

FROM LORD RAYLEIGH TO SHANNON: HOW DO WE DECODE SPEECH?

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

jba@auditorymodels.org

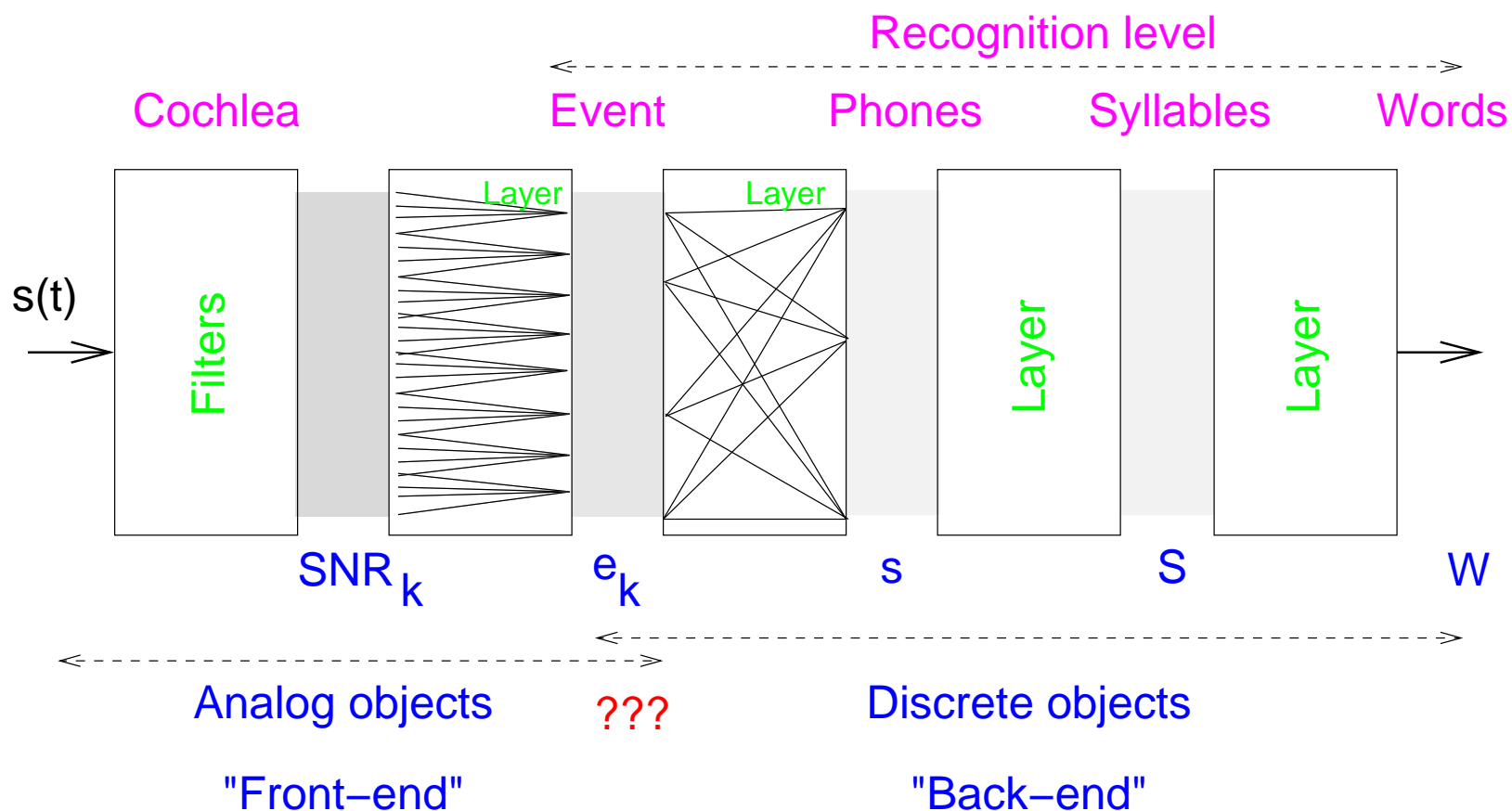
<http://auditorymodels.org/CUNY/>
<http://auditorymodels.org/jba/PAPERS/ICASSP/>

WHAT I WANT TO SHOW:

- Biological systems are the ultimate information processors
- **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**
- Humans have an intrinsic **robustness to noise and filtering**
 - **Robustness** is not due to **semantic context effects**

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

DEFINITIONS

✓ phone	A consonant (C) or vowel (V) sound
word	A meaningful phone or phone sequence (i.e., cat \equiv CVC)
phoneme	The replaceable set of phones which leave a word meaning invariant
recognition	Probability measure P_c of correct phoneme identification
✓ articulation	Recognition of “nonsense words”
✓ intelligibility	Recognition of words (i.e., meaningful speech)
robustness	Relative recognition with filtering and noise
confusion matrix	Table of identification frequencies $N_{sr} \equiv N_{r s}$
articulation matrix	A <i>confusion matrix</i> composed of nonsense sounds
✓ articulation event	A discrete subunit of articulation [e.g., Voicing: /ba/ vs. /pa/]
trial	A single presentation of a set of events
state	A values of a set of events at some instant of time
state machine	A machine (program) that transforms from one state to another
noiseless state machine	A deterministic state machine
context	Coordinated combinations of events within a trial
message	Specific information transmitted by a trial
p_n information density	Probability of event n , of N possible events $I_n = \log_2(1/p_n), \quad n = 1, \dots, N$
entropy	Average information: $H = \sum_{n=1}^N p_n I_n$

KEY HSR STUDIES

- The first articulation experiments date from Lord Rayleigh's 1908 and George Campbell 1910 phoneme identification experiments
- A basic probabilistic approach was developed by Stewart & Fletcher 1921
 - Detailed review of Fletcher's AI theory: Allen IEEE 1994
- French and Steinberg 1947 WWII studies
- Shannon's Information theory 1948+
- G.A. Miller, Heise and Lichten 1951; G.A. Miller & Nicely 1955
- *Language and communication* G.A. Miller, 1951 McGraw Hill
 - Miller first introduces IT to language modeling, following Shannon
- Boothroyd JASA 1968; Boothroyd & Nittrouer JASA 1988
- Bronkhorst et al. JASA 1993, 2002
- Van Petten et al. 1994
- Detailed review chapter Allen 2003

MOTIVATION

- Results of Lippmann 1997, sorted by Error Ratio

Corpus	Size in Words	Conditions	% Error		Error Ratio
			Machine	Human	
Alphabetic	26	20-talkers 8-listeners	5.0 ^{isolated}	1.6 ^{continuous}	3
Resource	1000	null grammar	17	2	8
WSJ-NAB	5000	quiet (trained)	7.2	0.9	8
Switchboard	14,000	spontaneous (tel. BW)	43	4	11
WSJ-NAB	5000	10 dB (trained)	12.8	1.1	12
WSJ-NAB	65,000	close mic	6.6	0.4	16
WSJ-NAB	65,000	omni mic	23.9	0.8	30
Resource	1000	word-pair grammar	3.6	0.1	36
WSJ-NAB	5000	quiet (not trained)	42	0.9	47
WSJ-NAB	5000	22 dB (not trained)	77.4	0.9	86
word spotting	20	judgment errors	24	0.3	80
TI-digit	10	connected	0.72	0.009	80

DEMO ScanMail examples

/Audio/ScanMailExample

TYPICAL ARTICULATION TEST RECORD

- Fletcher's method of nonsense phone error analysis

ARTICULATION TEST RECORD

DATE March 1928
3-16-28 SYLLABLE ARTICULATION 51.5%

TITLE OF TEST PRACTICE TESTS CONDITION TESTED 1500~ LOW PASS FILTER

1500 Hz lowpass filtering

NO.		OBSERVED	CALLED	OBSERVED	CALLED	OBSERVED	CALLED
1	THE FIRST GROUP IS	má'v	ná'v	pó'z	po'th	Kób ✓	Kób
2	CAN YOU HEAR	pōch ✓	pōch	nē'z	nē'zh	shē'th	siz
3	I WILL NOW SAY	seng ✓	seng	jō'ch ✓	jō'ch	fū'ch ✓	fū'ch
4	AS THE FOURTH WRITE	chūd ✓	chū'd	thā'm ✓	thā'm	thā'l ✓	thā'l
5	WRITE DOWN	run ✓	run	hab ✓	hab	po'th ✓	po'th

DATA

$$S \equiv P_c(\text{syllable}) = 0.515$$

$$v \equiv P_c(\text{vowels}) = 0.909$$

$$c \equiv P_c(\text{consonants}) = 0.74$$

MODELS

$$\hat{S} = cvc = 0.498 \quad (\text{CVC syllable model})$$

$$s \equiv P_c(\text{phone}) = (v + 2c)/3 = 0.796$$

$$s^3 = 0.505 \quad (3 \text{ phone syllable model})$$

THE METHOD

- The data bases they used were formed from
 - statistically balanced
 - nonsense
 - CVC, CV and VC syllable lists
where C represents a consonant and V a vowel
- The syllable lists were spoken, and the listeners recorded what they heard
- Probabilities-correct c and v for the sound-units were computed
- The average $\{C,V\}$ speech-unit articulation probability s was computed from the composition of $\{C,V\}$ units in the data base
(i.e. $s = (2c + v)/3$ for CVC's, $s = (c + v)/2$ for CV's)
 - Measure s looks like a sufficient statistic

WHAT THEY FOUND

- **Nonsense phones** are recognized as **independent units**:
 - The probability of correct recognition for the average phoneme s accurately predicts the nonsense syllable score S_{cvc} , where

$$\begin{aligned} S_{cvc} &= c^2 v \\ &= s^3 \end{aligned}$$

*This is a necessary but insufficient condition for *independence*

- **These statistical models are highly accurate**
- **!!! Remember: This only applies to “nonsense words” !!!**

QUESTIONS?

THE NEXT STEP

- Next they dissected $s \equiv P_{correct}(phone)$ into frequency bands!

