

# Cochlear nonlinearities and phoneme recognition

*Finding the features in individual consonants.*

Jont Allen

Feipeng Li

**Univ. of IL, Beckman Inst., Urbana IL**



# Notable quotes

*We need to know more about human speech processing and natural speech variation*

–Sadaoki Furui (ASRU 2009)

*This is so true!*

–Jont Allen

Question your assumptions:

- Elephant in the room:  
Human CV speech is *not* variable.
- CV speech is *not* redundant.
- Why we don't know anything about the topic?  
\$ spent on basic speech research  $\rightarrow 0$ ,

# Outline of talk

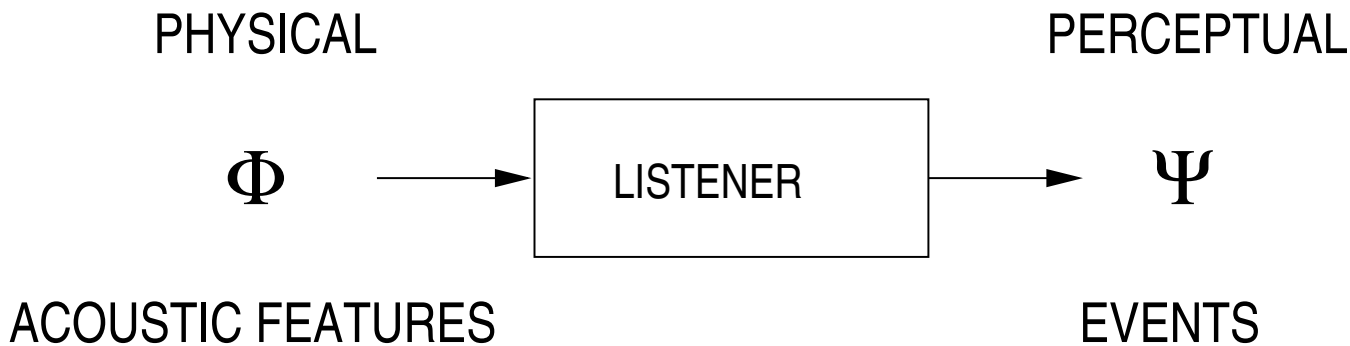
1. Intro + Objectives (5 mins)
  - The research goal is to
    - Identify the *elemental HSR events* in
    - Example consonants
2. Historical overview (5 mins  $\Sigma 10$ )
  - Rayleigh (1910) to Shannon (1948)
3. Methods (15 mins  $\Sigma 25$ )
  - -Information Theory; -Signal processing
  - -Psychophysics; -Articulation Index;
4. Results (30 mins  $\Sigma 55$ )
  - Confusions; Primes and Morphs;
  - Speech Modifications; Conflicting cues
5. Summary + Conclusions (5 mins  $\Sigma 60$ )

# I – Introduction (5 mins)

- Statement of the problem:
  - A fundamental understand the Human Speech code
- Short-term Goal:
  - Identify the key features in **individual** CV utterances
    - -Plosives (e.g., /p, t, k/ and /b, d, g/)
    - -Fricatives (e.g., /θ, ʃ, tʃ, s, h, f/ and /z, ʒ, v, ð/)
    - -With vowels /o, ε, ɪ/
- Applications:
  - **Reduce variability** in ASR at frontend
  - Hearing Aids, Cochlear Implants
  - Smart Telcom products
  - TTS (Text to speech)
  - Intelligibility modifications (Robustness problem)
    - Speech enhancement in noise

# Objective

- To develop rigorous procedures for analyzing and modifying speech in noise
- To identify perceptual features, denoted **events**



- Based on two basic measures:
  - AI-Gram (speech audibility measure)
  - Confusion matrix (CV discrimination measure)
- We will show that **onset and durational timing cues** form the consonant events

## II – Historical HSR Studies (5 mins)

- Lord Rayleigh's 1908 and George Campbell 1910
  - First electronic articulation experiments
- Harvey Fletcher's 1921 Articulation Index AI
  - Accurate predictions of nonsense syllable scores
  - French and Steinberg 1947 first publish AI
- Shannon The theory of Information 1948+
  - G.A. Miller, Heise and Lichten Entropy  $\mathcal{H}$  1951
  - G.A. Miller & Nicely CM  $P_{h|s}(SNR)$  1955
- Context:
  - G.A. Miller 1951 *Language and communication*
  - G.A. Miller 1962 5-word Grammer  $\equiv$  4 dB of SNR
  - Boothroyd JASA 1968; Boothroyd & Nittrouer 1988
  - Bronkhorst et al. JASA 1993

