

**SPEECH IS THE PROBLEM, HEARING IS THE SOLUTION**  
**CHANNEL CAPACITY AND THE “FRONT END”**

**Jont B. Allen**  
**ECE Univ. of IL**  
**Beckman Inst., Urbana IL**

October 7, 2003

[jba@auditorymodels.org](mailto:jba@auditorymodels.org)

## MY VIEW

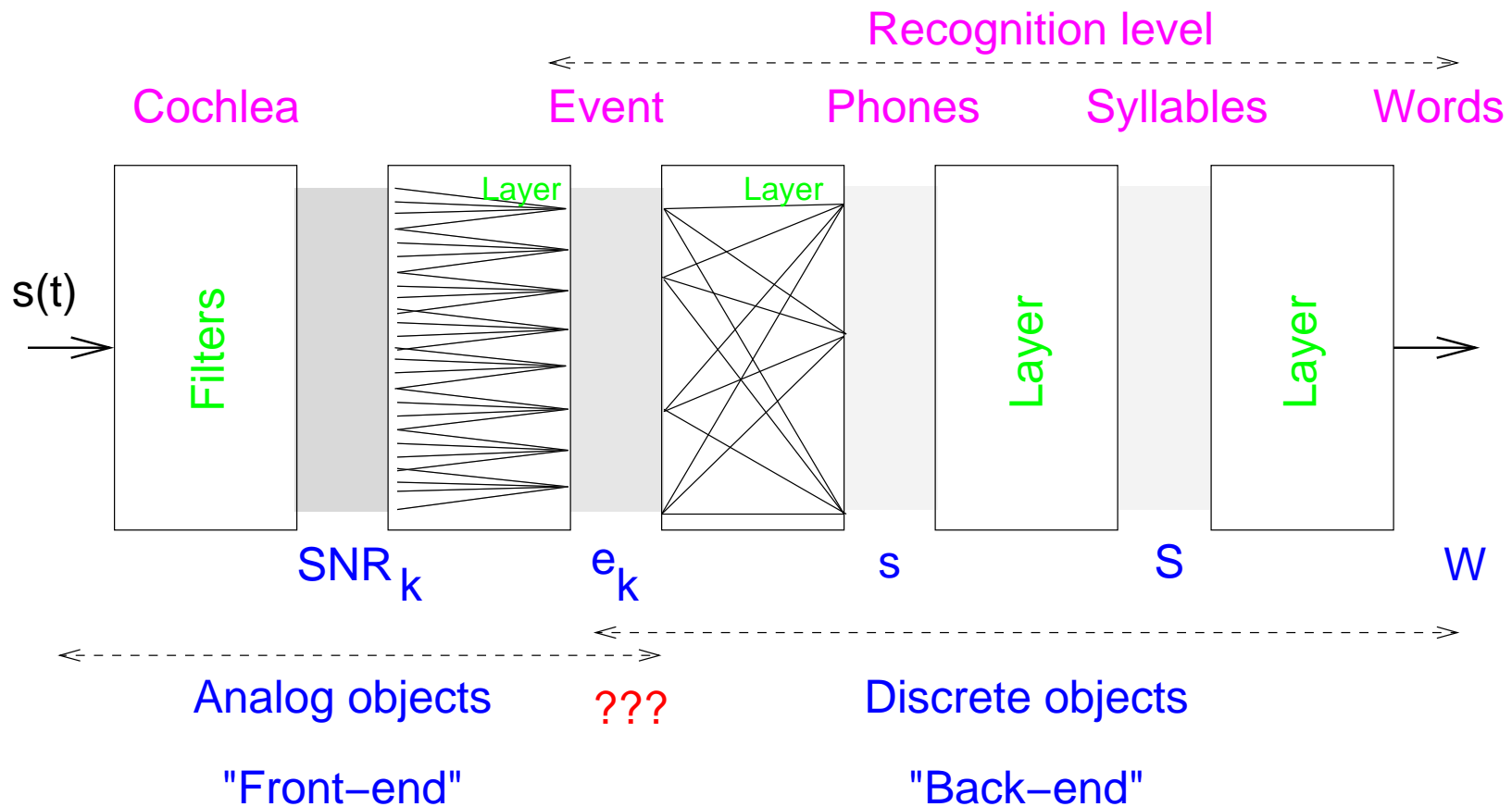
- “We have met the enemy, and he is us.” –*Pogo (Walt Kelly)*
  - Why aren’t we benchmarking against human performance?
- Biological systems are the ultimate information processors  
And, we need to learn from them

## WHAT I WANT TO SHOW:

- Humans have an intrinsic **robustness** to noise and filtering
  - **Robustness** is not due to **semantic context effects**
- **HSR** is a bottom–up, divide and conquer strategy
  - We recognize speech based on a hierarchy of **context layers**
  - As in **vision**, **entropy decreases** as we **integrate context**

## HOW WE RECOGNIZE SPEECH?

- Hierarchical “bottom up” analysis
- Accurate statistical models of performance at each stage



- Entropy drops (i.e., context is integrated) in stages

## MODEL OF BAND EVENT ERRORS

- When the SNR is varied they found that the event-error is

$$e_k = e_{min}^{SNR_k/K}$$

where  $SNR_k$  is the signal to noise ratio in dB, divided by 30, such that

$$0 \leq SNR_k \leq 1$$

- 

$$SNR_k \equiv \left\{ \begin{array}{ll} 0 & 20 \log_{10}(snr_k) < 0 \\ 20 \log_{10}(snr_k)/30 & 0 < 20 \log_{10}(snr_k) < 30 \\ 1 & 30 < 20 \log_{10}(snr_k). \end{array} \right\}$$