SPECIFIC DEFINITIONS

SYMBOL	DEFINITION
α	gain applied to the speech
$c(\alpha) \equiv P_c(\text{consonant} \alpha)$	consonant articulation
$v(\alpha) \equiv P_c(\text{vowel} \alpha)$	vowel articulation
$s(\alpha) = [2c(\alpha) + v(\alpha)]/3$	average phone articulation for CVC's
$e(\alpha) = 1 - s(\alpha)$	phone articulation error
f_c	high- and low-pass cut-off frequency
$s_L(\alpha, f_c)$	s for low-pass filtered speech
$s_H(\alpha, f_c)$	s for high-pass filtered speech

FLETCHER'S TWO BAND FORMULATION

- Split the speech into **low and high bands**, having articulations

$$s_L(\alpha, f_c) \text{ and } s_H(\alpha, f_c)$$

- **Fletcher** proposed a **linearizing transformation** of the phone articulations

$$(s_L) + (s_H) = (s)$$

- This is a nonlinear transformation of probabilities
- There was no guarantee that such a transformation exists
However, Fletcher's intuition was correct

WHAT THEY FOUND

- For nonsense {C,V} syllables the phone articulation transformation is:

$$s = \frac{\log(1 - e)}{\log(e_{min})},$$

with $e_{min} = 0.015$ (1.5% error, or 98.5% correct)

– This relationship must have taken years to discover!

- Solving for $e \equiv 1 - s(\mathcal{A})$:

$$e = e_{min}^{\mathcal{A}(s)} = e_{min}^{\mathcal{A}(s_L) + \mathcal{A}(s_H)} = e_{min}^{\mathcal{A}(s_L)} e_{min}^{\mathcal{A}(s_H)}$$

- In terms of the error probabilities $e = 1 - s$, $e_L = 1 - s_L$ and

$$e_L = 1 - s_L:$$

$$e = e_L e_H.$$

FLETCHER'S TWO BAND EXAMPLE

- If we have 100 spoken sounds, and 10 errors are made while listening to the low band, and 20 errors are made while listening to the high band, then

$$e = 0.1 \times 0.2 = 0.02,$$

namely 2 errors will be made when listening to the full band, so

$$s = 1 - 0.02 = 0.98$$

$$S = s^3 = 0.941$$

- This is an unexpected, simple, and amazing result
 - What does this mean? Why does it turn out this way?

DEMO of the the McGurk effect

THE FLETCHER-STEWART MULTI-CHANNEL MODEL

- Fletcher 1921 generalize the two-band case to $K = 20$ frequency bands

$$\begin{aligned} 1 - s &= e_1 e_2 \cdots e_k \cdots e_K \times e_{\text{visual}} \\ &= (1 - s_1)(1 - s_2) \cdots (1 - s_K) \times (1 - s_{\text{visual}}) \end{aligned}$$

where

$$e_i \equiv 1 - s_i$$

–This formula forms the basis of [articulation index](#) theory

–Why $K = 20$ bands?

Each band equals 1mm along the basilar membrane

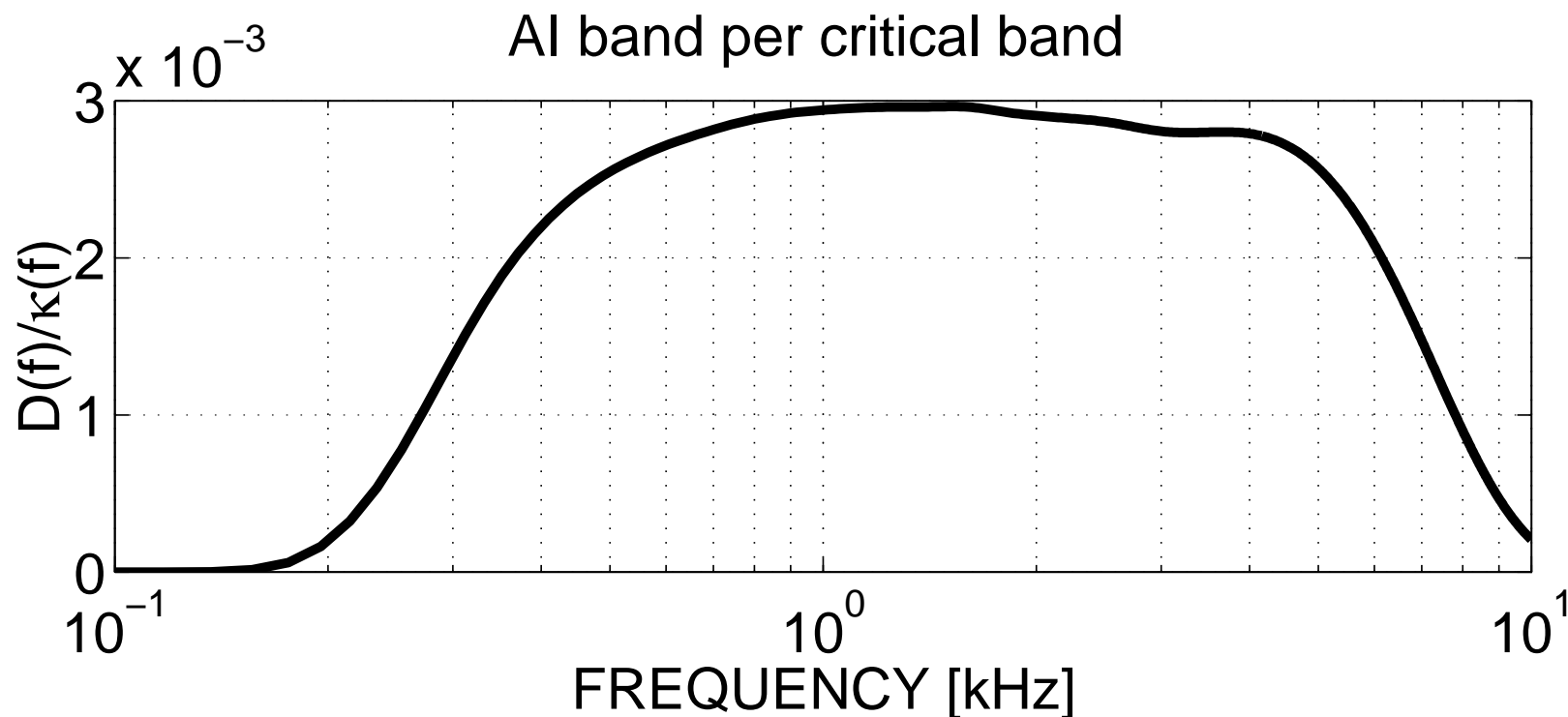
- I have added a visual channel, to account for the McGurk effect (Channel 21)
- Probability of error e_i models [events](#), as in the visual example

DENSITY OF ARTICULATION PER CRITICAL BAND

- This plot is the ratio of $D(f)/\kappa(f)$, where $D(f)$ is the articulation density

$$D(f_c) \equiv \frac{\partial \mathcal{A}_L}{\partial f_c}, \quad K \text{ AI bands}$$

$\kappa(f)$ = the critical ratio [\propto cochlear filter bandwidth (ERB)]



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

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

Thus

$$0 \leq SNR_k \leq 1.$$

- Total error:

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

- The speech SNR in dB (not the energy) determines the event errors e_k , and thus the phoneme articulation

$$s = 1 - e_1 e_2 \cdots e_K$$

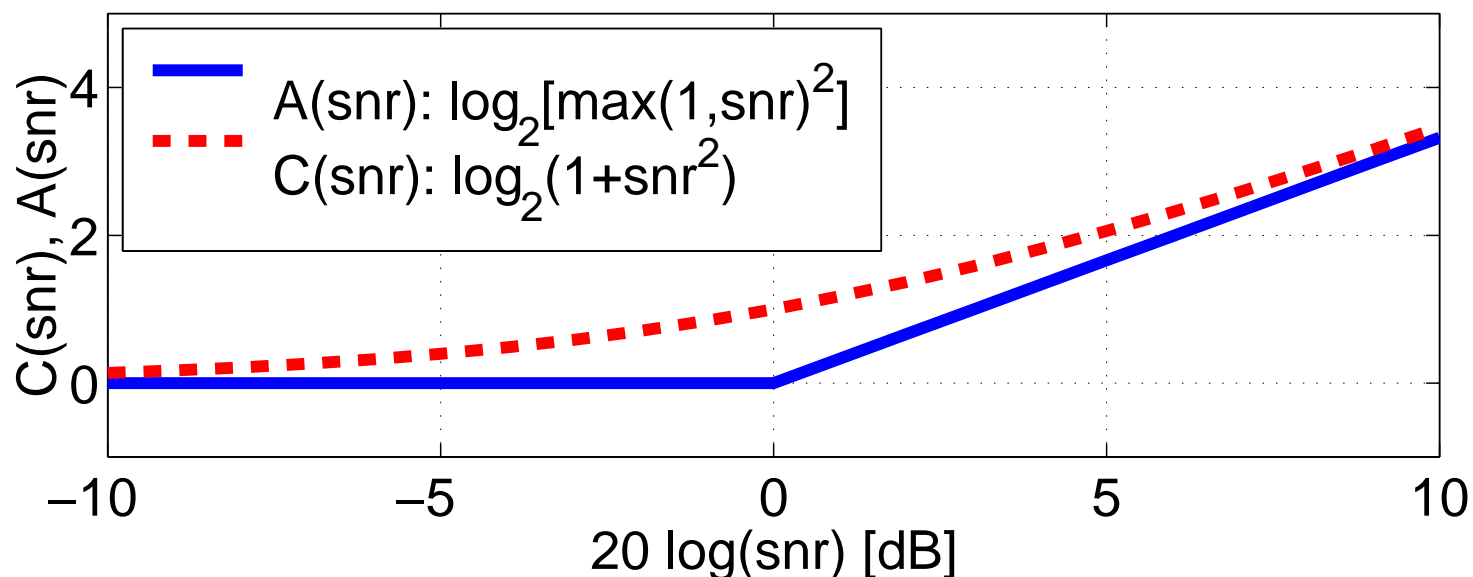
AI AS A CHANNEL CAPACITY

- Since $\sum_k (\log snr_k) = \log(\prod_k snr_k)$

$$\mathcal{A} \equiv \frac{1}{K} \sum_k SNR_k \propto \log \left(\prod_k snr_k \right)^{1/K} \quad (1)$$

