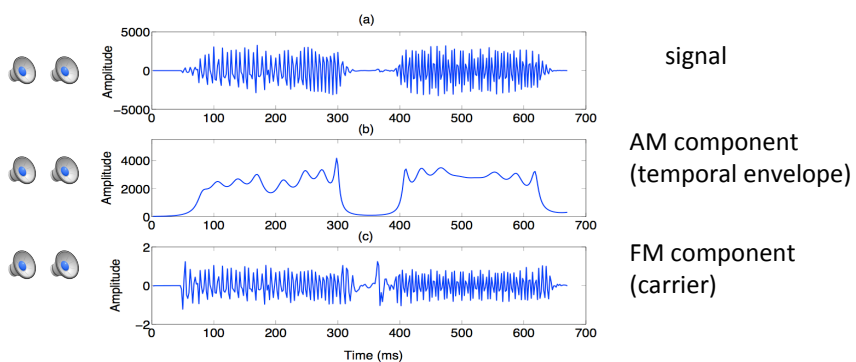


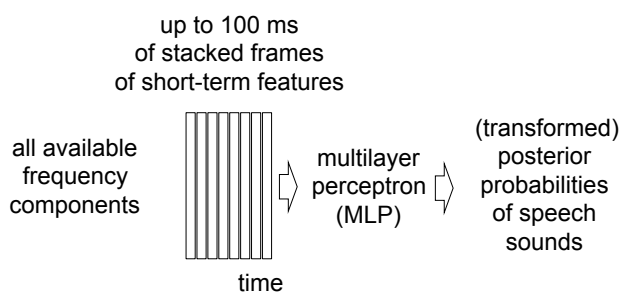
## Autoregressive model of Hilbert envelope of the signal



channel vocoder based on AM or FM components

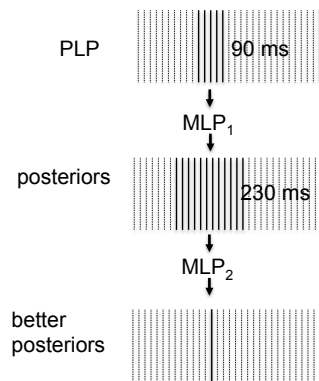
## Artificial Neural Nets for Deriving Speech Features

## conventional artificial neural net

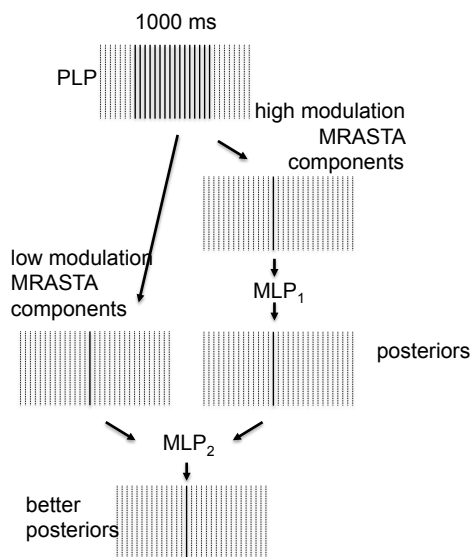


## DEEP: Some Hierarchical Nets

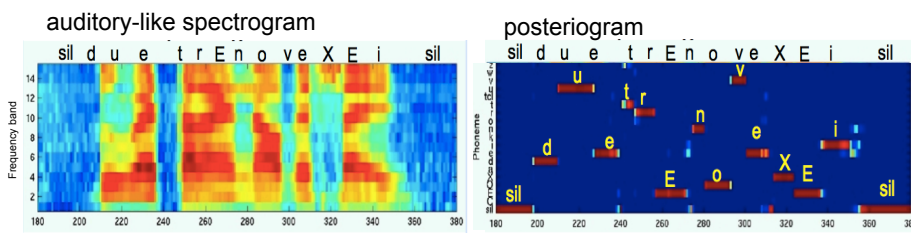
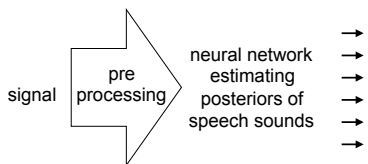
serial hierarchy (with Joel Pinto)