# Speech feature research

- 1910-1980: Bell Labs
- 1940-1960: Haskins Lab
- 1960-1990: MIT
- 1980-2010: ASR at AT&T, IBM, BBN, University research

## Cochlear research

- 1910-1950: Bell Labs
- 1960-2010: MIT + Harvard HSTB
- 1980-2010: NIH funded University research

# Speech feature research

- 1910-1980: Bell Labs
- 1940-1960: Haskins Lab **Synthetic speech**
- 1960-1990: MIT **Consonant features unknown**
- 1980-2010: ASR at AT&T, IBM, BBN, University research **Not designed to be robustness to noise**

## Cochlear research

- 1910-1950: Bell Labs
- 1960-2010: MIT + Harvard HSTB
- 1980-2010: NIH funded University research



# III – Methods (15 mins)

- Information Theory **IT**  $\equiv$  Articulation index **AI**
  - Confusion matrix **CM** scores:  $P_{h|s}(SNR)$
  - **AI** to model mean phone errors  $\sum_h P_{h|s}(SNR)$
- Psychophysics
  - Real consonant-vowel CV speech
  - Several types of additive noise
  - Large number of trials
    - >20 talkers and >20 listeners
- Signal processing
  - AI-gram (crude cochlear model)
  - Frequency, time, intensity truncation **3<sup>d</sup>-DS**
  - Short-Time Fourier Transform **STFT** modifications

# The CM and the *elemental-event*

- Miller-Nicely's 1955 articulation matrix  $P_{h|s}(SNR)$ , measured at [-18, -12, -6 shown, 0, 6, 12] dB SNR

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>ð</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>ð</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      |