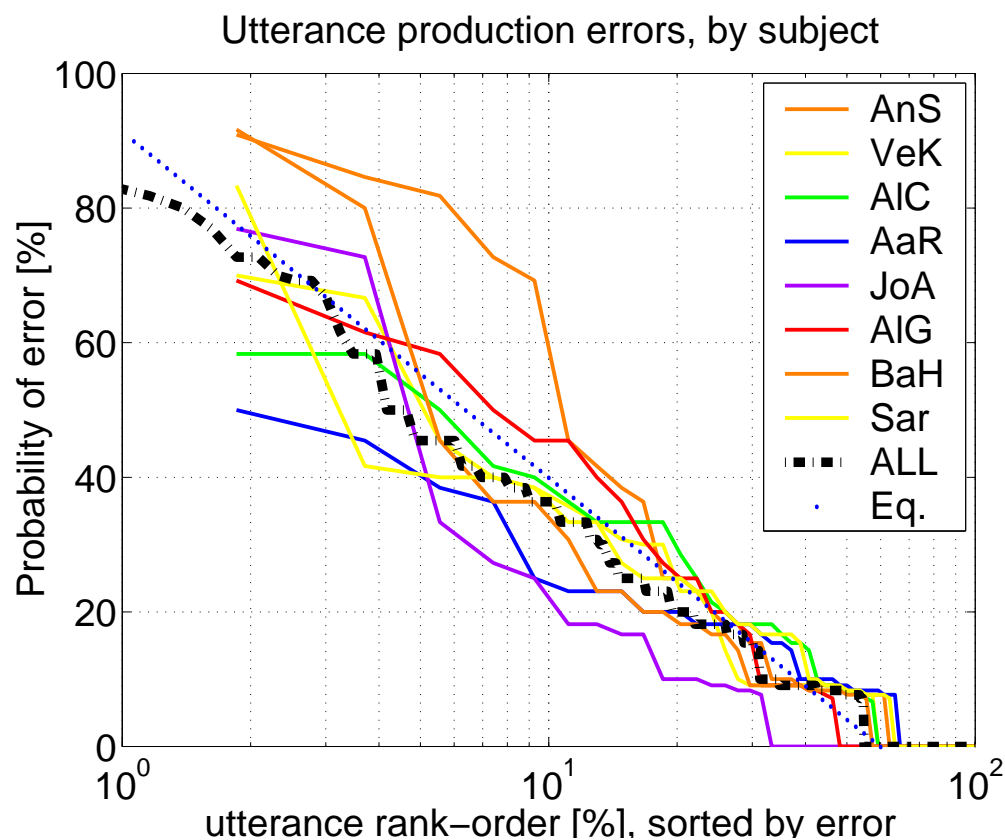
- and from Shannon (for the Gaussian channel)

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



TALKER PRODUCTION 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



SOURCES OF ERROR

- Talker production errors

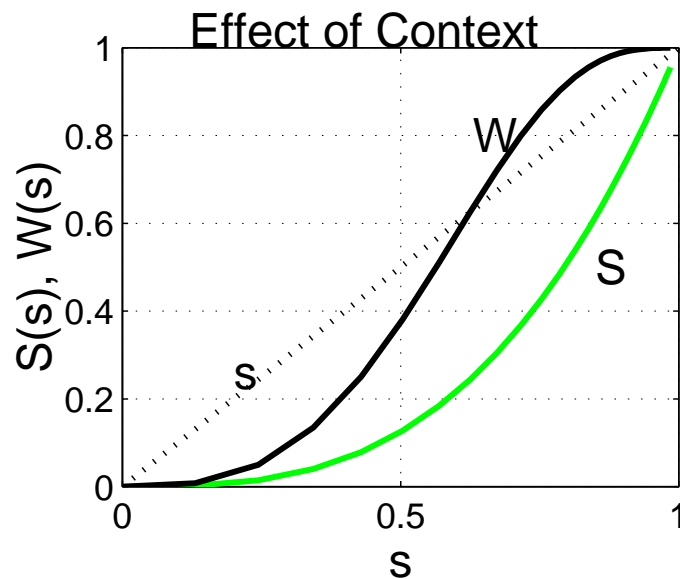
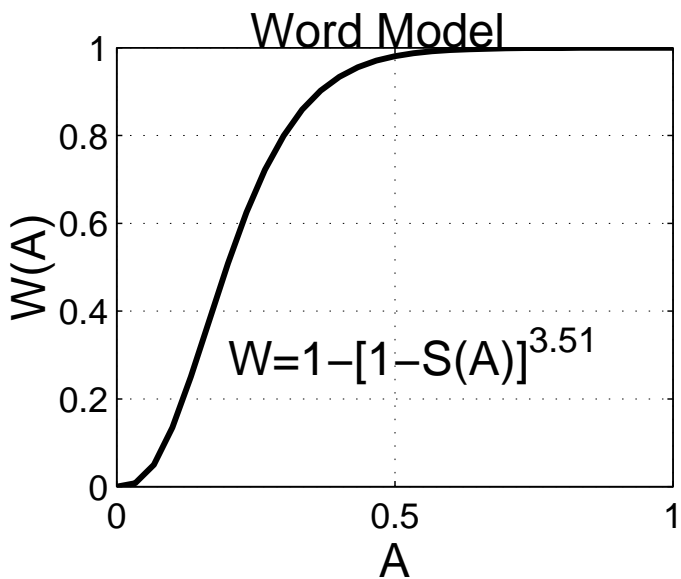
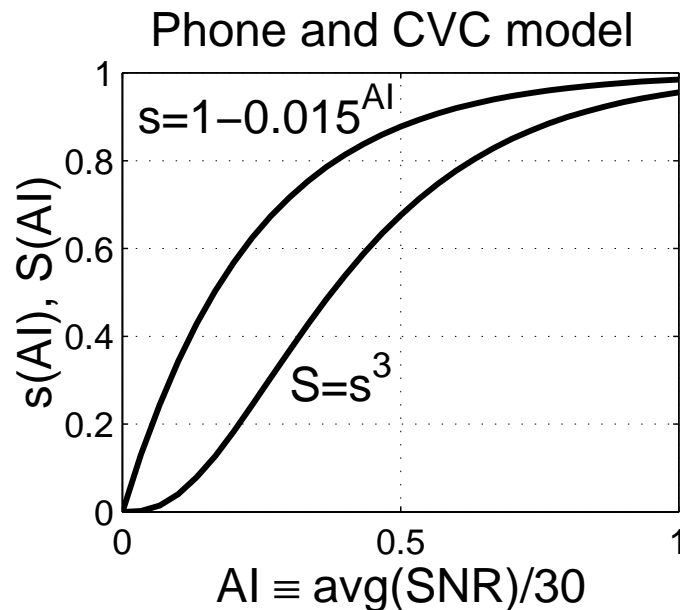
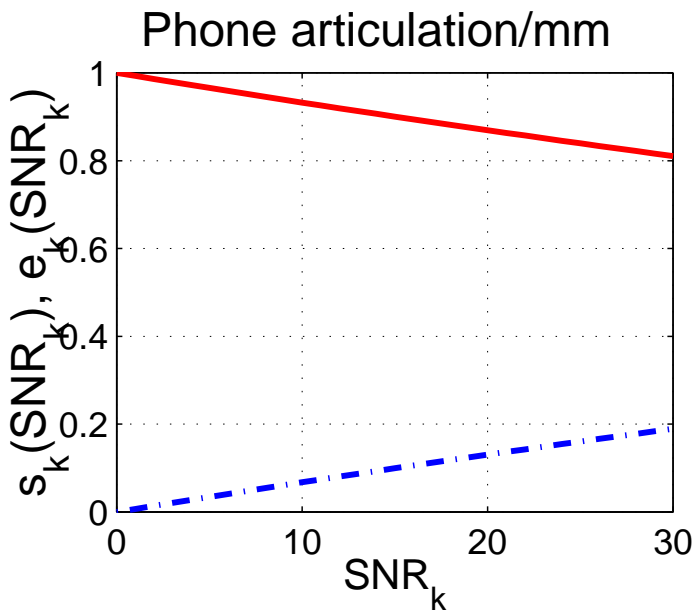
- Production errors are defined by **token utterances error** over listeners
- Once poor utterances are identified, they may be selectively removed
 - * This method allows us to control the gross error rate
- With this method we can obtain a 100% score in the clear
 - * The price for this is a reduced $N_{utterances}$

- Listener errors (after selectively removing production errors)

- Listener bias may be determined from individual confusion matrices
- This bias can be a function of the production error threshold
 - * The main effect is on L_2 listeners

EXAMPLE CALCULATIONS

Wide-band channel vs. SNR



THE RECOGNITION CHAIN

- The cochlear critical bandwidth defines the SNR_k
- The *event-error* model: $e_k \propto e_{\min}^{SNR_k}$ (SNR in dB units)
- The *average-phone articulation* model:

$$s = 1 - e_1 e_2 \cdots e_k \cdots e_K$$

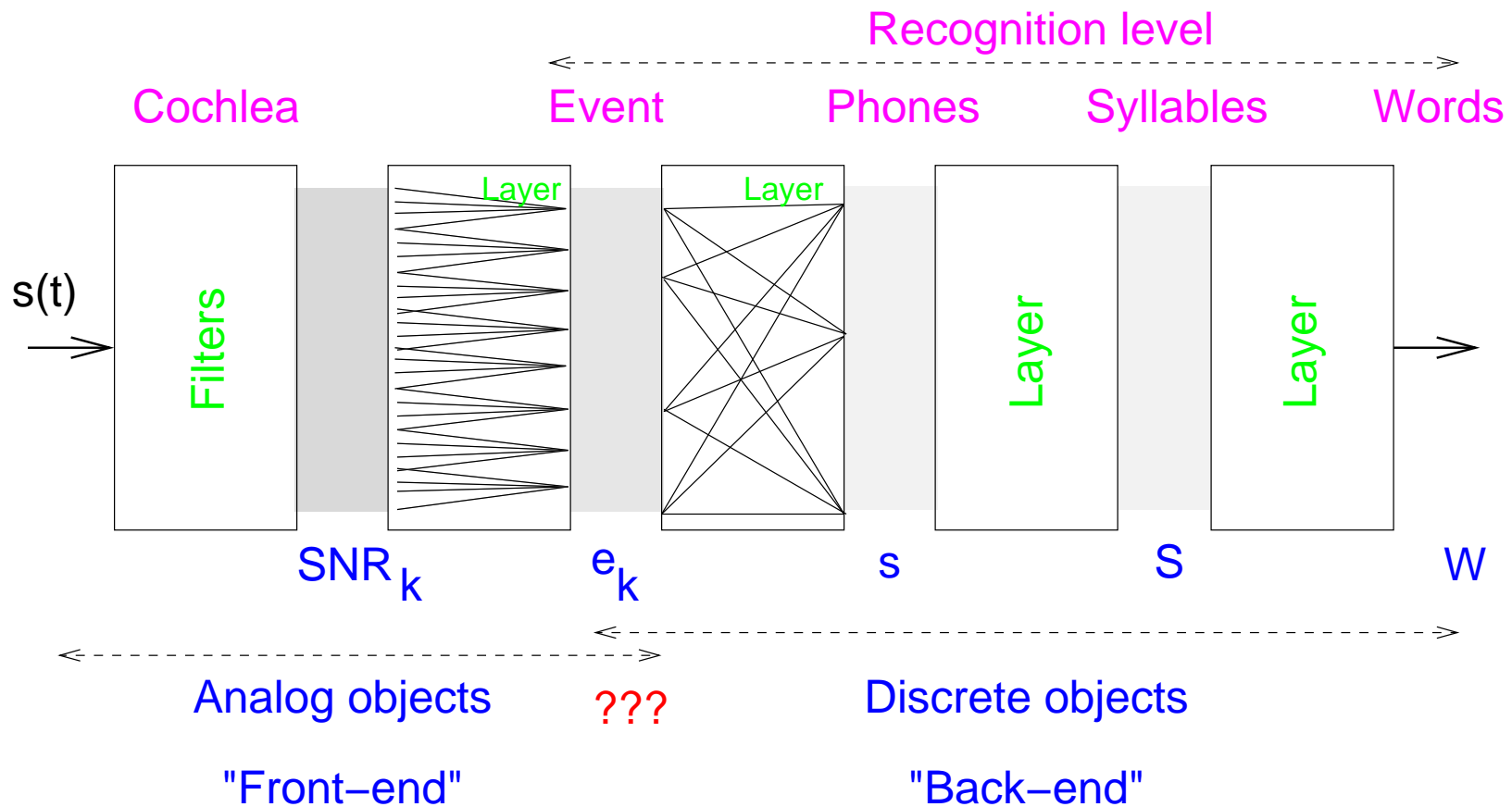
- The nonsense CVC *syllable articulation* model: $S = s^3$
- Heuristic *degree of freedom context models* Boothroyd (see discussion [Allen 1994](#))
 - Word: $W = 1 - (1 - S)^j$
 - Sentence: $I = 1 - (1 - W)^k$
 - Sentence with context: $C = 1 - (1 - I)^l$
- Layers of context:
 - j depends on the ratio of words to pseudo-words in the corpus,
 - k depends on the number of salient words in a sentences,
 - l depends on the word salience and topic context.

COMPOSITION LAWS

- Rules regarding $\prod_i P_{\text{error}}^{(i)}$ versus product $\prod_i P_{\text{correct}}^{(i)}$?
 - **Parallel processing:** $P_e = \prod_k e_k$
 - * Errors in many bands have no effect
 - * One band with small error (i.e., $e_k = 0$) dominates
e.g., $e = e_L e_H$; $e = e_1 e_2 \cdots e_K$; the McGurk example
 - **Serial processing:** $P_c = \prod_k s_k$
 - * All items of a string must be correct for success
e.g., $S_{cvc} = cvc \approx s^3$; $S_{cv} \approx s^2$
- HSR seems to be a problem in **combinatorics**,
of elementary **pre-phonetic** events.

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

SUMMARY OF MODEL RESULTS

- Hierarchical probability relations:

band *SNR* →

band errors (events) →

phoneme errors →

syllable errors →

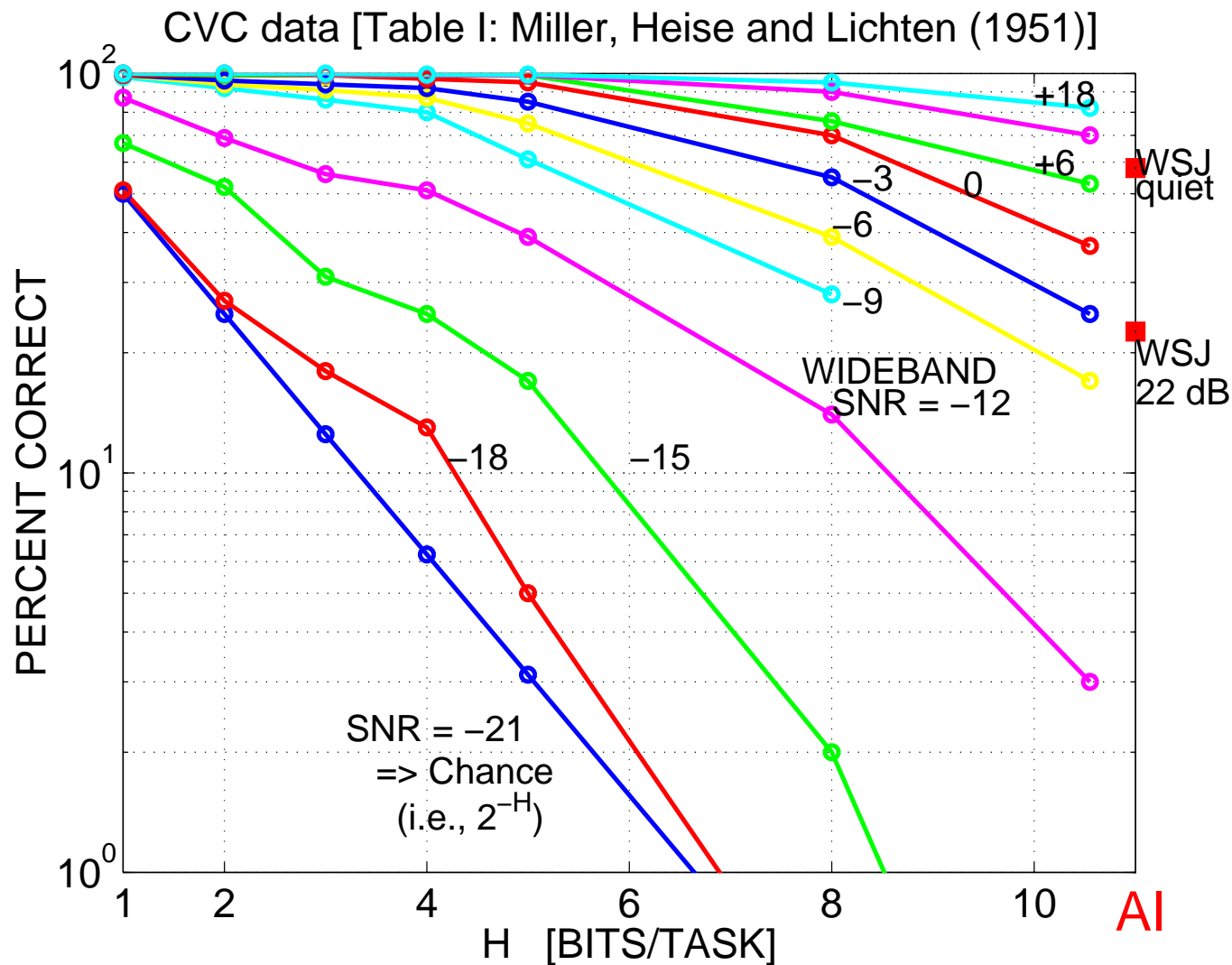
nonsense word errors →

true word errors, etc.

- The HSR error is established **well before** language is accessed!
HSR error depends only on the *SNR* in bands