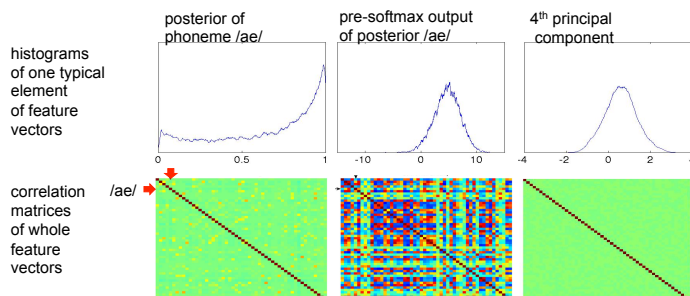
serio-parallel hierarchy (with Fabio Valente)



## Artificial Neural Nets

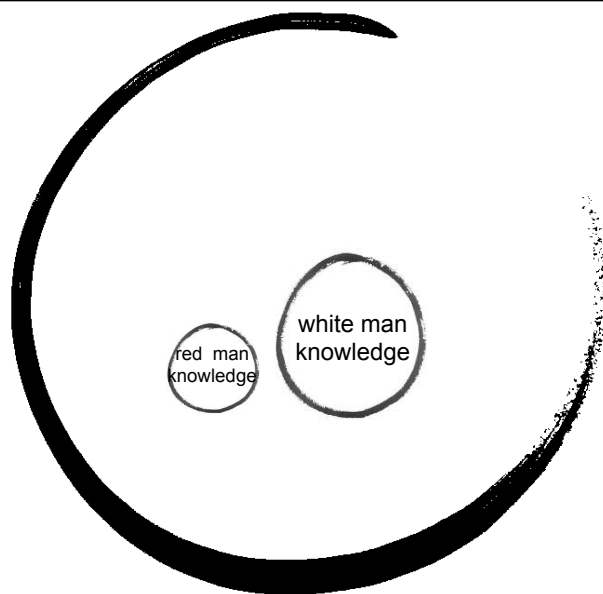


## TANDEM



# Unknown Unknowns

117



The problem is not what you do not know,  
the problem is what you do not know that you do not know

## Machine Learning

Create models of the “world”

1. from labeled (annotated) training data
2. **from prior knowledge of what is possible and how likely it is**

Find the model that best accounts for the observed data

Assumption: the future is the same as the past

- both the training and the test data are independently and identically obtained samples from the same probability distribution

Unexpected events are hard to deal with because

1. not seen in the training
2. **low (zero) prior probability**

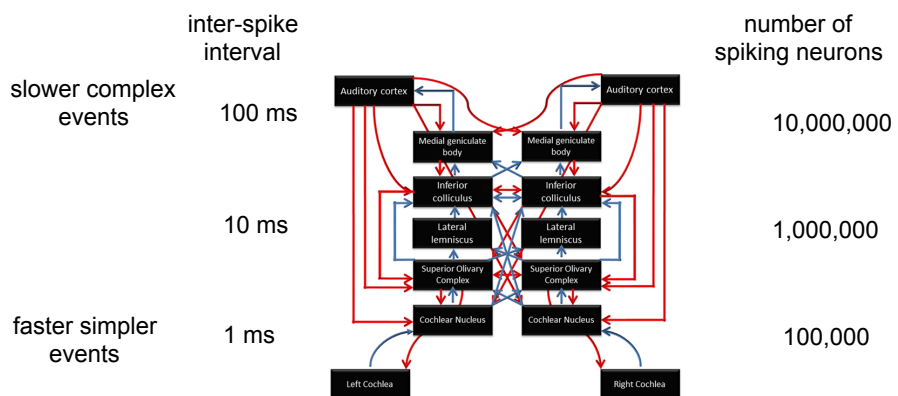
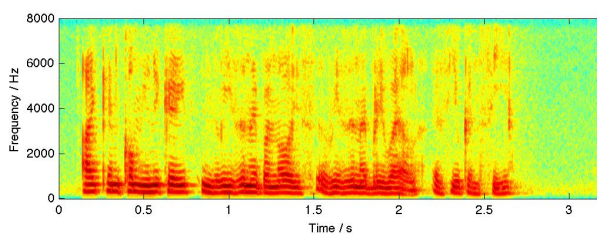
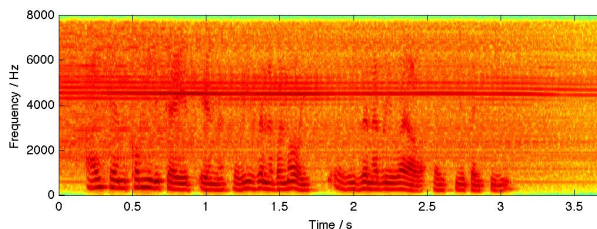
Successfully surviving natural systems attend well to the unexpected

## Power of Priors (Language Model)



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

# Unexpected Noise



Deep: many layers  
 Long: cortical event every 100 ms or so  
 Wide: many possible descriptions of an event in auditory cortex

**DEEP:** Information in the signal should be extracted in stages, from description of signal features to description of phonetic events.