UNVOICED

VOICED

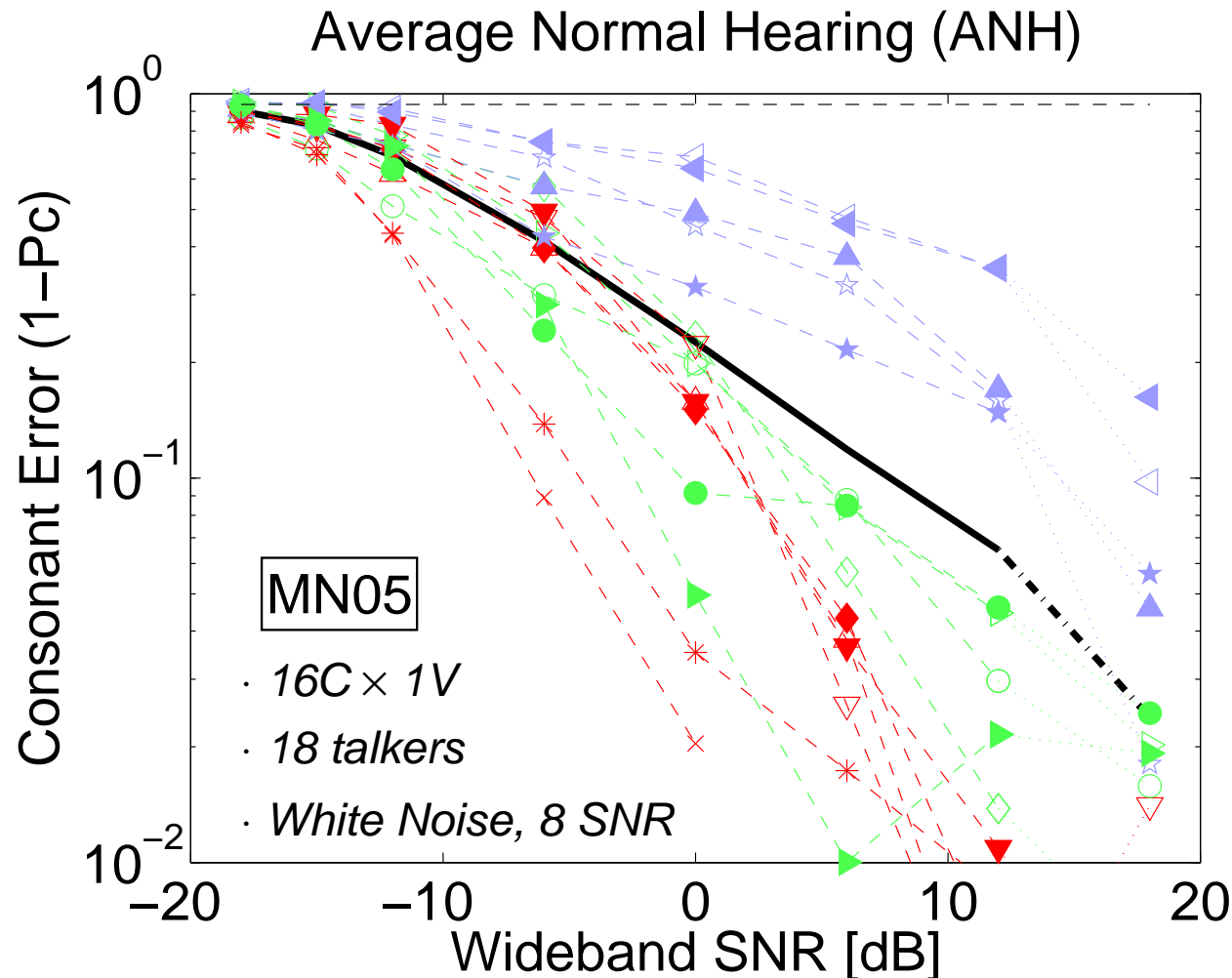
NASAL

RESPONSE

- Confusion groups  $\equiv$  *inhomogeneous elemental-events*

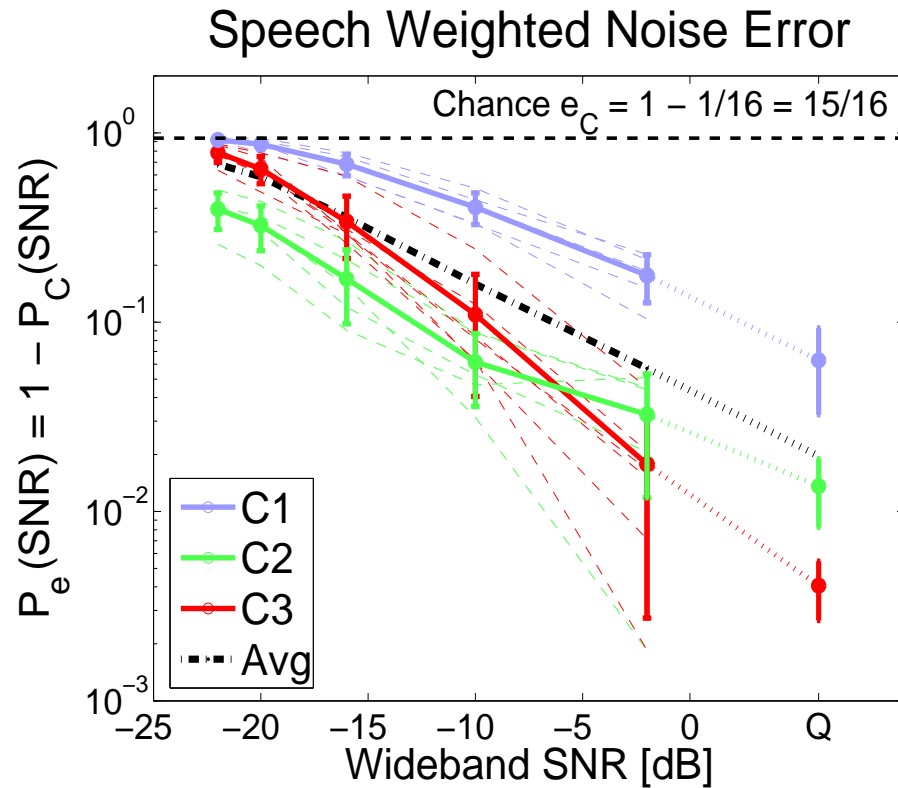
# Average phone scores vs. SNR

- Consonant chance performance is -20 dB-SNR in **white noise** Phatak Allen, 2007