- Total error:

$$e = e_1 e_2 \cdots e_K = e_{min}^{(SNR_1 + SNR_2 \cdots SNR_K)/K}$$

- The speech SNR (not the energy) determines the event errors  $e_k$  and thus the phoneme articulation  $s = 1 - e_1 e_2 \cdots e_K$

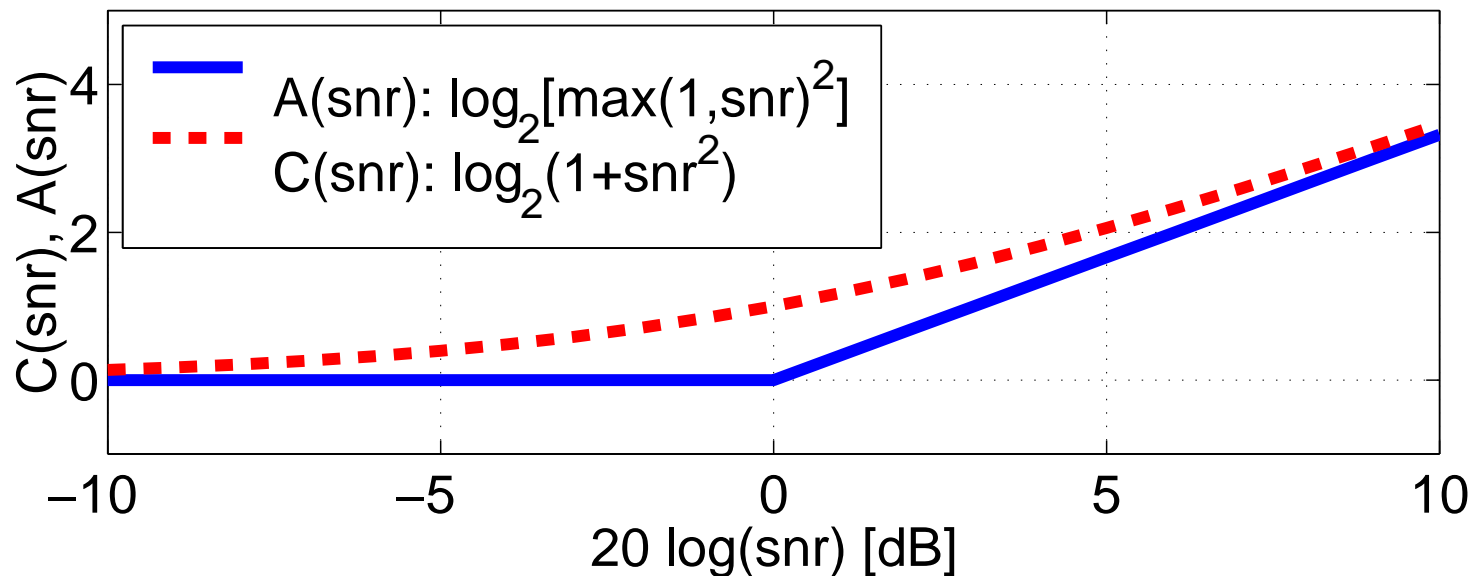
## AI AS A CHANNEL CAPACITY

- Since  $\sum_k (\log \text{snr}_k) = \log(\prod_k \text{snr}_k)$

$$\mathcal{A} \equiv \frac{1}{K} \sum_k \text{SNR}_k \propto \log \left( \prod_k \text{snr}_k \right)^{1/K} \quad (1)$$