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



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

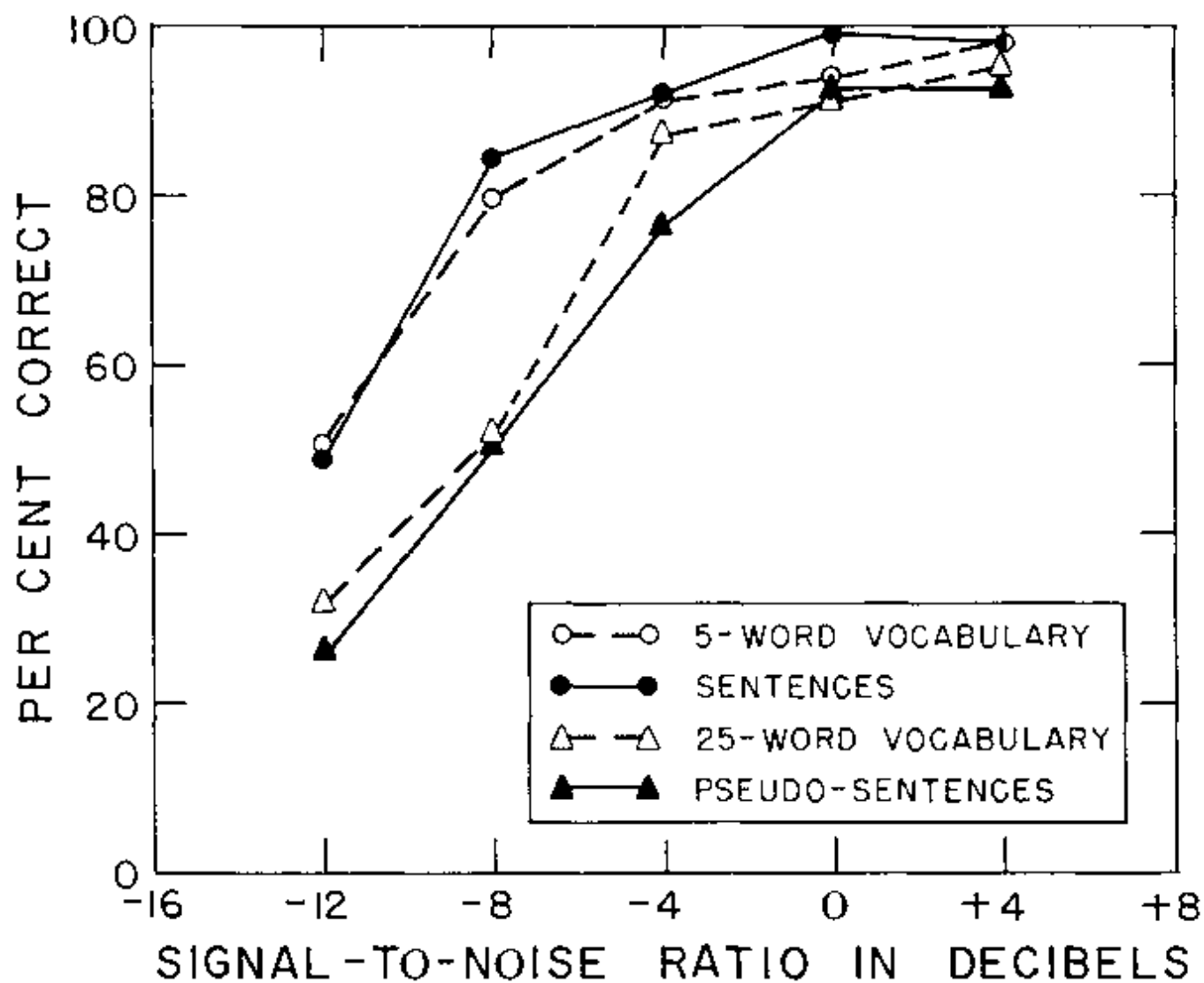
Example: **He left no black socks**

25 words reverse order

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

GRAMMATICAL CONTEXT

- Results of tests

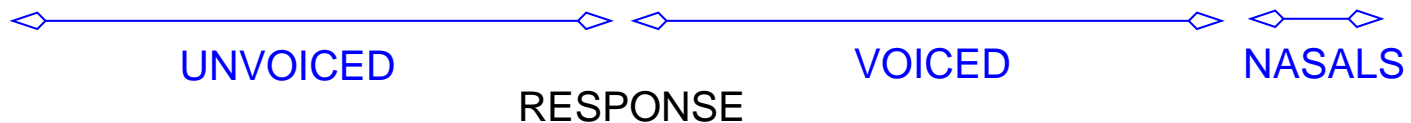


CONFUSION MATRIX PARTITIONING

- 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

MILLER'S BINARY FEATURES

- Miller & Nicely derived binary consonant features [i.e., events]

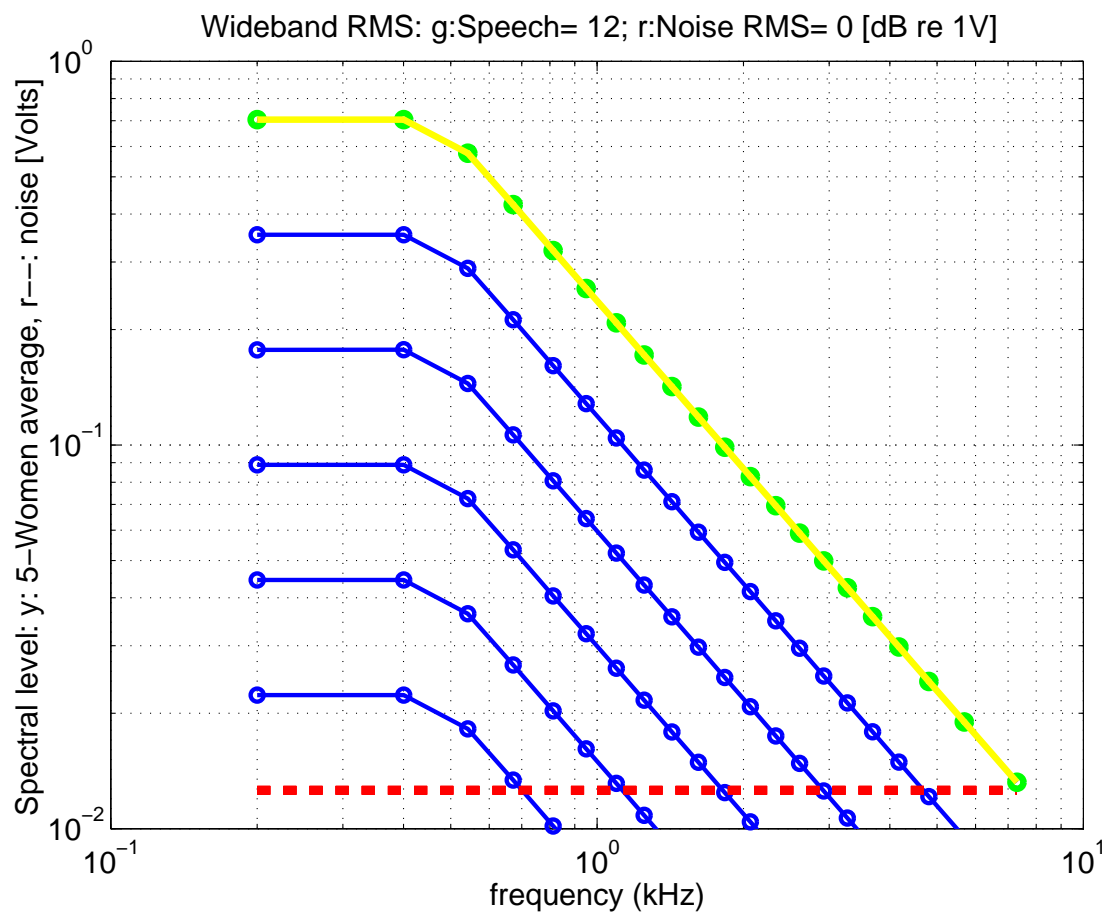
TABLE XIX. Classification of consonants used to analyze confusions.

Consonant	Voicing	Nasality	Affrication	Duration	Place
<i>p</i>	0	0	0	0	0
<i>t</i>	0	0	0	0	1
<i>k</i>	0	0	0	0	2
<i>f</i>	0	0	1	0	0
<i>θ</i>	0	0	1	0	1
<i>s</i>	0	0	1	1	1
<i>ʃ</i>	0	0	1	1	2
<i>b</i>	1	0	0	0	0
<i>d</i>	1	0	0	0	1
<i>g</i>	1	0	0	0	2
<i>v</i>	1	0	1	0	0
<i>ð</i>	1	0	1	0	1
<i>z</i>	1	0	1	1	1
<i>ʒ</i>	1	0	1	1	2
<i>m</i>	1	1	0	0	0
<i>n</i>	1	1	0	0	1

“... the *impressive thing* to us was that ... the *[binary] features* were *perceived* almost *independently* of one another.” –Miller & Nicely 1955