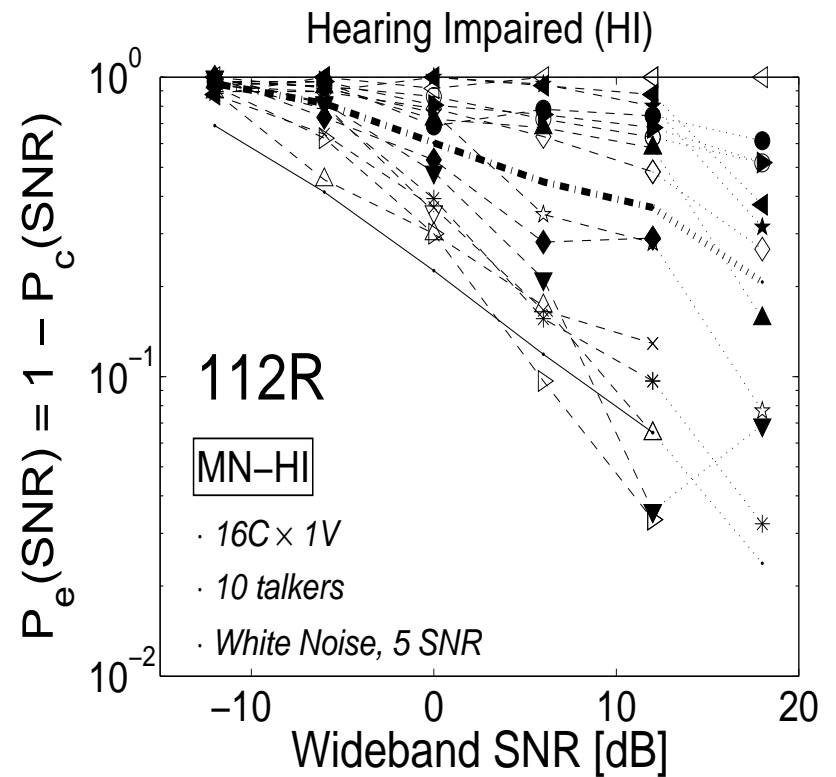


# Consonant Variability

- Avg. Consonant error  $P_{h|s}(SNR)$  strongly **heterogeneous!**
- NH listeners above chance at  $< -25$  dB SNR in SWN
- HI  $P_e(SNR) \gg ANH P_e(SNR)$



