourlard 1990

ANN/HMM hybrid system



Can use huge numbers of trainable parameters (regularization, cross-validation)
 Can form complex class boundaries
 Hierarchical layer-by-layer processing

Accept correlated and arbitrarily distributed features

From cepstrum to concatenated spectral vectors (increasing context)
 Extreme case of TRAP (1 sec long individual temporal trajectories of spectrum)

Making friends with HMM community

TANDEM features for HMM (gaussianized and decorrelated posteriors from ANN)

Late eighties and early nineties brought back the Rosenblatt's perceptron, this time with nonlinearities on each processing layer and with straightforward adaptation of net weights.

Initially, because of the lack of data and computing resources, it was difficult to use large ANN models and the nets were competitive only in simple architectures such as single-state phoneme model based ASR.

However, the ANN advantages such as the ability of forming a complex class boundaries, hierarchical layer-by-layer processing, and finally the ability to deal with wide class of input features, were always there.

It encouraged us to introduce larger and larger temporal contexts, to abandon the cepstrum based envelope for the individual spectral energy trajectories, sometimes as long as 1 second,

ANN people also contributed to the competitive GMM/HMM systems by providing them TANDEM features derived by ANN, which were forced by ETSI (European Telecommunications Standards Institute) requirement of using HMM/GMM based classifier in the competition.

Recent developments

Many capable people work with increasingly larger (and deeper) nets training on **very large datasets** using powerful hardware.

From **locally optimized modules** (features extraction, likelihood estimation, time alignment, sound-to-letter, ...) to **global optimization** of the whole system.

signal x → **deep neural network** → words

Very powerful paradigm, but still

- needs labeled data,
- assumes that present is similar to past (generalization),
- prefers recognition of **more likely items** rather than **information rich rare items**.

Ever increasing amounts of training data may obviate true progress and disadvantages less "fortunate" groups.

Work towards decreasing the necessary amounts of training data !

I would like to stop this brief and necessarily incomplete history description at the point where most papers in this conference will probably be, that is at the current explosive development in use of ANN in their current deep-learning form.

Many very competent and capable researchers are introducing complex huge net architectures trained on unbelievable amounts of training data using more and more powerful hardware.

The most important advance is that the module-by-module optimization is gradually replaced by the global optimization.

The end-to-end approaches take the input data, sometimes even in the original raw signal form and directly deliver a sequence of letters describing the message.

Still, some issues from the early days of ASR remain: the need for some labeled data is here, the assumption that the operation domain is similar to the domain on which the system was trained,

and the preference of frequent and therefore expected words over the information rich infrequent words is still here.

The need for training data is still exponentially increasing and the system performance is increasing with it.

We may not know anymore if the improvements come from some fundamental improvement of the underlying principles or from more data.

~ 300 days of labeled data

~ 100 years of unlabeled data

Parthasarati et al 2019

137 billion free parameters

100 billion words in training

Shazer et al 2017

What should be the signal x ?

No knowledge is better than wrong knowledge
but **more we know, less we need to learn.**

We speak in order to be heard and need to be
heard in order to be understood.

Roman Jakobson



Features x (early decision making)

Alleviated relevant info is lost forever, irrelevant info that is left in may
create problems during the inference.

**Emulate some relevant properties of human hearing in
feature extraction ?**

More we know. less we need to learn.

And what we know for sure is that speech is

perceived through
hearing. So we initially
decided to emulate some
basic properties of
human hearing in the
feature extraction
module of ASR systems

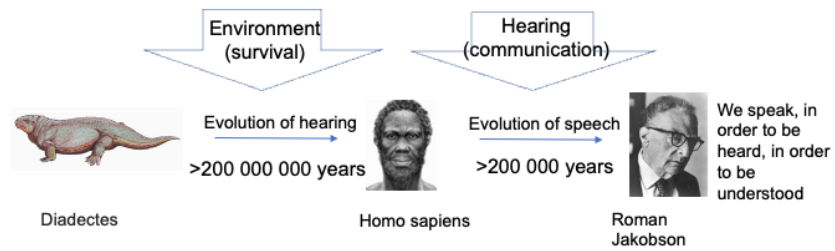
The feature module was one of obvious
points in the system to start with.

Features need to represent linguistic
information about message in speech and
not irrelevant information (speaker
specific, acoustic environment,...)

whatever is alleviated during feature

computation, is lost forever
whatever is left in, makes the
recognition (training) more difficult

Evolution of speech



If speech evolved to be heard, then knowledge of hearing which is relevant to speech processing is imprinted on speech.

Use speech data to derive properties of hearing which are relevant for speech processing!

Large amounts of labeled speech data are becoming available. It would be a shame to use it only for training of ASR algorithms.

Should Airplanes Flap Wings ?

"Airplanes do not flap wings but have wings nevertheless,....."