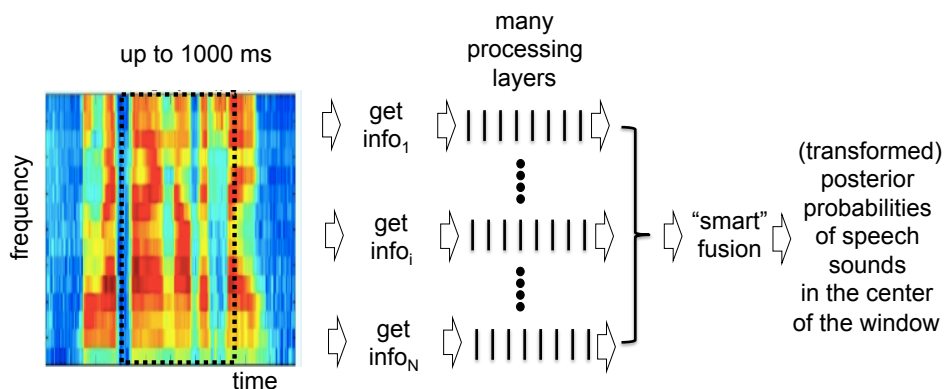
**LONG:** Information about underlying speech sounds is spread in time for more than 200 ms

**WIDE:** There are many ways to form parallel processing streams using different signal projections and different prior assumptions.

Not all processing streams get always corrupted and we need to find ways to find the uncorrupted processing streams.

Information in speech is coded hierarchically (**deep**)  
in temporal dynamics (**long**)  
and in many redundant dimensions (**wide**)

### Deep, Long, and Wide Neural Nets



# Longer is Better

Phonetic classifier accuracy as a function of a time span of an analysis

Fant, Cole, Roginski NIPS 1992

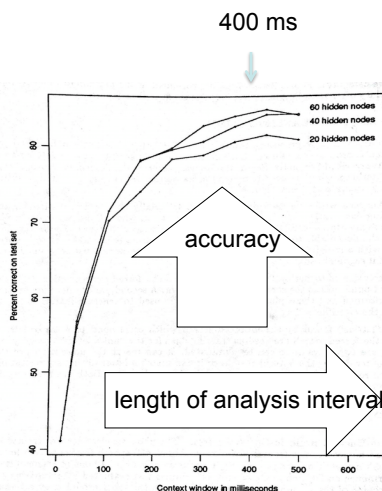
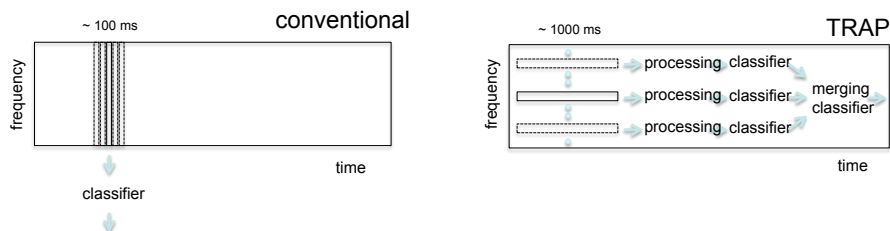
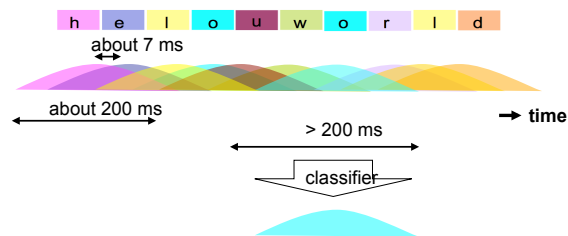


Figure 1: Performance of the phonetic classifier as a function of PLP context and number of hidden units.

## LONG: Classifying Temporal Patterns of Spectral Energies with Sangita Sharma, Pratibha Jain, Honza Cernocky, Pavel Matejka, Petr Schwartz ....