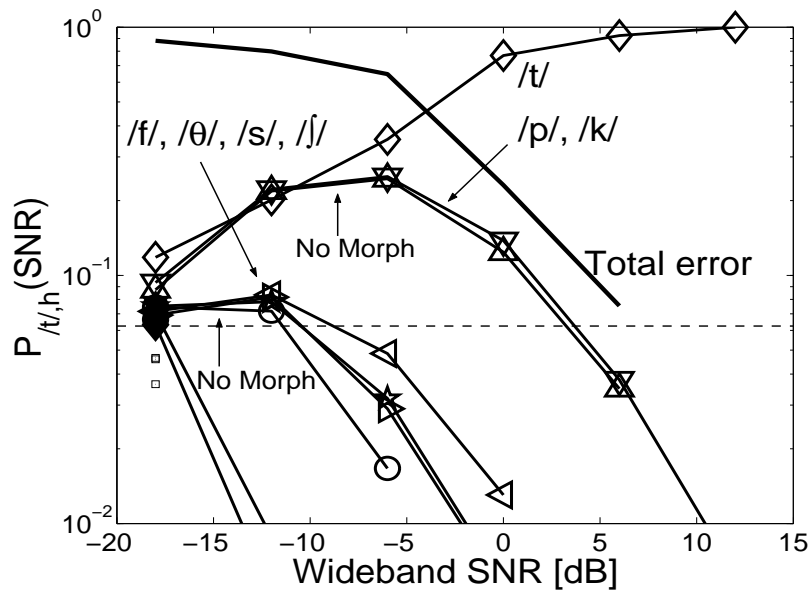
(a) Consonant errors ANH



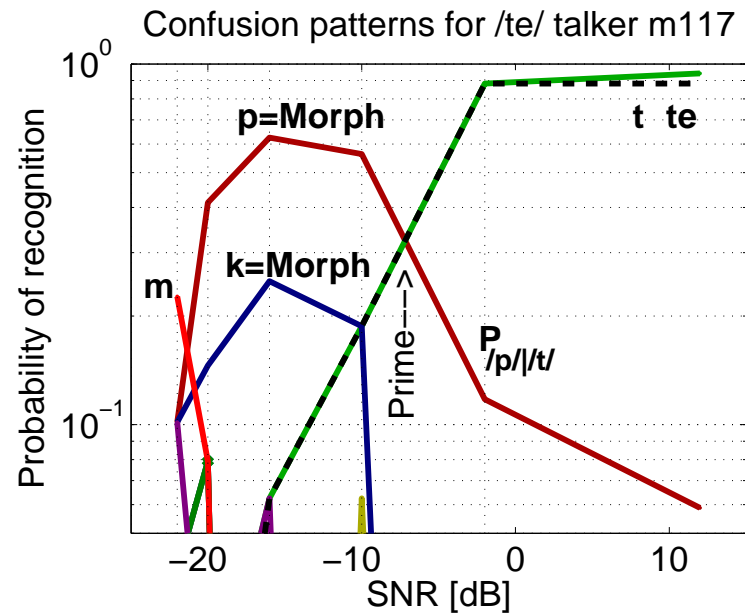
(b) Consonant errors HI

# Row of CM $P_{h|/t/}$

- Utterance phone scores are heterogeneous!



(c) Average over all /t/s.

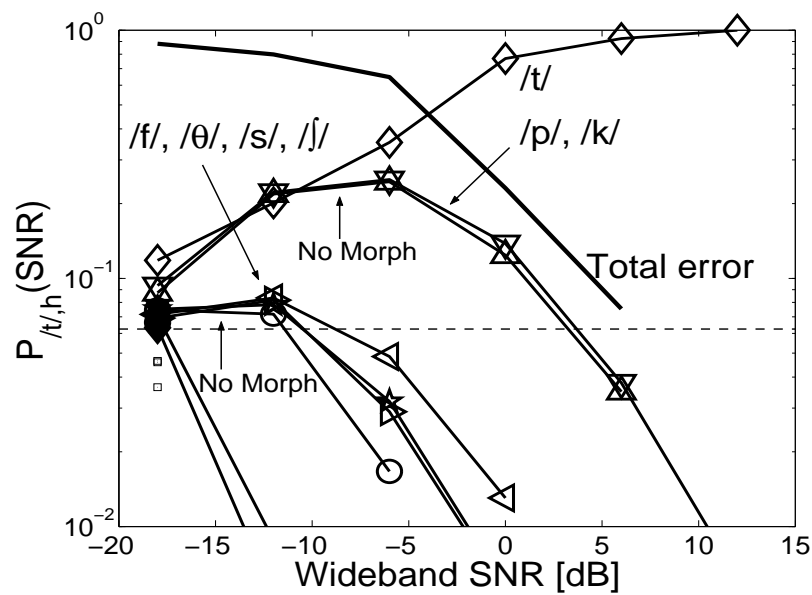


(d) Talker m117 /te/  $P_{h|/ta/}$ (SNR)

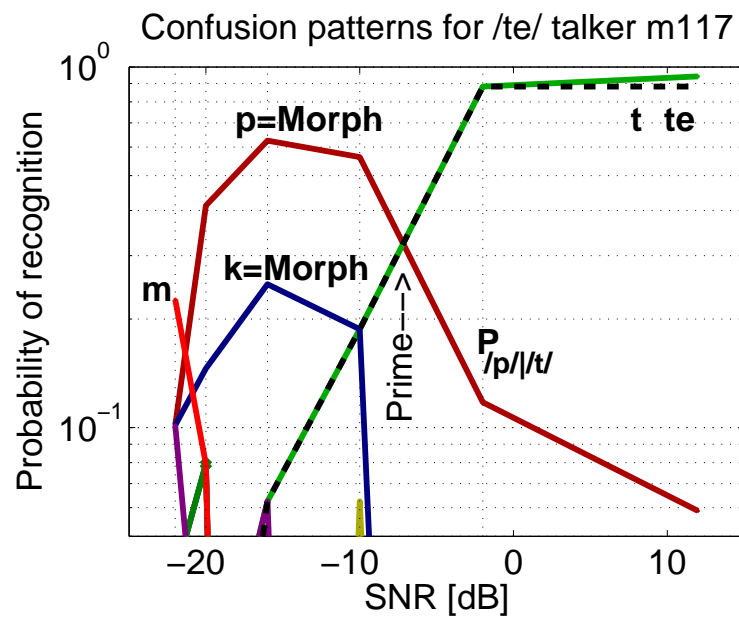
- Phone groups are due to shared sub-phonemic units
  - CV Morphs
  - Morphing sentences

# Row of CM $P_{h|/t/}$

- Utterance phone scores are heterogeneous!