Of course, we should try to incorporate the knowledge that we have **of hearing, speech production, etc.**, into our systems,.....but we need to estimate the parameter values from the data. There is no other way

F. Jelinek, Five speculations (and a divertimento) on the themes of H. Bourlard, H. Hermansky, and N. Morgan, *Speech Communication* 18, 1996, pp. 242–245

We have seen that emulating some basic properties of hearing such as nonequal spectral resolution of relative insensitivity to unchanging spectral patterns can be useful in ASR.

However, there is always a question which properties are relevant for speech and which serve other hearing purposes (or no purposes at all, such as appendix).

Fred Jelinek correctly suggests that properties that are emulated should be supported by speech data.

Indeed, large amounts of labeled speech data are becoming available. It would be a shame to use it only for blind training of ASR algorithms.

Spectral Envelope

Speech Production



27

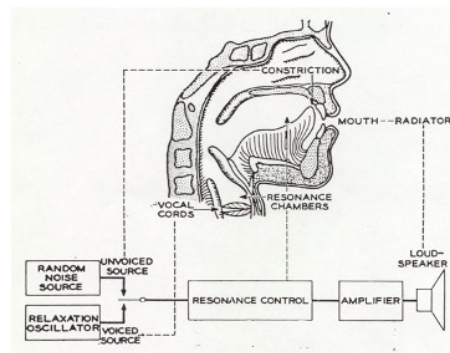
The messages are carried in movements of vocal tract of a speaker. These movements would not be heard but humankind invented means for making them to be heard by exciting the vocal tract by sound sources with rich audible spectrum and relative flat spectral envelope - the vibration of vocal cords, air frictions in narrow passages in the tract (including the non-vibrating vocal chords in production of whispered speech) or sound pulses in release of air in plosive sounds. These rich source spectra are modulated by changing transfer function of the vocal tract. Thus the message information is carried in the changes of the audible spectra of speech.

Message is carried in changes in vocal tract shape,
which modulate spectral components of speech