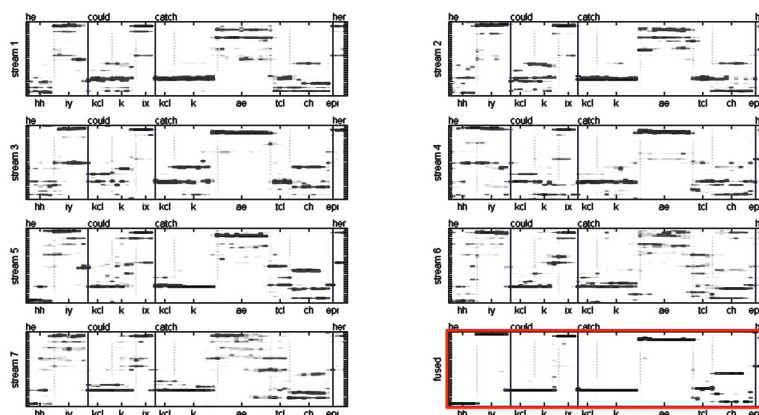


Each temporal pattern contains most of coarticulation span of speech sound in its center.



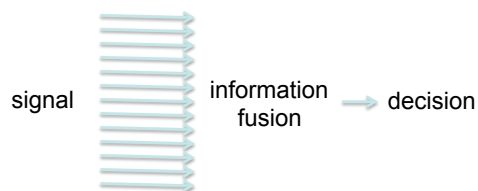


## Fusion of streams of different carrier frequencies



## Wide: Multi-stream Processing

Information in speech is coded in many redundant dimensions.  
Not all dimensions get corrupted at the same time.



- Parallel information-providing streams, each carrying different redundant dimensions of a given target.
- A strategy for comparing the streams.
- A strategy for selecting "reliable" streams.

### Stream formation

- Different perceptual modalities
- Different processing channels within each modality
- Bottom-up and top-down dominated channels

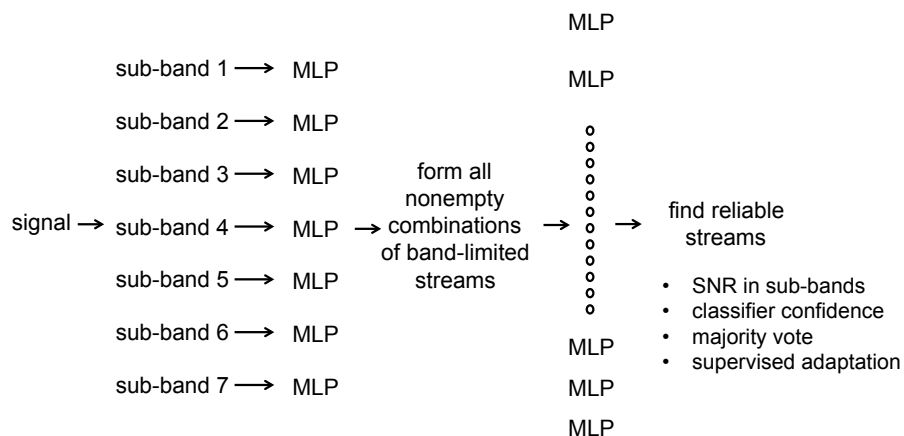
### Comparing the streams ?

- various correlation (distance) measures

Selecting reliable streams ?????

### Early Attempts for Multi-Stream Recognition

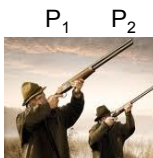
with Sangita Sharma and Misha Pavel



### Monitoring Performance

Fletcher et al  
Boothroyd ann Nittrouer  
Allen

$$P(\epsilon) = \prod_i P(\epsilon_i)$$



P<sub>1</sub> P<sub>2</sub>



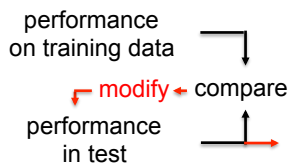
$$P_{\text{miss}} = (1-P_1)(1-P_2)$$

**observer** - false positives and negatives are possible

$$P_{\text{miss\_observed}} \neq (1-P_1)(1-P_2)$$

Do listeners know when they know ?

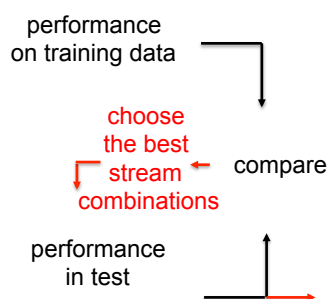
How to make machine know when it knows ?