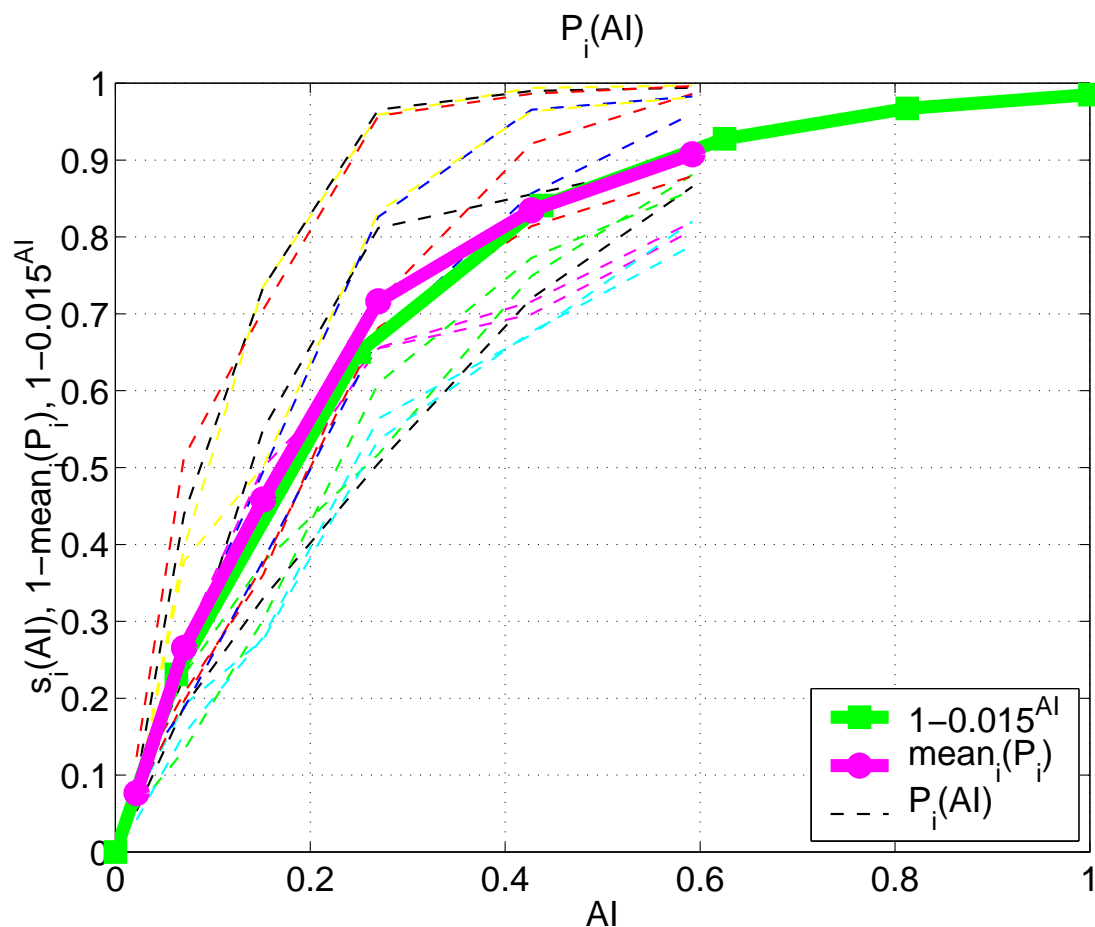
FINDING THE AI FOR MILLER NICELY TALKERS

- Average spectrum for female talkers



TRACE OF MILLER NICELY AND THE AI

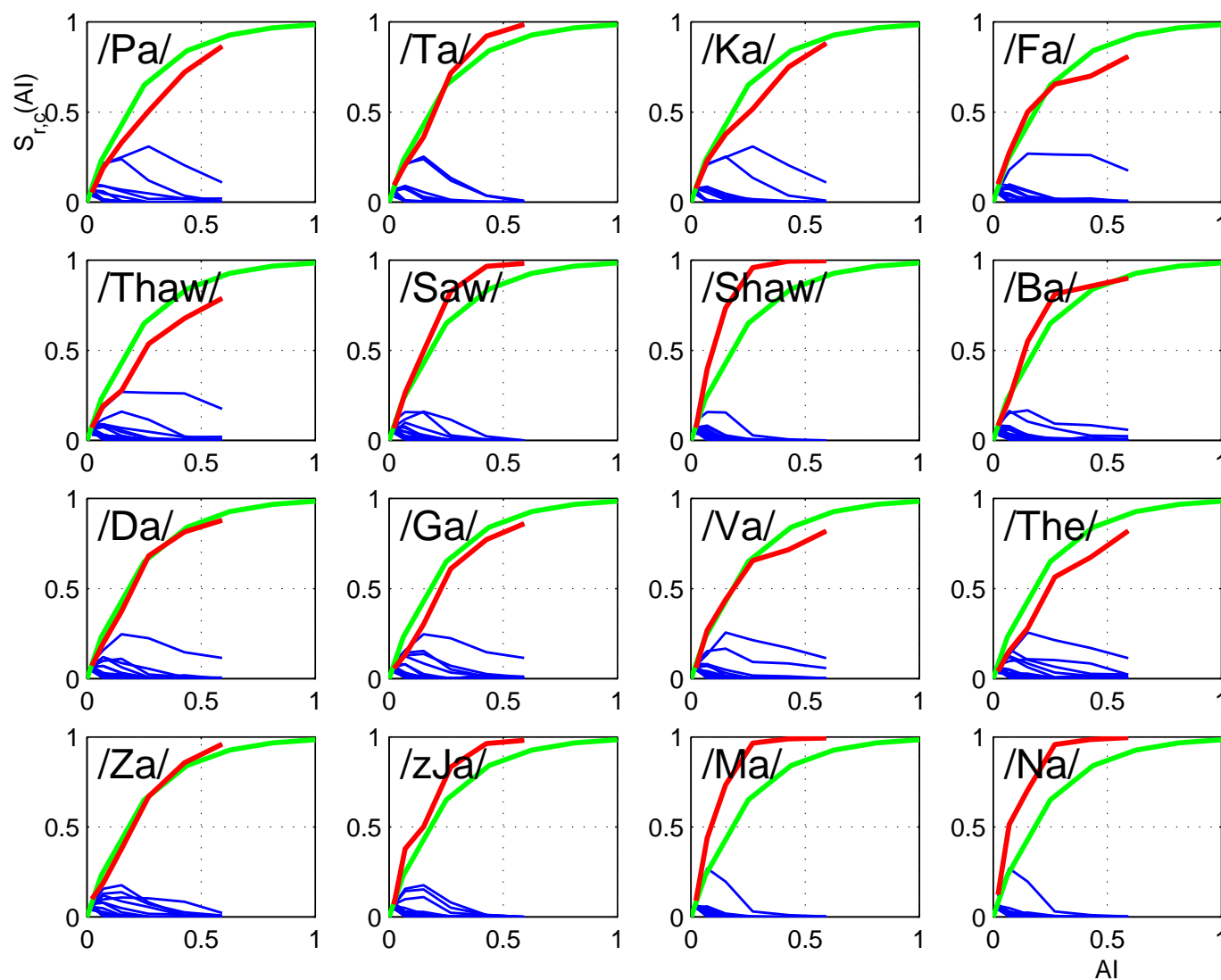
- Next we look at the average PI function vs. AI



SYMMETRIC COMPONENT OF $P_C(SNR)$

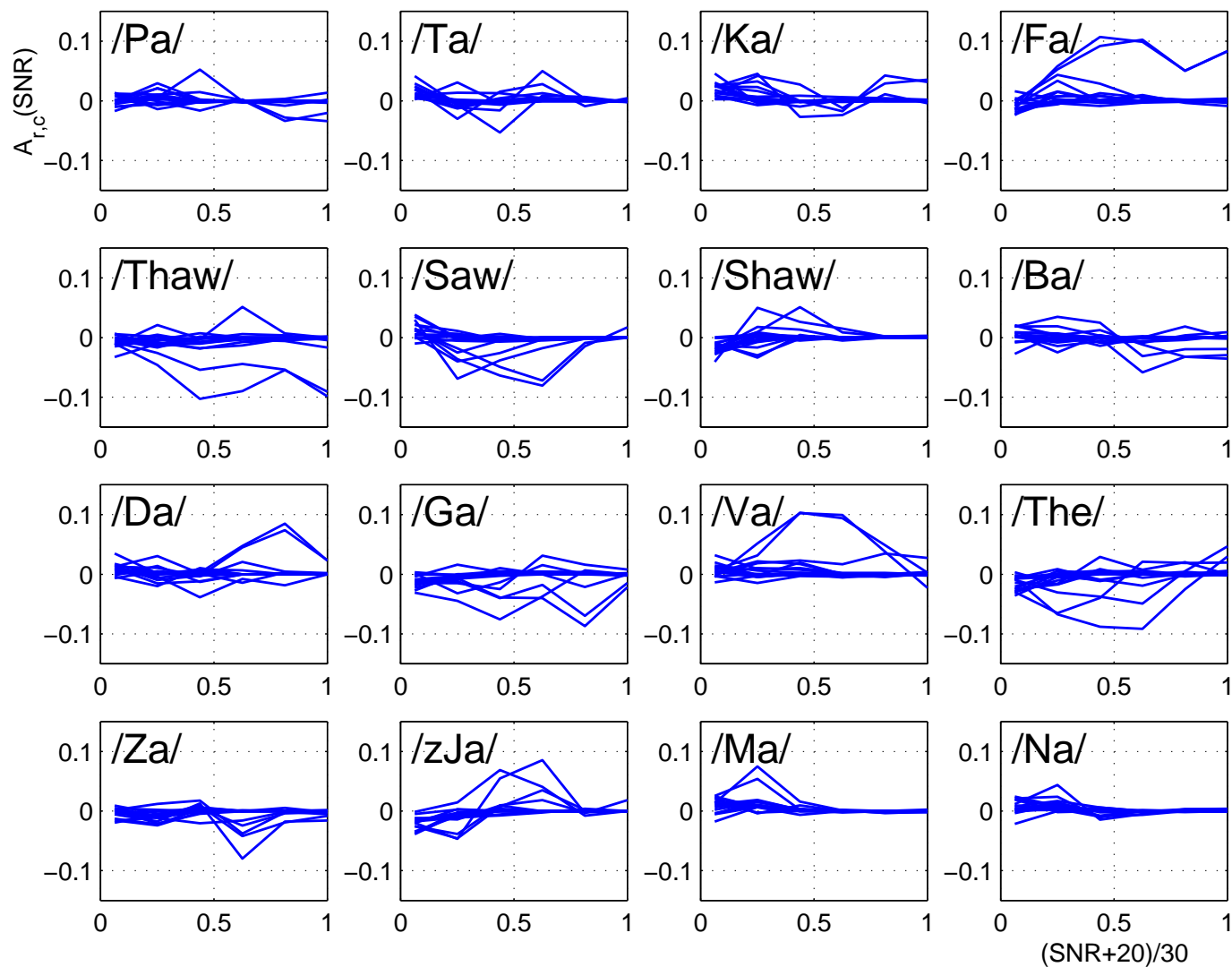
- We stand to learn from linear operations on $P_{ij}(snr)$

Symmetric: $S_c(snr) \equiv [P_{ij}(snr) + P_{ji}(snr)]/2$



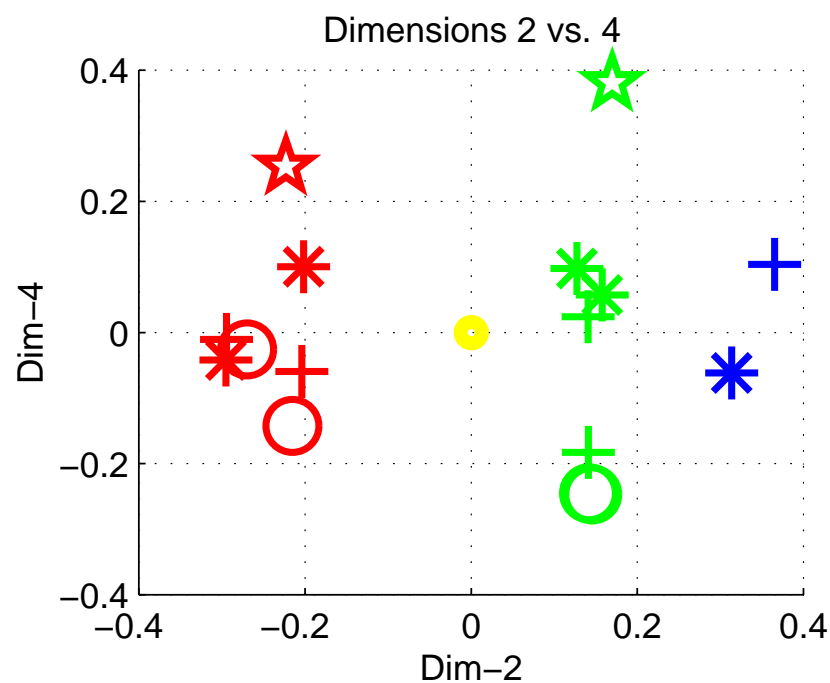
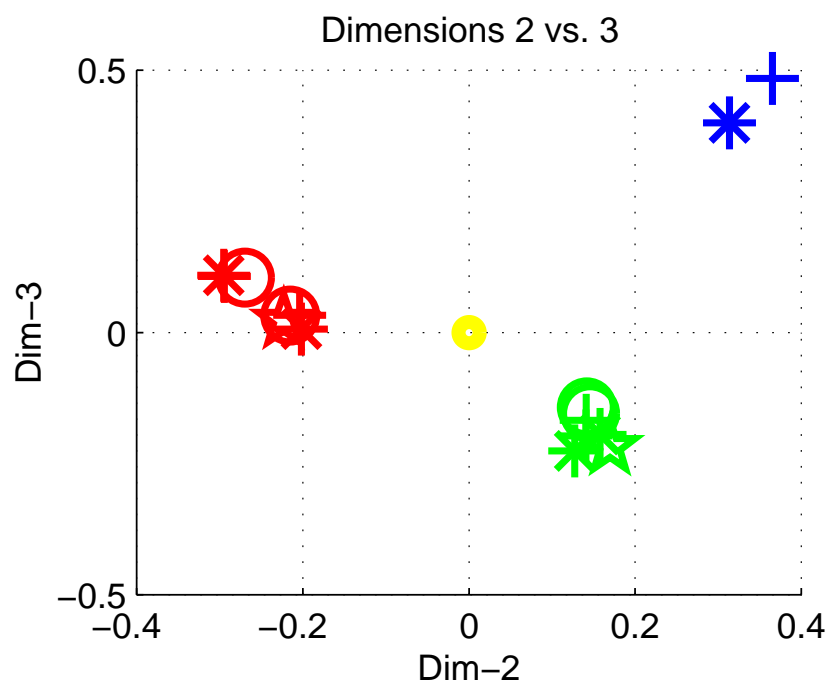
SKEW-SYMMETRIC COMPONENT OF $P_C(SNR)$