- and from Shannon

$$C = \int_{-\infty}^{\infty} \log_2[1 + \text{snr}^2(f)] df, \quad (2)$$

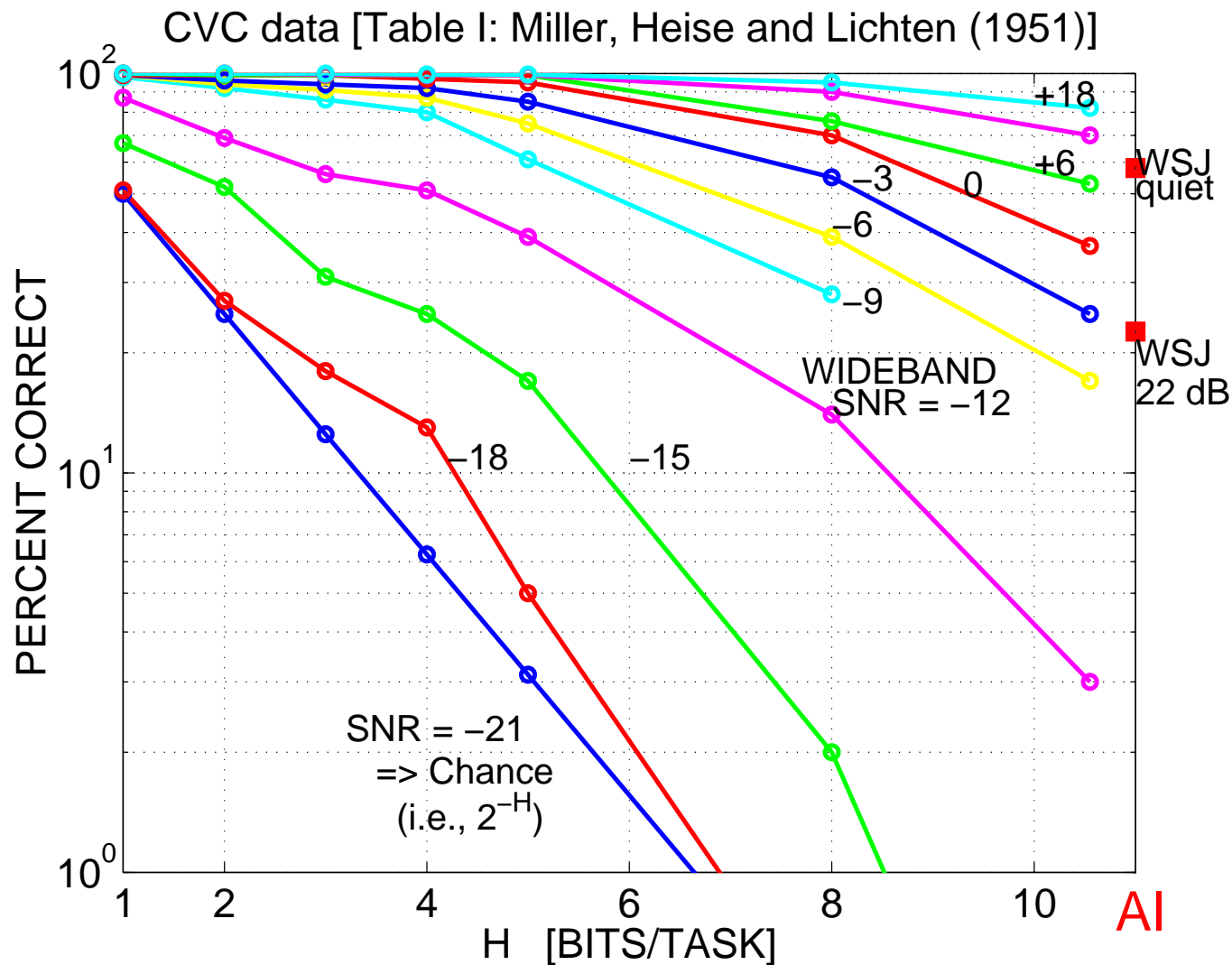


## AI IS A CHANNEL CAPACITY

- In conclusion:
  - The channel capacity is the maximum information rate that can asymptotically be sent over a channel without error
  - The AI is basically a channel capacity

# SPEECH ENTROPY VS. THE WIDEBAND SNR

- $P_c(\mathcal{H}, SNR)$  Miller, Heise and Lichten 1951
- Many of the results of MHL51 expand on the AI model



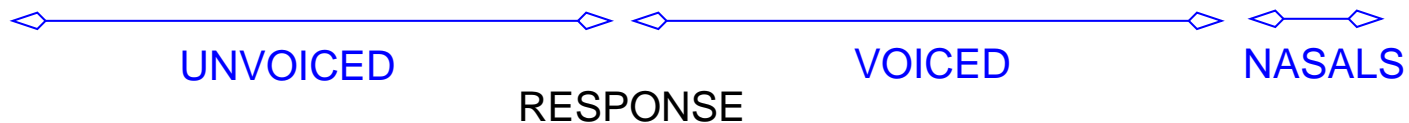


# CONFUSION MATRIX PARTITIONING

- Miller & Nicely 1955 Confusion Matrix (Table III)
  - MN55 established a natural phone hierarchical clustering:

TABLE III. Confusion matrix for  $S/N = -6$  db and frequency response of 200–6500 cps.