H. Dudley: The Carrier

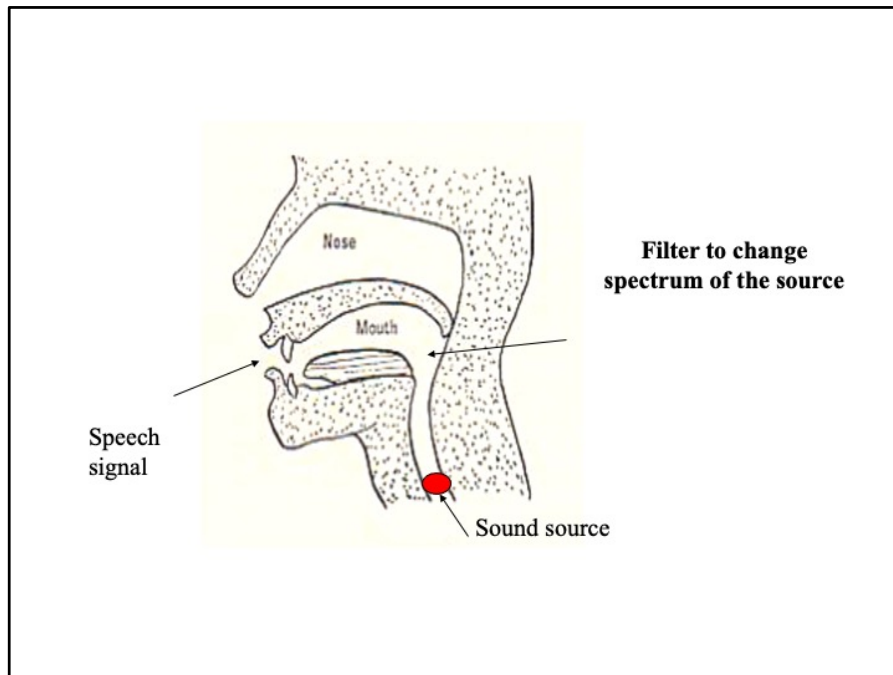
Nature of Speech, BSTJ 1940

Homer Dudley

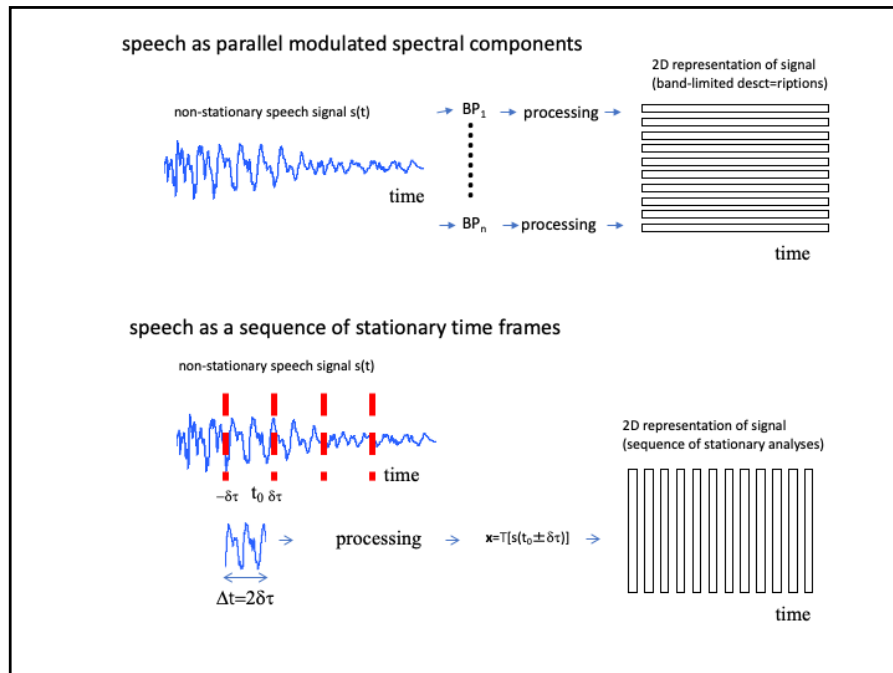


One important concept was the concept of carrier nature of speech. It says that the bulk of the information about the message in speech is in **changes of vocal tract shape**, i.e. in slow movements of

vocal organs such as tongue or lips. These movements are made audible by exciting the vocal tract by the source which rings in overtones (vibration of vocal cords, air friction in tract constrictions, sudden release of air after the constriction release, ,,,). Notice there was no mentioning of vocal tract resonances (formants) in Dudley's concept.

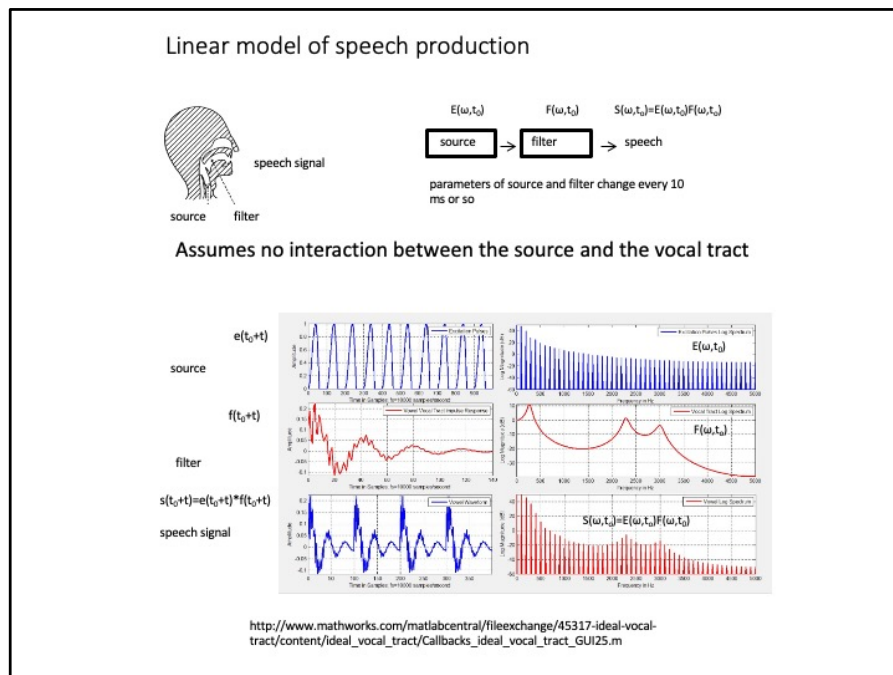


Two basic elements of the vocal tract are the **tract cavities** which filter the spectrum of the **sound source**. So far, all was fine. Spectrally rich source signal gets modulated by changing filtering properties of the tract shape. However, mathematical analysis is process in not easy. Some simplifi=yng assumptions may be needed.



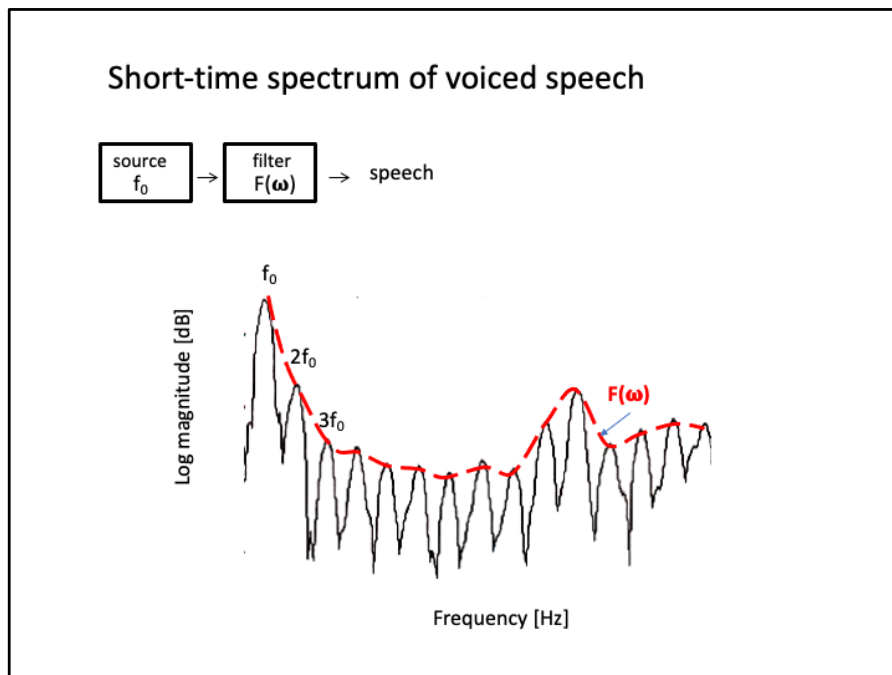
In the original concept (Homer Dudley 1940) the speech signal was seen as a number of parallel band-limited signals modulated by different parts of the changing frequency response of the moving vocal tract.