- Skew: $A_c(snr) \equiv [P_{ij}(snr) - P_{ji}(snr)]/2$



SVD REPRESENTATION OF THE PERCEPTUAL SPACE

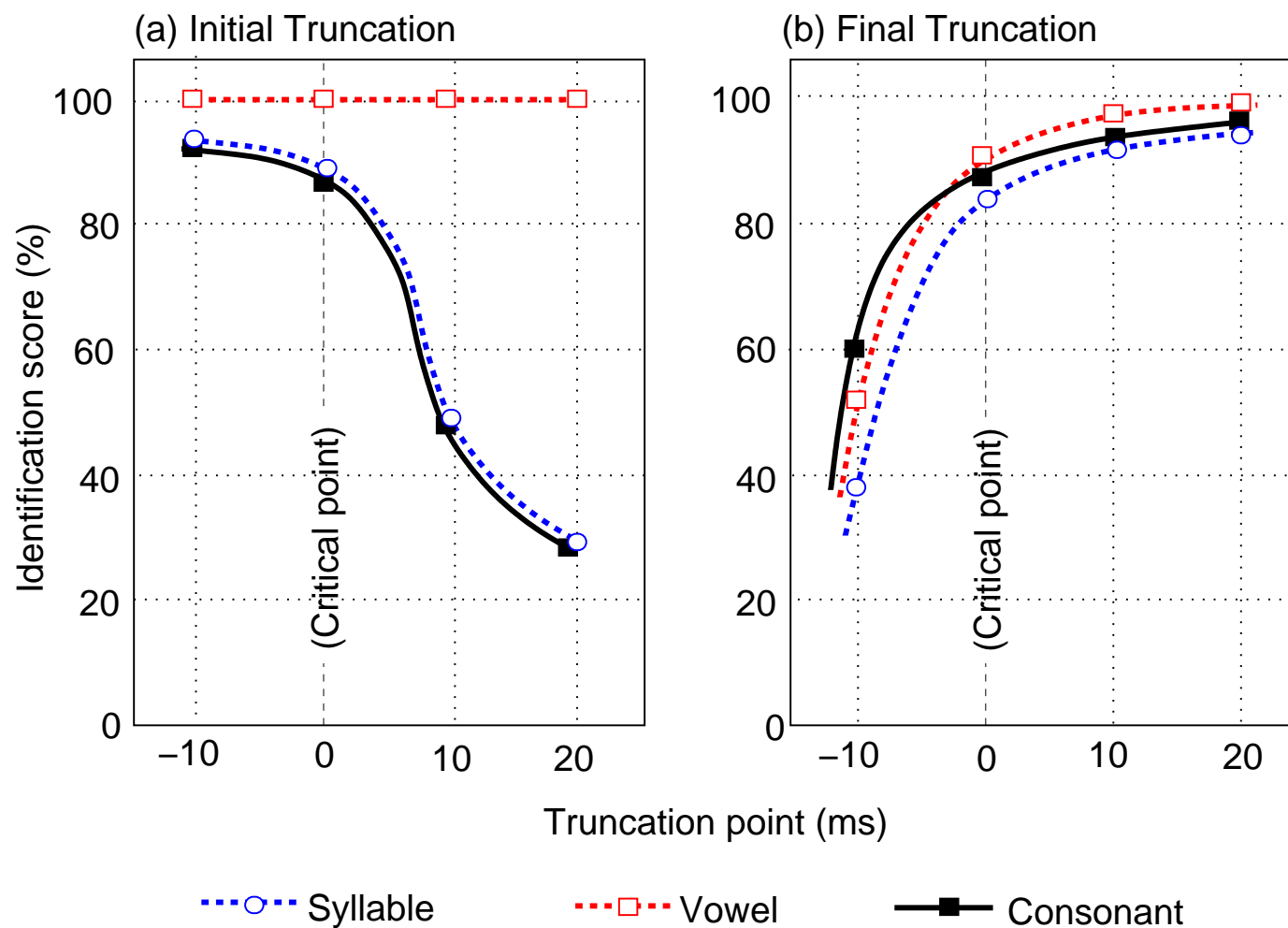
- 4^{dim} SVD perceptual representation of the confusion matrix



DEMO

TEMPORAL RESOLUTION OF PHONE RECOGNITION

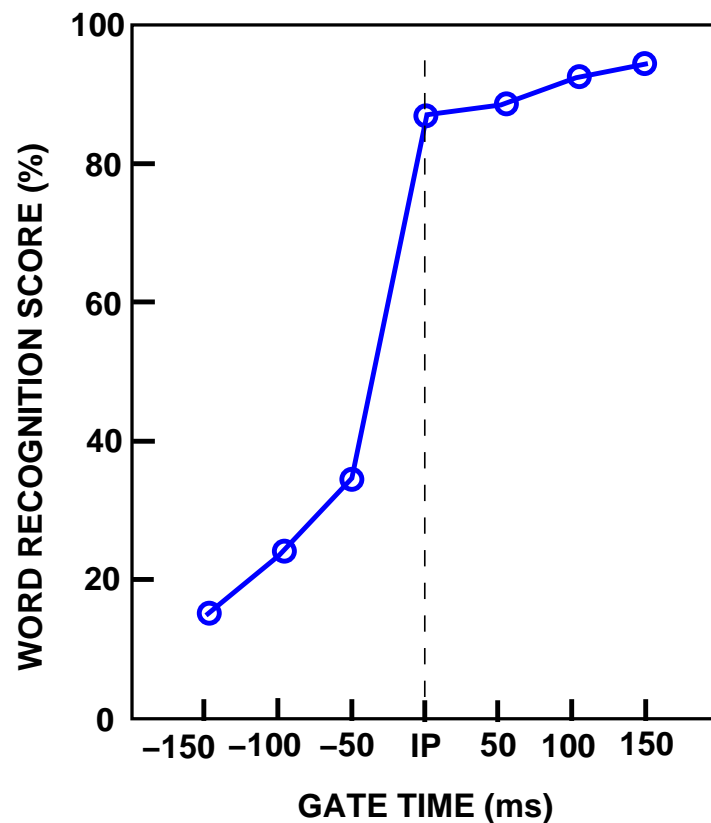
- Phones are recognized in on a 10 ms time scale (Furui 1986)



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*.”