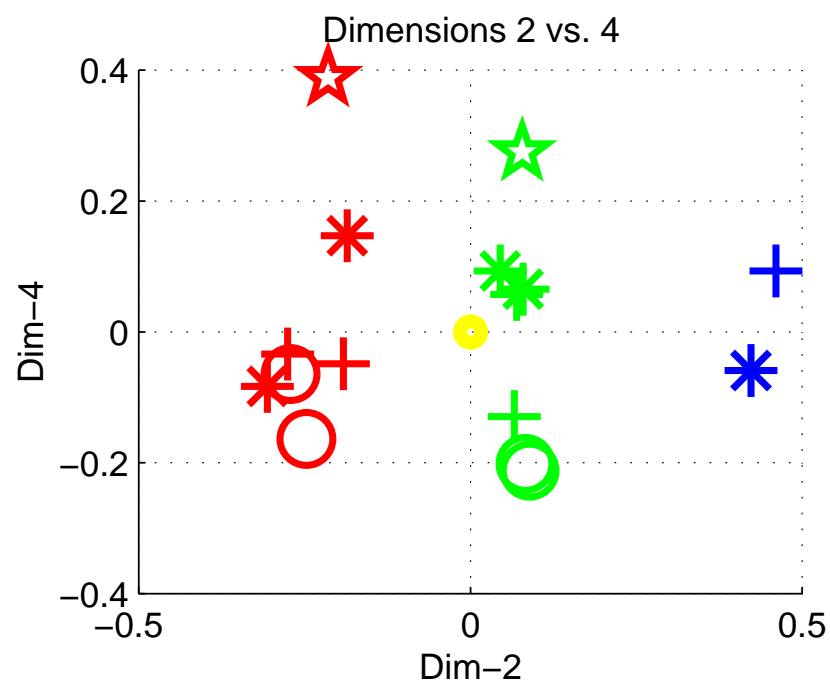
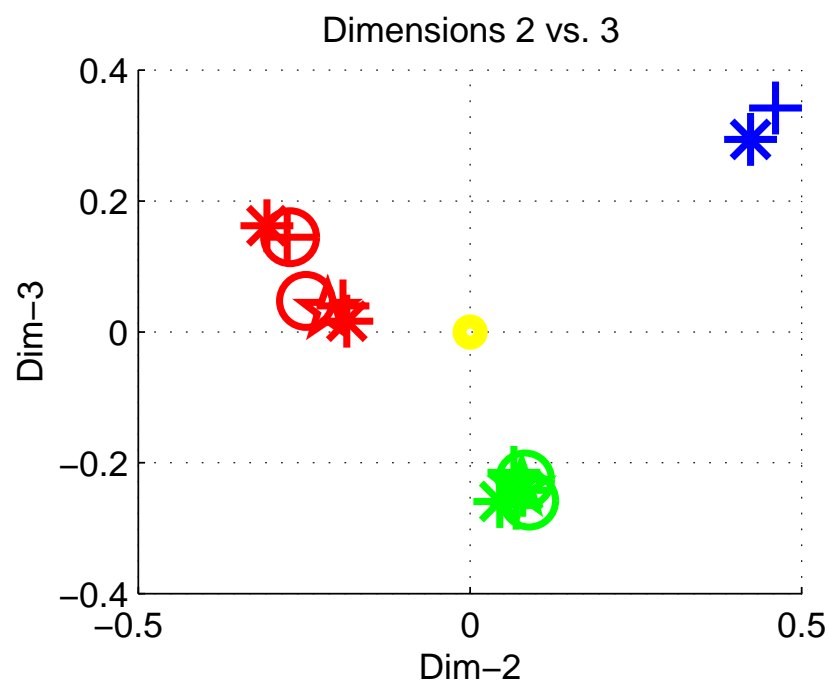
	<i>p</i>	<i>t</i>	<i>k</i>	<i>f</i>	<i>θ</i>	<i>s</i>	<i>ʃ</i>	<i>b</i>	<i>d</i>	<i>g</i>	<i>v</i>	<i>ø</i>	<i>z</i>	<i>ʒ</i>	<i>m</i>	<i>n</i>
<i>p</i>	80	43	64	17	14	6	2	1	1		1	1			2	
<i>t</i>	71	84	55	5	9	3	8	1				1	2		2	3
<i>k</i>	66	76	107	12	8	9	4					1			1	
<i>f</i>	18	12	9	175	48	11	1	7	2	1	2	2				
<i>θ</i>	19	17	16	104	64	32	7	5	4	5	6	4	5			
<i>s</i>	8	5	4	23	39	107	45	4	2	3	1	1	3	2		1
<i>ʃ</i>	1	6	3	4	6	29	195		3							1
<i>b</i>	1			5	4	4		136	10	9	47	16	6	1	5	4
<i>d</i>							8	5	80	45	11	20	20	26	1	
<i>g</i>					2			3	63	66	3	19	37	56		3
<i>v</i>				2		2		48	5	5	145	45	12		4	
<i>ø</i>					6			31	6	17	86	58	21	5	6	4
<i>z</i>					1	1	1	7	20	27	16	28	94	44		1
<i>ʒ</i>								1	26	18	3	8	45	129		2
<i>m</i>	1							4			4	1	3		177	46
<i>n</i>					4			1	5	2		7	1	6	47	163



“This breakdown of the confusion matrix into five smaller matrices . . . is equivalent to . . . five communication channels . . . .” –Miller & Nicely 1955

# SVD REPRESENTATION OF THE PERCEPTUAL SPACE

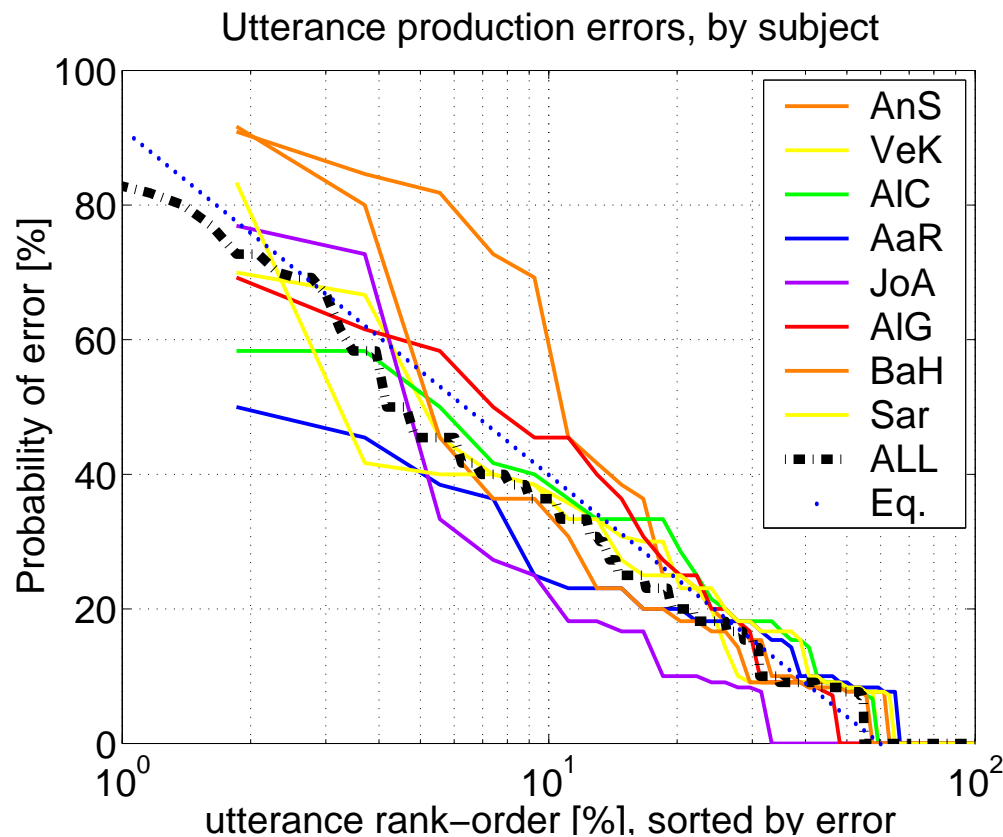
- $4^{dim}$  SVD perceptual representation of the confusion matrix