As digital signal processing techniques started to dominate speech processing, the concept gradually shifted to a time-frame processing where time segment of speech are short enough so that the signal can be considered to be stationary. Parameters are extracted from these segments independently of the other neighbouring segments. The 2-D time-frequency representation of the signal is a sampled representation. The frame rate of the 2D representation implies the highest frequency of change preserved in this representation.



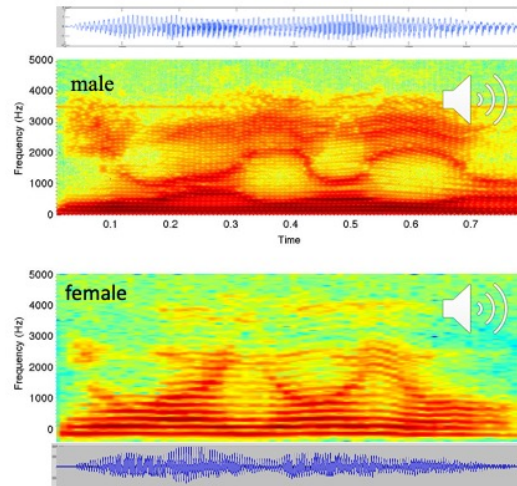
In early sixties of the last century in works of Gunnar Fant in KTH Stockholm and Kens Stevens at MIT studied extensively the nature of both the source and the acoustic filter under the assumption that there is no interaction between them (hence the linear model). Some interaction of course exists - the filter gets damped when the glottis in the source opens and the lung cavity gets into the action but this is mostly neglected. Further, the whole system is assumed to be stationary for short segments of time (10-20 msec) – another assumption which is not entirely correct. However, these assumptions allow for cleaner mathematical formulation of the speech production.

Digital techniques for short-time Fourier analysis sealed this mathematically clean but rather approximate concept, which mostly remains in speech processing until these days. Speech signal is first chopped into short speech segments, each segment is analysed independently of other segments, and these discretized analysis results are describing the dynamic speech signal.

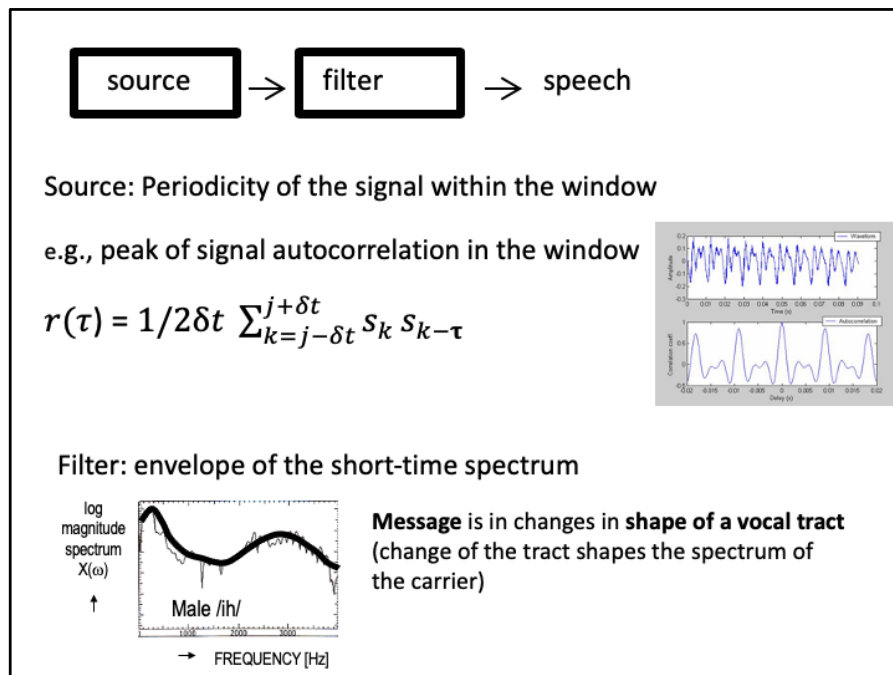


Spectrum of voiced speech consists of spectral envelope $F(\omega)$ which characterize the acoustic filter and the fine structure which characterizes the speech source. The filter has spectral peaks where the vocal tract resonates (formants). The fine structure has peaks spaced in integral multiples of the distributed in integral multiples of the fundamental frequency f_0 which is the frequency of the vibration of vocal chords.