## Finding Reliable Streams

Streams which yield the best performance on the test data

Classifier can never work better than it does on the data on which it was trained



## Evaluating Performance

How often sound classes occur and how often do they get confused?

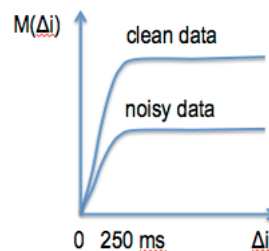
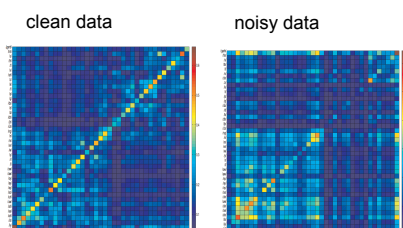
$$AC = \frac{1}{N} \sum_{i=1}^N (\mathbf{p}_i)^r (\mathbf{p}_i^T)^r$$

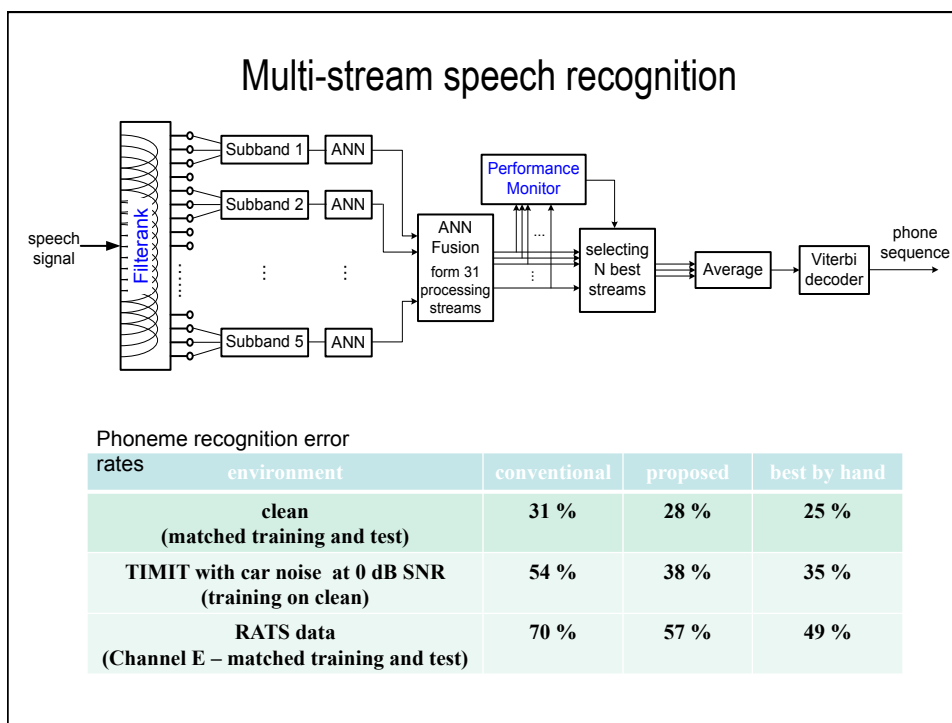
$\mathbf{p}_i$  – vector of sound posteriors at  $i$ -th time instant  
 $N$  – time interval of the evaluation  
 $r$  – th power element-by-element (currently  $r=0.1$ )

How much sound classes differ and how fast do they change?

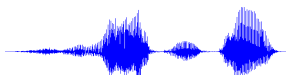
$$M(\Delta i) = \frac{\sum_{i=0}^{N-\Delta i} D(\mathbf{p}_i, \mathbf{p}_{i+\Delta i})}{N - \Delta i}$$

$\Delta i$  – time delay  
 $D(\cdot)$  – symmetric KL divergence





## Towards Increasing Error Rates



Signal processing, information theory, machine learning, ...