DEMO

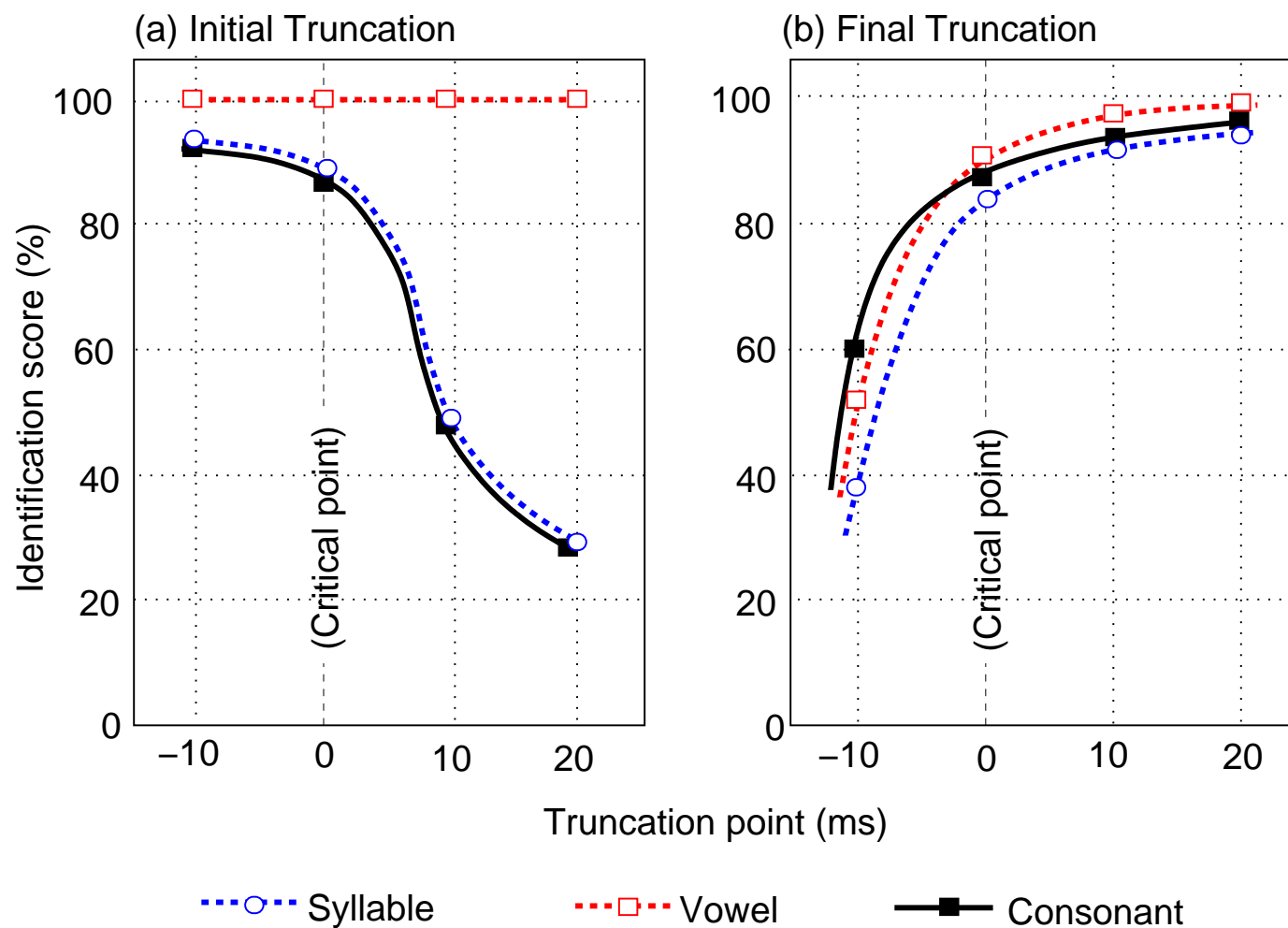
## DISTRIBUTION OF UTTERANCE ERRORS

- What determines  $s_{max} = 1 - e_{min}$ ?
- Utterance *talker mispronunciations*, as defined by 32 listeners
- Errors are distributed like **Zipf's Law** [ $\dots N/N_T \approx 0.6e^{-4.48P_e}$ ]  
 35% of the utterances have **no** error  
 33% have  $> 10\%$  error,  $10\% > 35\%$  error,  $5\% > 50\%$  error



# TEMPORAL RESOLUTION OF PHONE RECOGNITION

- Phones are recognized in on a 10 ms time scale (Furui 1986)