Spectrogram from
short-time fourier
transform

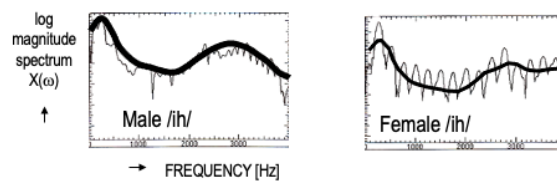


Spectrum from Fourier transform has equal spectral resolution at all frequencies



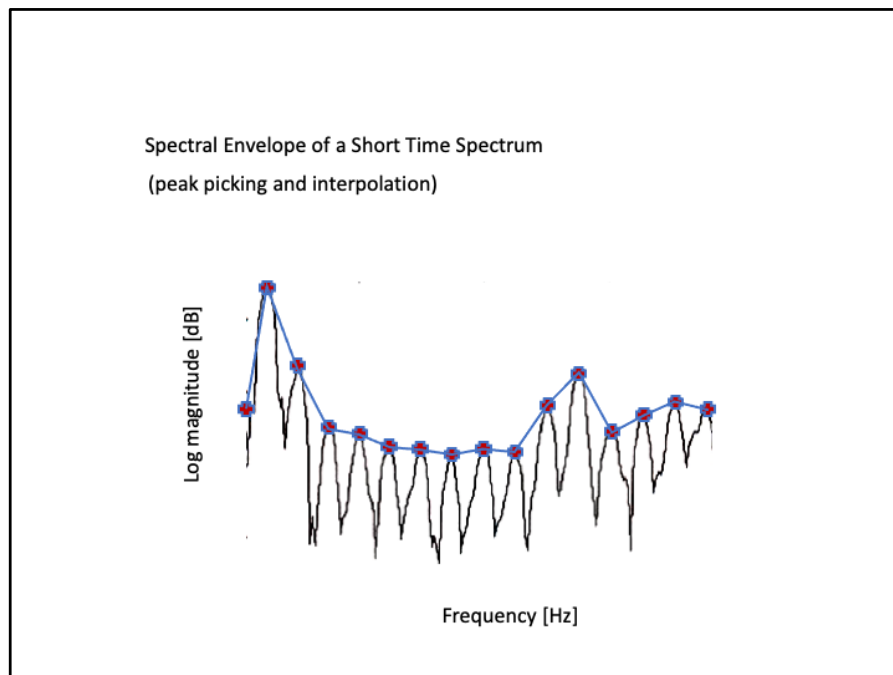
Simplest emulation of the voiced source can be characterized by a train of appropriately shaped pulses with the period of the signal waveforms (can be estimated by different means, the most straightforward one is, e.g., the signal autocorrelation). The filter is characterized by the spectral envelope of the short-time Fourier spectrum, the source is represented by the spectral fine structure.

Spectrum and its envelope



- The same vowel
 - Similar spectral envelope (vocal tract shape?)
 - Different fine structure (source of excitation)
- **Estimate spectral envelope if you want information about the phonetic quality of the sound**

Fine structure of speech spectrum differs considerably among speakers of different genders, female speech has typically higher fundamental frequency so that the fine structure components are more widely spread, male speech with its lower fundamental frequency has harmonic peaks closer to each other. However, spectral envelopes of speech segments with the same phonetic values are more similar among speakers.



One straightforward way of estimating spectral envelope is to find tips of harmonic peaks (their are spaced in integral multiples of the fundamental frequency) and to interpolate between the found peaks.

How to see fast (high frequency) and slow (low frequency) signal componets ?

- Derive spectrum of the signal (Fourier transform)



Addition

- easy to separate in a domain where the contributions of the slow changing vocal tract spectral envelope and the fast changing spectrum of source occupy different locations
- speed of changes of logarithmic spectral componets will show in Fourier transform of the log spectrum

Logarithmic spectrum of speech is seen as a superposition of the flat spectrum of the source, which changes fast in frequency, and the slower changing spectrum of the spectral envelope. Since these additive components of the logarithmic speech spectrum are changing with different rates in frequency, are more clearly separable when taking Fourier transform of the logarithmic spectrum