↓

signal processing

↓

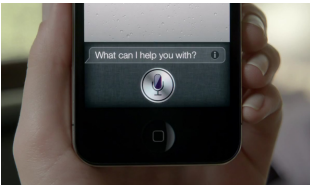
pattern classification

↓

decoder

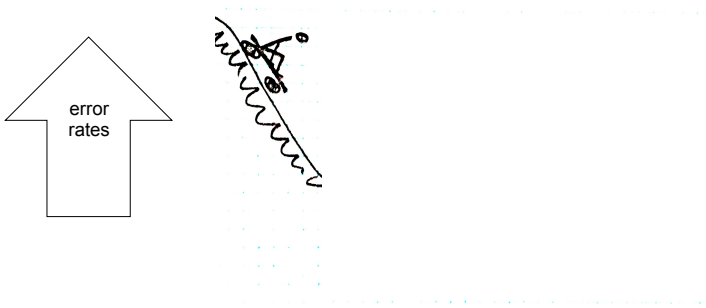
↓

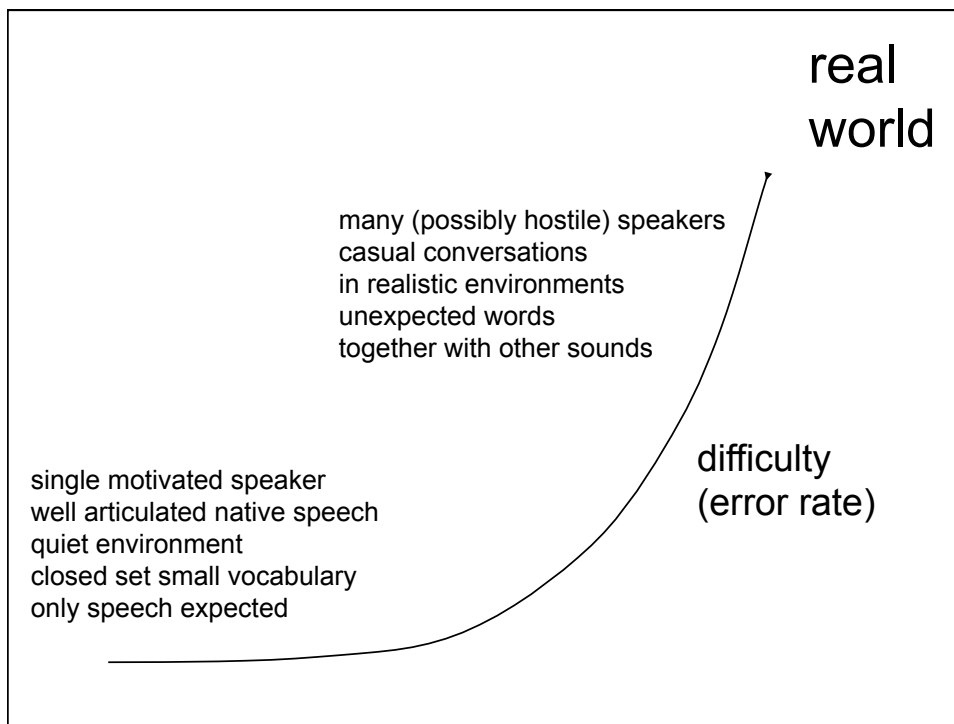
message



Why to rock the boat? We have good thing going.


Why to rock the boat?  
We have good thing going.





Repetition, fillers, hesitations, interruptions, unfinished and non-grammatical sentences, new words, dialects, emotions, ...

Current DARPA and IARPA programs, research agenda of the JHU CoE HLT, industrial efforts (Google, Microsoft, IBM, Amazon,...)



Signal processing, information theory, machine learning, ...

&

neural information processing, psychophysics, physiology, cognitive science, phonetics and linguistics, ...

**Engineering and Life Sciences together !**

## How to Get There ?

Fred Jelinek



Speech recognition  
...a problem of maximum likelihood decoding

**information and communication theory, machine learning, large data,....**

Roman Jakobson



We speak, in order to be heard, in order to be understood

**human communication, speech production, perception, neuroscience, cognitive science,...**

Gordon Moore



The complexity for minimum component costs has increased at a rate of roughly a factor of two per year...

**tools**

John Pierce



**..devise a clear, simple, definitive experiments. So a science of speech can grow, certain step by certain step.**

However, also John Pierce:

*(Speech recognition is so far (1969) field of) mad inventors or untrustworthy engineers (because machine needs) intelligence and knowledge of language comparable to those of a native speaker .*

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

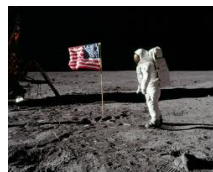
## Why Am I Working in Machine Recognition of Speech?



Why did I climbed Mt. Everest?  
Because it is there !  
-Sir Edmund Hilary

Spoken language is one of the  
most amazing accomplishments of  
human race.

Implement .... *intelligence and knowledge of  
language comparable to those of a native  
speaker !*



Don't Follow Leaders,  
Watch the Parking Meters...