## GRAMMATICAL CONTEXT

- Five groups of five words that form grammatical sentences:

Don	Brought	His	Black	Bread
He	Has	More	Cheap	Sheep
Red	Left	No	Good	Shoes
Slim	Loves	Some	Wet	Socks
Who	Took	The	Wrong	Things

- Tests:

5 word lists

25 word

25 words with grammatical context

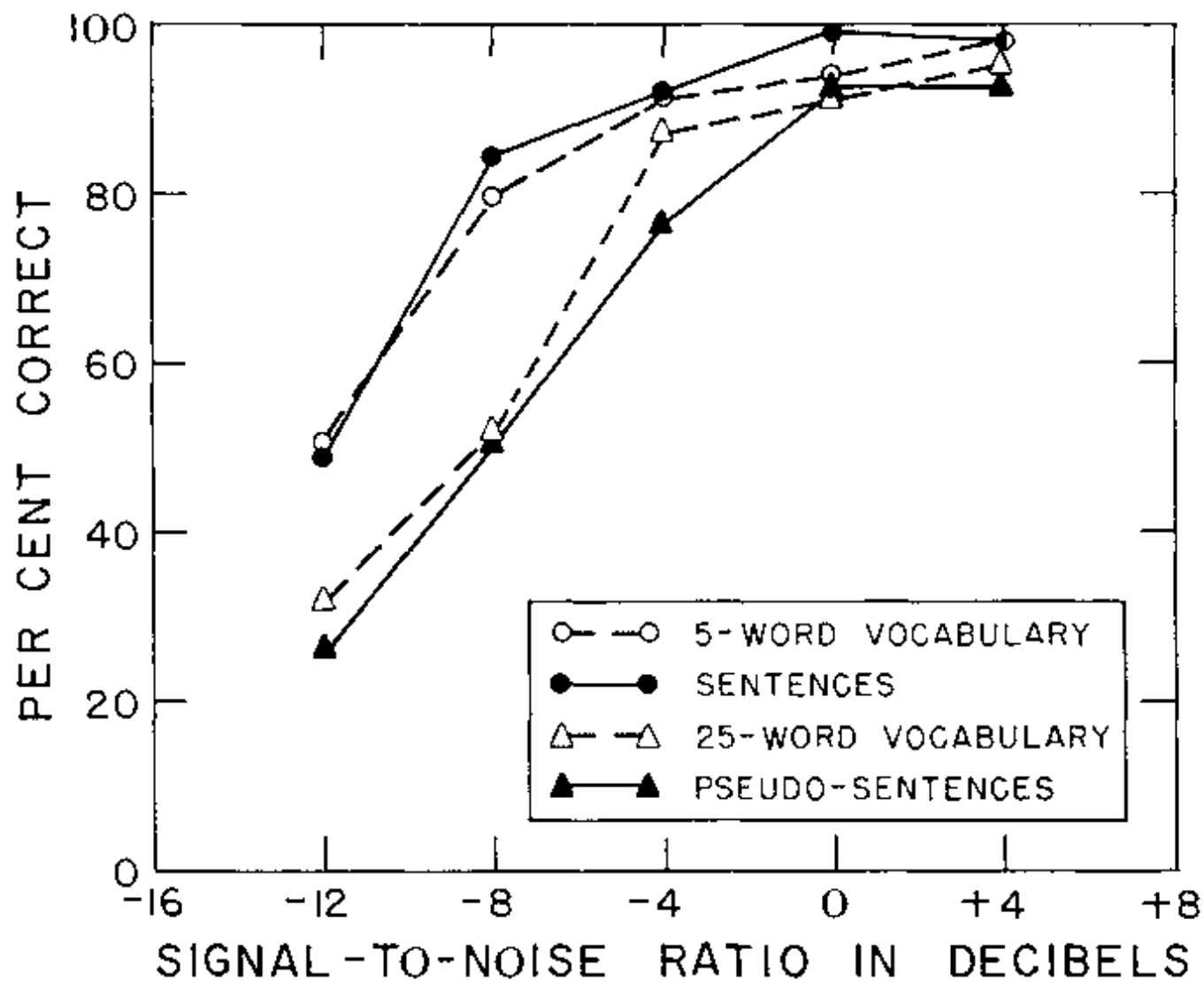
Example: **He left no black socks**

25 words reverse order

Example: **Socks black no left he.**