Cepstral Analysis



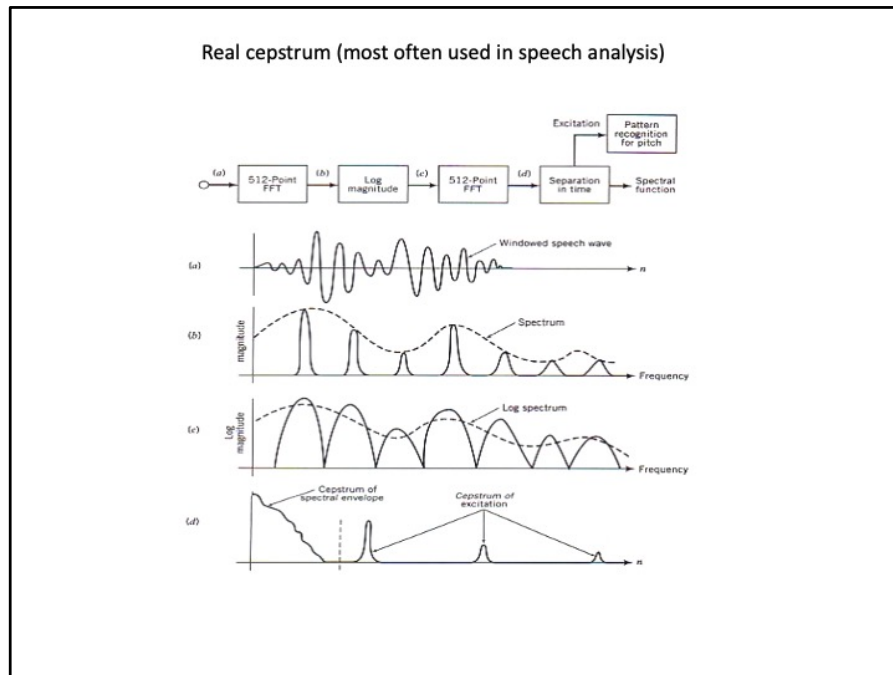
$$|E(\omega)| \cdot |V(\omega)| = |X(\omega)|$$

$$\log |E(\omega)| + \log |V(\omega)| = \log |X(\omega)|$$



Addition of slowly changing vocal tract spectral envelope and the fast changing vocal source spectrum

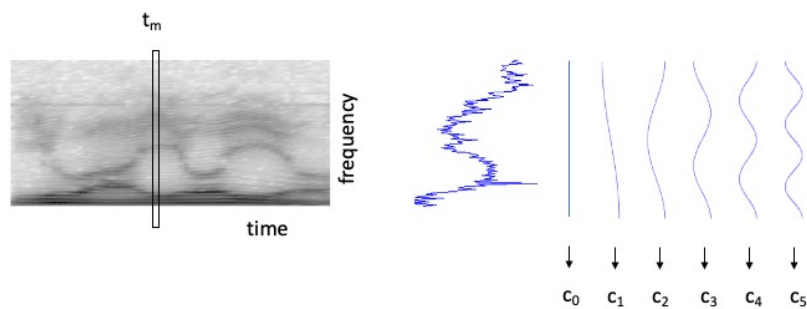
Mathematically, the log spectrum of speech is given as a sum of log spectra of the source and the spectral envelope



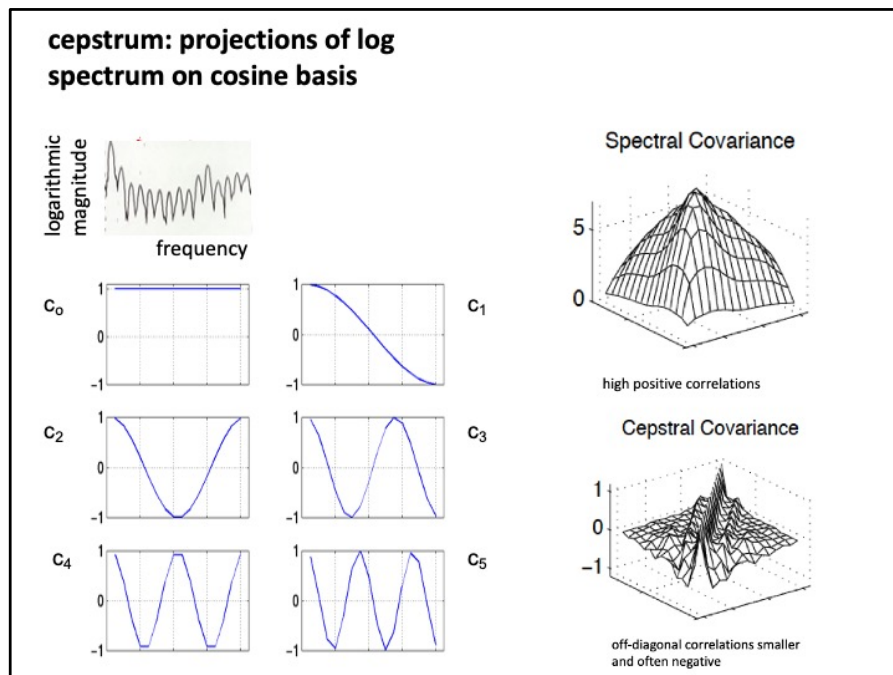
After taking fourier transform of the logarithmic speech spectrum, we obtain so called **cepstrum** where the contrution of the spectral envelope is at the low cepstral time coefficients and the contribution of the source are at higher cepstral coefficients.