(e) Average over all /t/s.

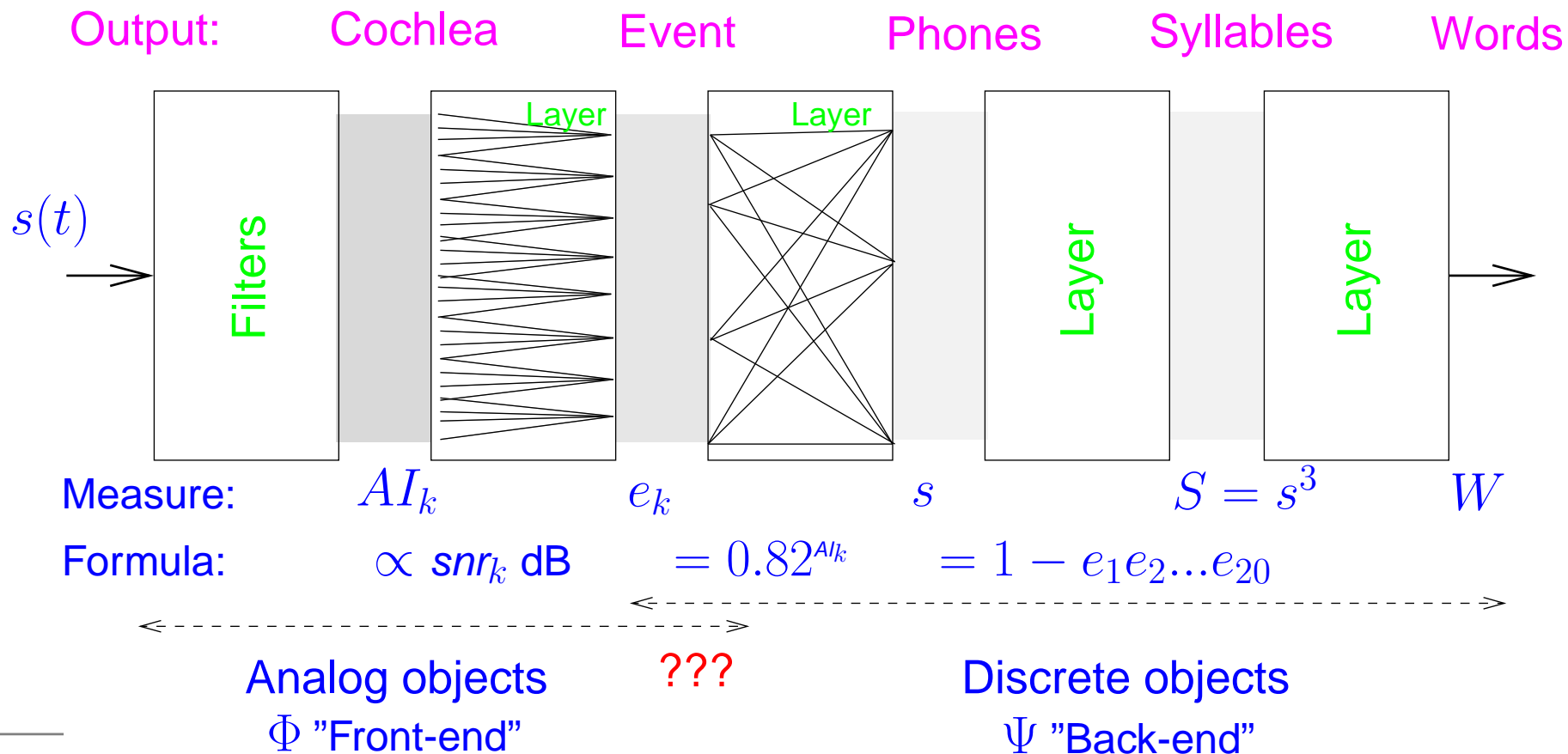


(f) Talker m117 /te/  $P_{h|/ta/}$  (SNR)

- Phone groups are due to shared sub-phonemic units
  - CV Morphs DEMO
  - Morphing sentences DEMO