## GRAMMATICAL CONTEXT

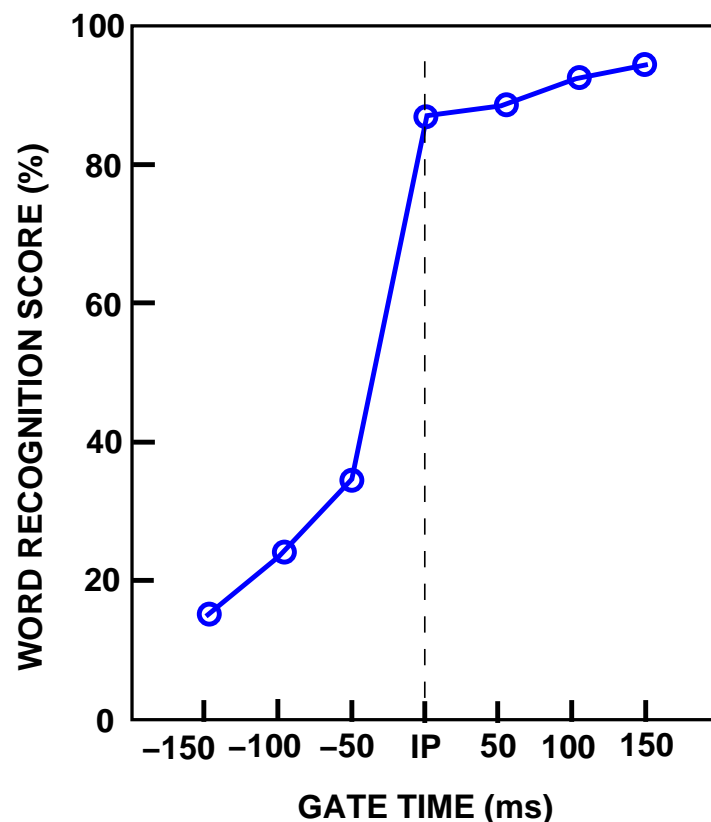
- Results of tests



## WORD SEMANTICS: IP DEFINITION

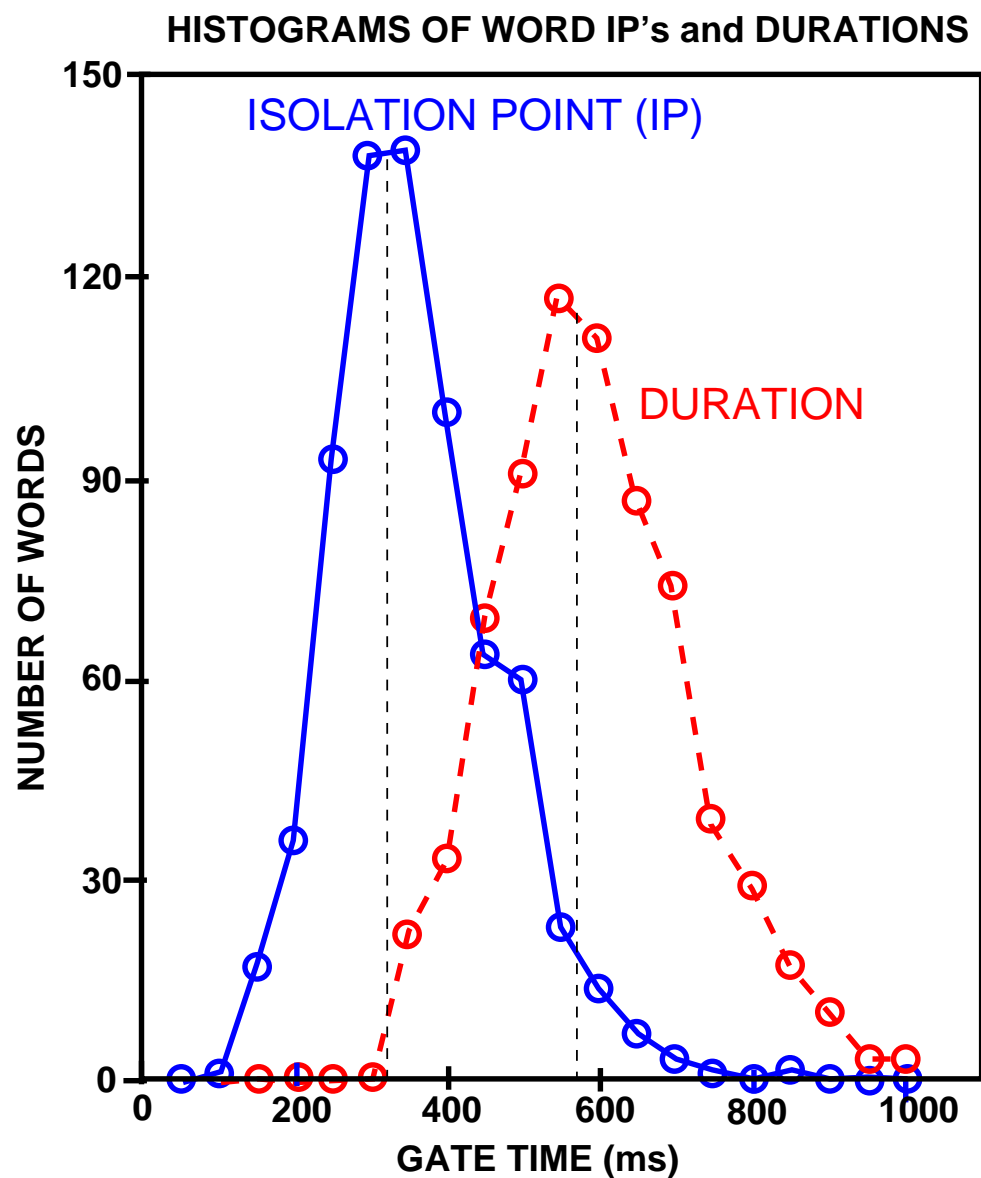
- 704 isolated words were truncated in 50 ms steps [Van Petten 1999](#)
- **Isolation point** is defined as *the time of the discontinuity in recognition*  
Expt. I – **Neutral sentences**: “The next word is *test-word*.”