Cepstrum:

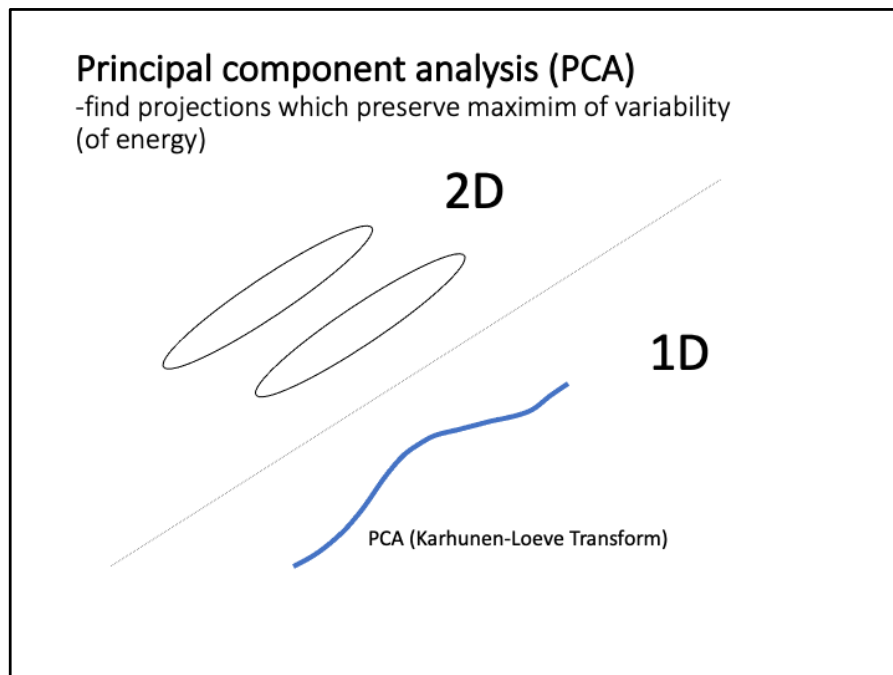
projection of log magnitude spectrum on cosine basis



Cepstrum can be seen as projection of the short-time spectral on the bases formed by integral multiples of half period of cosine functions with decreasing period (increasing frequency)

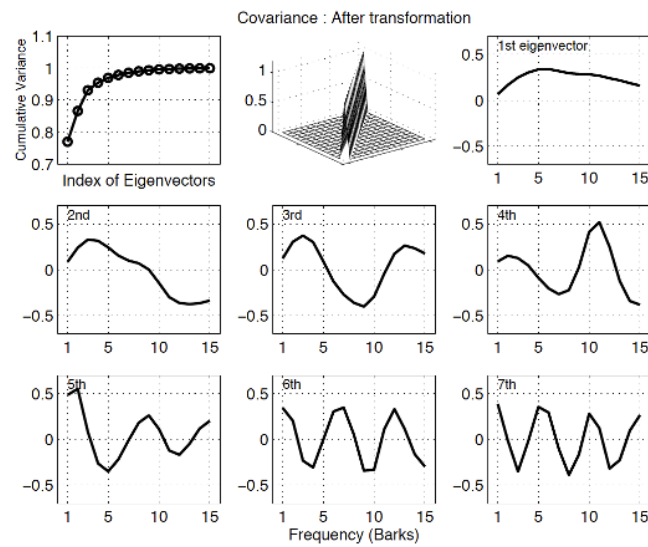


Another reason why the cepstrum can be useful (especially in pattern classifications) is that the cepstral projections approximately decorrelate the speech spectra. This allows for simpler (diagonal covariance) classifiers.



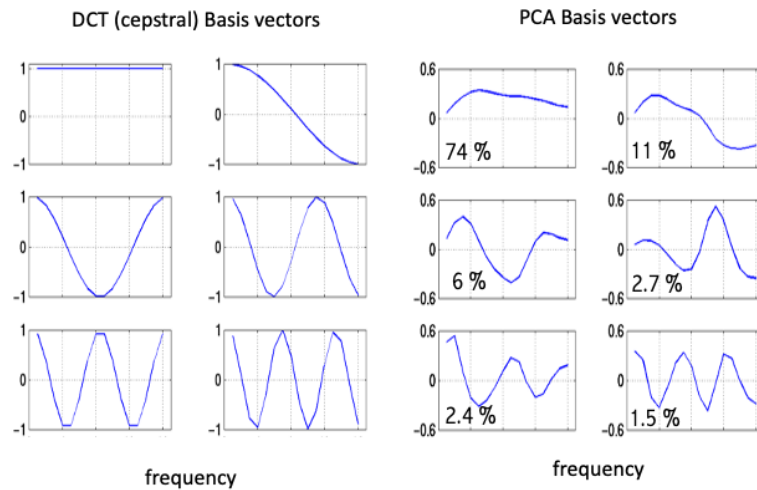
The perfect way of decorrelating vector space is to project the space in its principal components. This can be easily derived by the principal component analysis (PCA) of the space, involving the computation of the covariance matrix of the space. PCA projects the space in directions of its maximum variability. If you believe that the information you are seeking is in its variability, PCA is a good way of reducing free parameters in your representation.