# Model of human speech recognition HSR

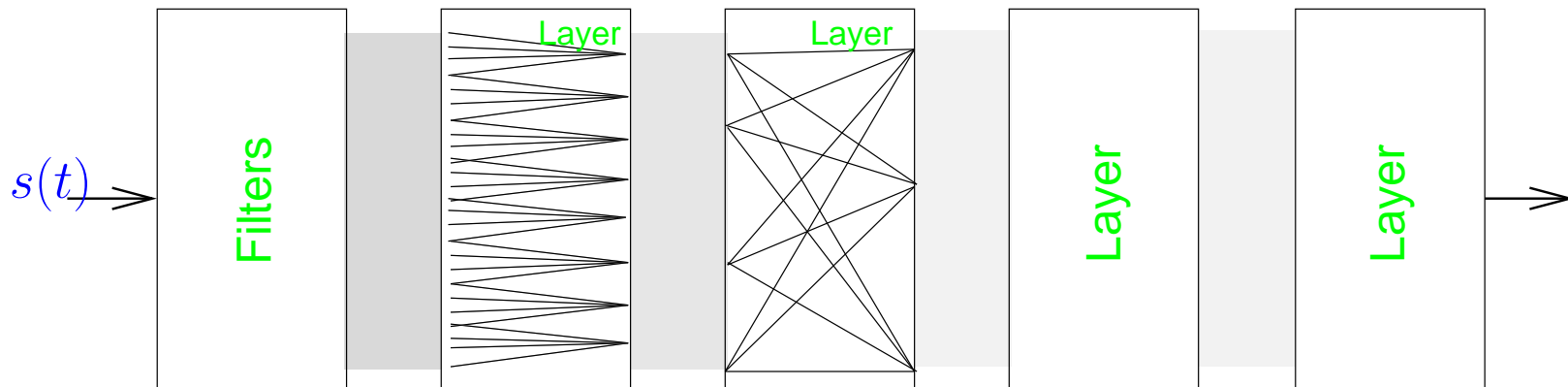
- Research Goal:
  - Identify *elemental HSR events*
  - An *event* is defined as a *perceptual feature*
  - *Event errors* are measured by band errors  $e_k$



# Definition and use of the $AI$

- The average error is:  $P_e(SNR) \equiv \prod_k e_k = 0.02^{AI}$
- $e_k = 0.822^{AI_k(snr_k)}$  cochlear  $k^{th}$  band-error
- $AI_k = \log_{10}(1 + 4snr_k^2)^{1/3}$  band channel-capacity
- $AI \equiv \overline{AI_k} = \frac{1}{20} \sum_{k=1}^{20} AI_k,$

Output:      Cochlea      Event      Phones      Syllables      Words



$AI_k \propto snr_k \text{ [dB]}$      $e_k = 0.82^{AI_k}$      $s = 1 - e_1 e_2 \dots e_{20}$      $S_{cv} = s^2$      $W$