ACCURACY OF IDENTIFICATION VERSUS GATE TIME



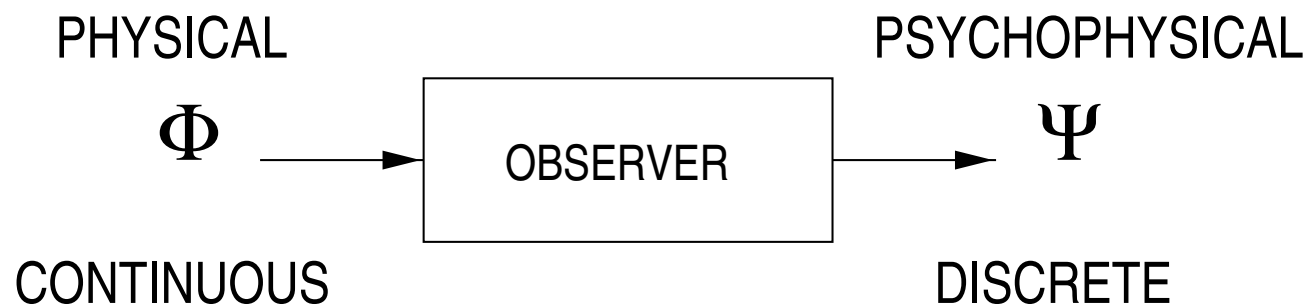
## WORD SEMANTICS: IP VS. DURATION

- Isolation point vs. word durations (real words, no sentence context)





## FROM CONTINUOUS TO DISCRETE



- $\Phi$ -domain signals

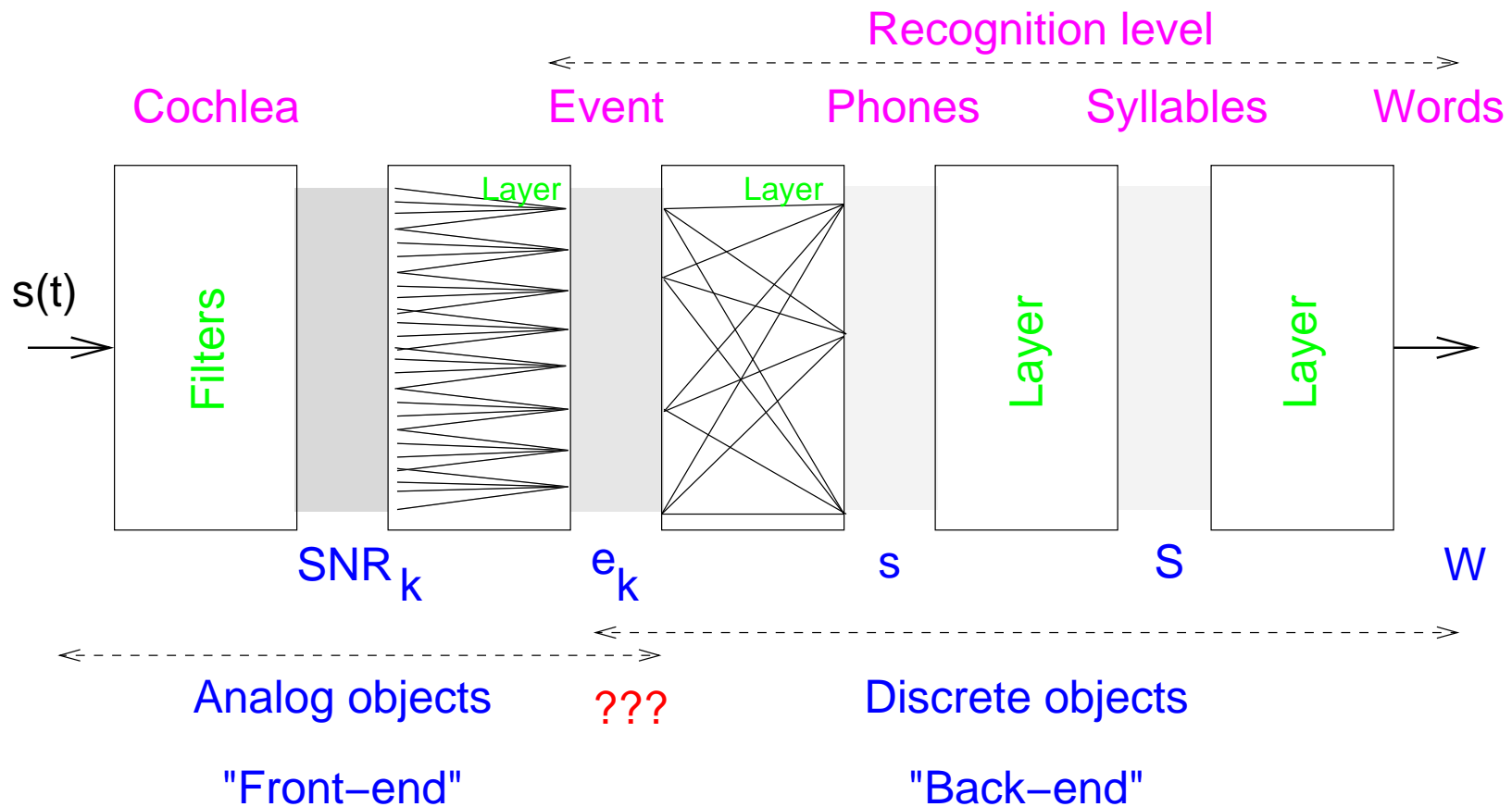
- Speech signal
- Cochlear filter outputs
- Neural rate
- Voltage in cochlear nucleus cells

- $\Psi$ -domain objects

- Words
- Syllables
- Phonemes
- Events [Miller's features]

# HOW WE RECOGNIZE SPEECH?

- Hierarchical “bottom up” analysis
- Accurate statistical models of performance at each stage



- Entropy drops (i.e., context is integrated) in stages