PCA on (auditory-like) log magnitude spectrum

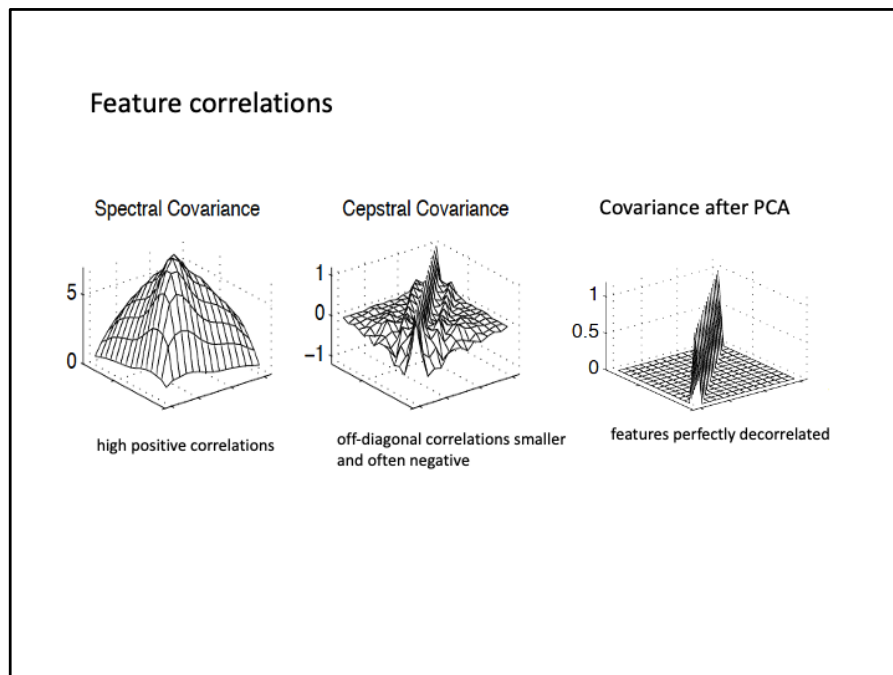


Computing PCA coefficients on (auditory-like to be discussed shortly) spectral space indeed yields projections which are similar to cosine projections in the cepstrum.

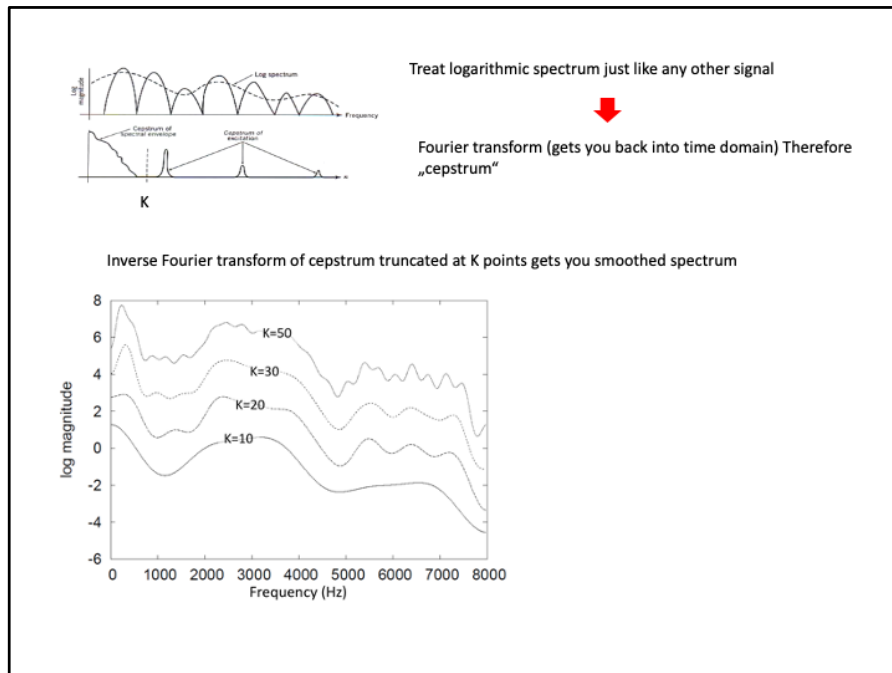
PCA Vectors from auditory-like logarithmic spectrum



Comparing the cepstral and PCA projections, we see the similarities



And comparing the covariance matrices of the cepstrally and PCA projected log spectral spaces also shows that cepstrum is a reasonable “poor man” approximation of the PCA projections (without the burden of deriving the PCA bases).



Computing cepstrum of the logarithmic spectrum gets us to the time domain. By projecting the cepstrum back to the spectral domain using the inverse fourier transform gets us back to the logarithmic spectral domain. Depending on how many cepstral coefficients we keep for the inverse fourier transform, we obtain different amounts of smoothing of the original logarithmic spectrum, m .