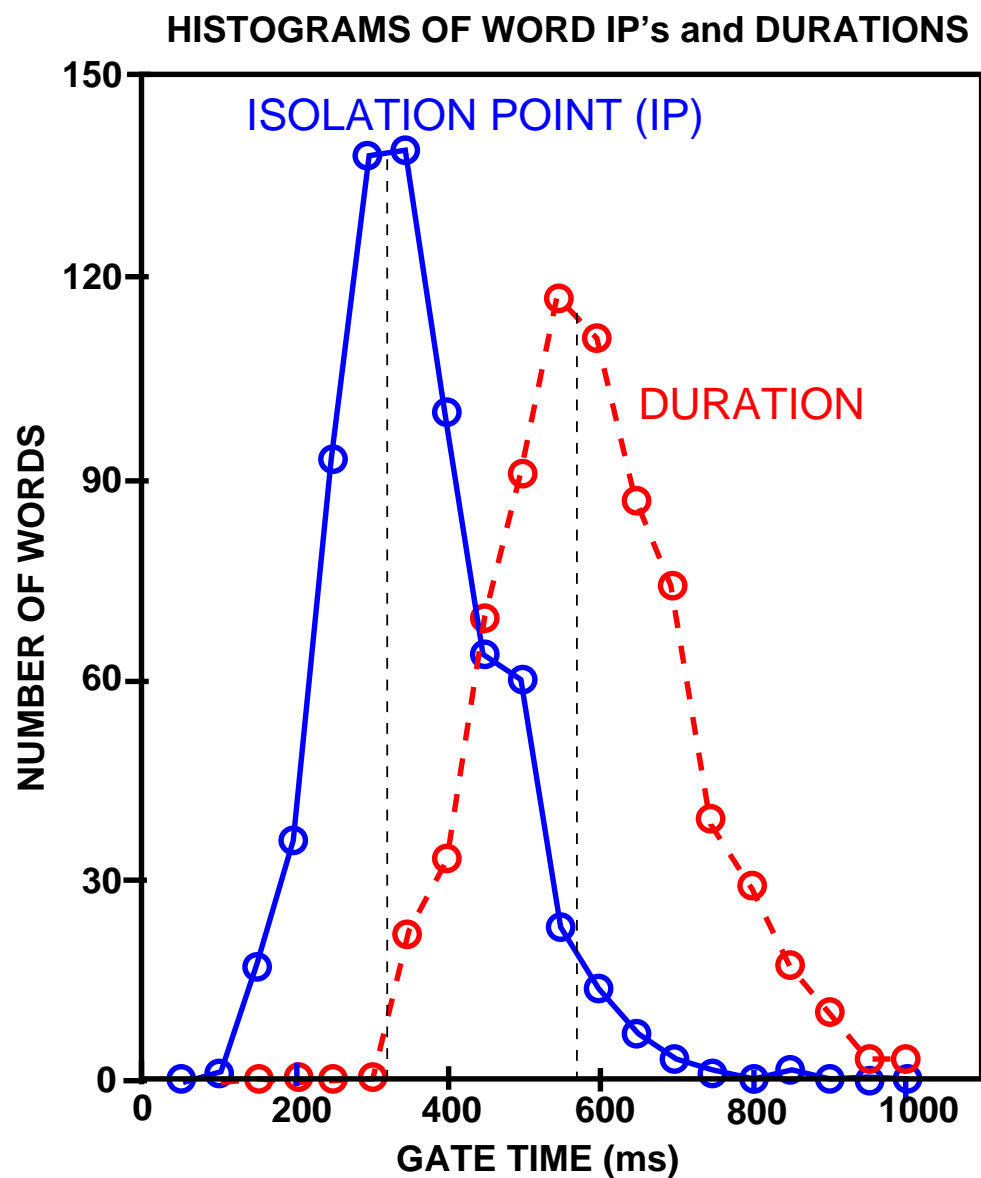
ACCURACY OF IDENTIFICATION VERSUS GATE TIME



- Categorical perception

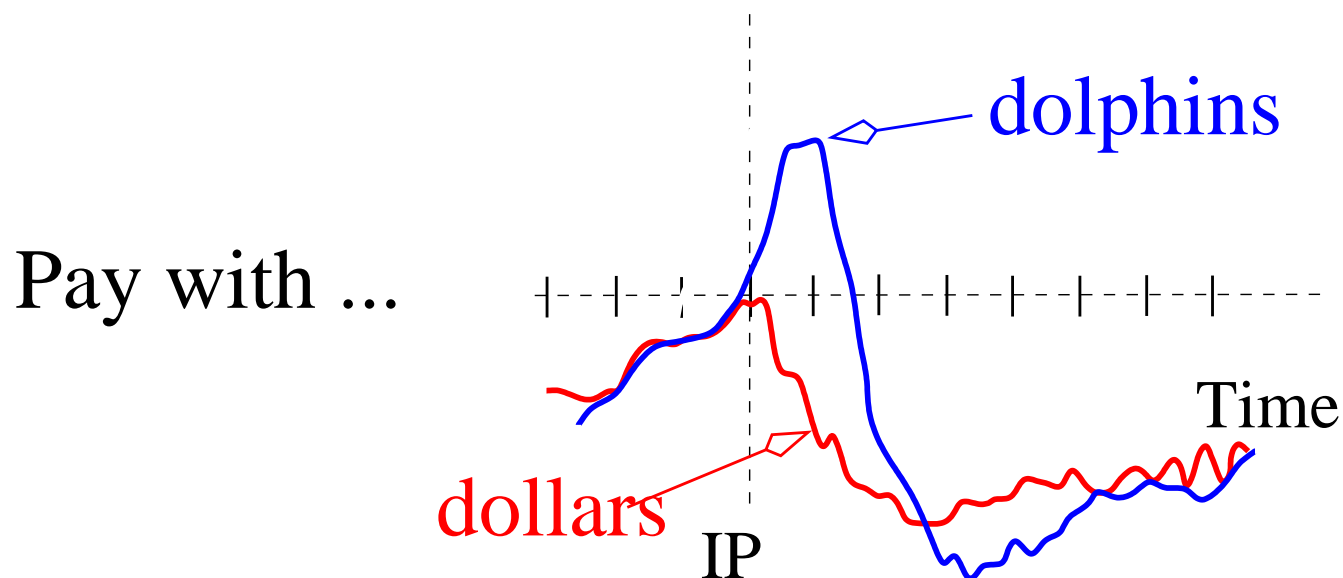
WORD SEMANTICS: IP VS. DURATION

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



ERP MEASURE OF CONTEXT RE IP

- Expt. II – Event related scalp potential (N-400 ERP) re IP, from Exp. I
Sentence semantics effects



— *Cohort congruous* dollars
 — *Cohort incongruous* dolphins

- Words are recognized on a syllable by syllable basis, within 50 ms
- Context is recognized on a syllable by syllable basis, within 200 ms

FROM CONTINUOUS TO DISCRETE



- Φ -domain signals

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

- Ψ -domain objects

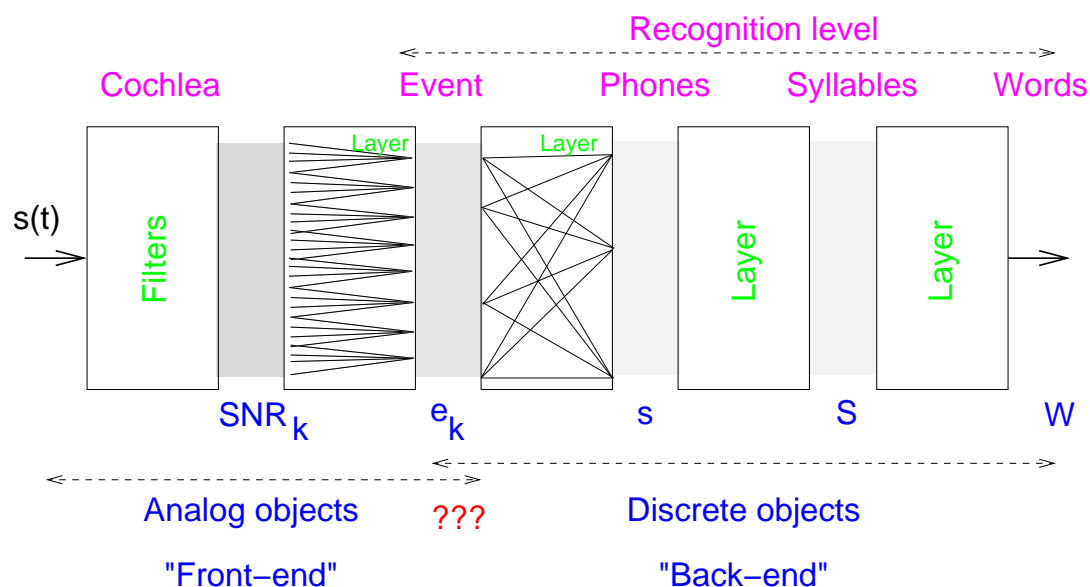
Words
Syllables
Phonemes
Events [Miller's features]

CATEGORIAL PERCEPTION

- Meaningful words are recognized before they end
- Syllables are recognized within 50 ms

SUMMARY

- Miller & Nicely found 5 independent channels, described by discrete **events** [Miller's features]
- Speech is recognized in layers:
 $SNR_k \Rightarrow \text{events} \Rightarrow \text{phones} \Rightarrow \text{syllables} \Rightarrow \text{words} \Rightarrow \dots$



- Language model performance is independent of noise robustness!
- To study HSR, entropy must be controlled
- Speech psychophysics is an important tool for studying HSR

FUTURE GOALS

- Use **psychophysics** to gain **insight** into **event** extraction
- The next break through:
 - More robust ASR
 - An **event** extracting hearing aid

This talk may be found at:

<http://auditorymodels.org/jba/PAPERS/ICASSP/>

References

- Allen, J. B. (1994). "How do humans process and recognize speech?," *IEEE Transactions on speech and audio* **2**(4):567–577.
- Boothroyd, A. (1968). "Statistical theory of the speech discrimination score," *J. Acoust. Soc. Am.* **43**(2):362–367.
- Boothroyd, A. and Nittrouer, S. (1988). "Mathematical treatment of context effects in phoneme and word recognition," *J. Acoust. Soc. Am.* **84**(1):101–114.
- Bronkhorst, A. W., Bosman, A. J., and Smoorenburg, G. F. (1993). "A model for context effects in speech recognition," *J. Acoust. Soc. Am.* **93**(1):499–509.
- Bronkhorst, A. W., Brand, T., and Wagener, K. (2002). "Evaluation of context effects in sentence recognition," *J. Acoust. Soc. Am.* **111**(6):2874–2886.
- Campbell, G. A. (1910). "Telephonic intelligibility," *Phil. Mag.* **19**(6):152–9.
- French, N. R. and Steinberg, J. C. (1947). "Factors governing the intelligibility of speech sounds," *J. Acoust. Soc. Am.* **19**:90–119.
- Furui, S. (1986). "On the role of spectral transition for speech perception," *Journal of the Acoustical Society of America* **80**(4):1016–1025.
- Miller, G. A. (1951). *Language and communication*. McGraw Hill, New York.
- Miller, G. A., Heise, G. A., and Lichten, W. (1951). "The intelligibility of speech as a function of the context of the test material," *J. Exp. Psychol.* **41**:329–335.
- Miller, G. A. and Nicely, P. E. (1955). "An analysis of perceptual confusions among some english consonants," *J. Acoust. Soc. Am.* **27**(2):338–352.
- Rayleigh, L. (1908). "Acoustical notes – viii," *Philosophical Magazine* **16**(6):235–246.
- Shannon, C. E. (1948). "The mathematical theory of communication," *Bell System Tech. Jol.* **27**:379–423 (parts I, II), 623–656 (part III).
- Shannon, C. E. (1951). "Prediction and entropy of printed english," *Bell System Tech. Jol.* **30**:50–64.
- Van Petten, C., Coulson, S., Rubin, S., Planten, E., and Parks, M. (1999). "Time course of word identification and semantic integration in spoken language," *J. of Exp. Psych.: Learning, Memory and Cognition* **25**(2):394–417.
- Wang, M. D. and Bilger, R. C. (1973). "Consonant confusions in noise: A study of perceptual features," *J. Acoust. Soc. Am.* **54**:1248–1266.