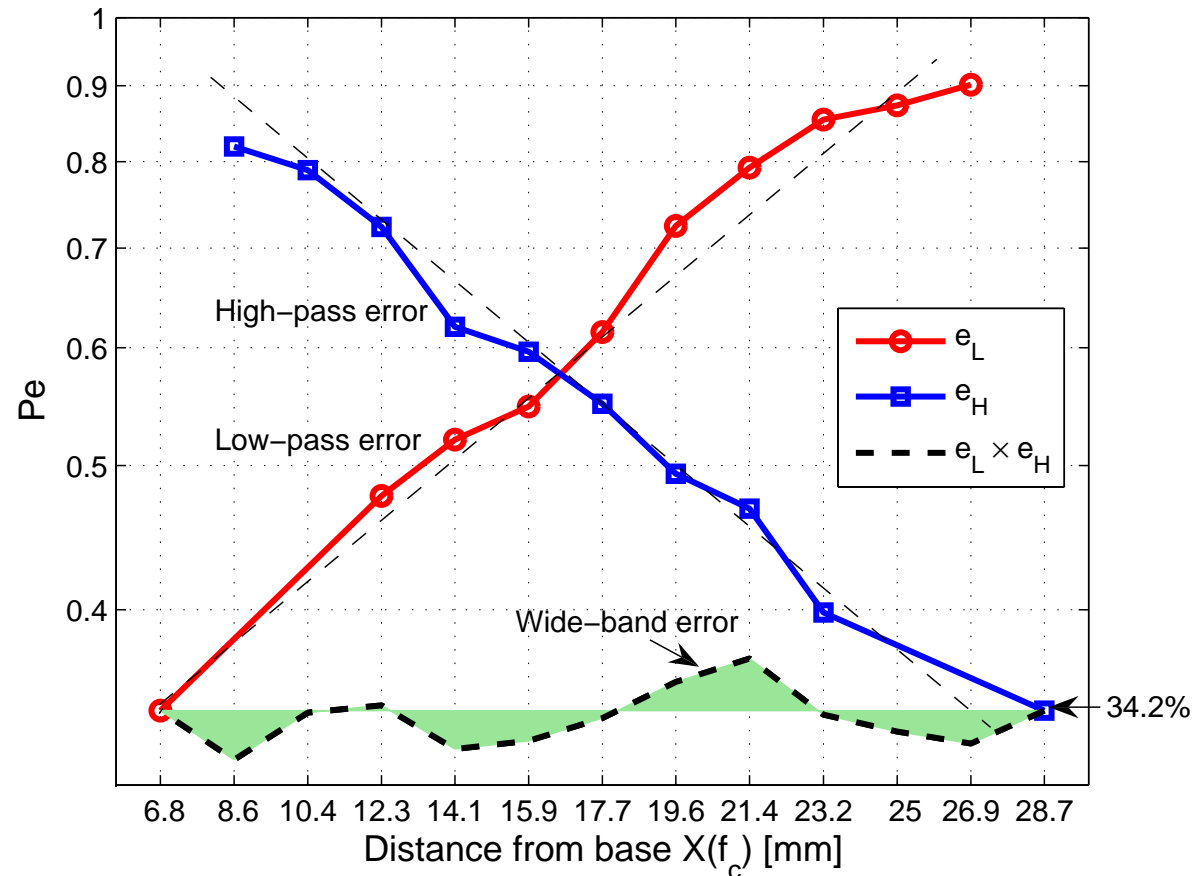
←----- Analog objects      ???      Discrete objects -----→



# Fletcher's Lowpass/Highpass result

- The AI is based on the *band-error product formula*

$$P_e(\text{snr}, f_c) \equiv e_{lp}(\text{snr}, f_c) \times e_{hp}(\text{snr}, f_c)$$

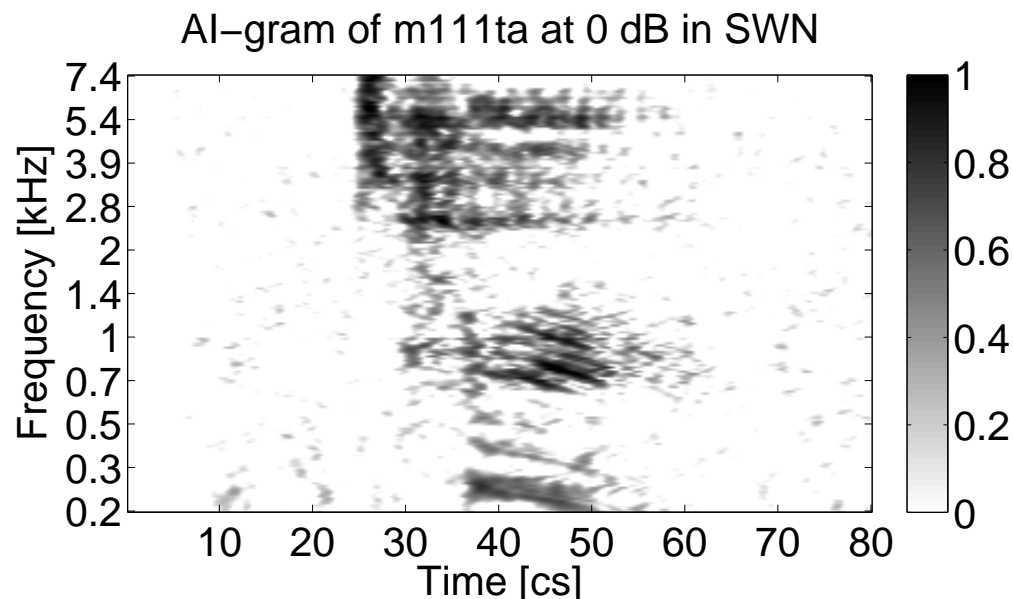


# Human listeners as a Shannon Channel

- The **Channel capacity theorem** gives the maximum information rate as:

$$C \equiv \int \log_2 (1 + \text{snr}^2(f)) df \quad (1)$$

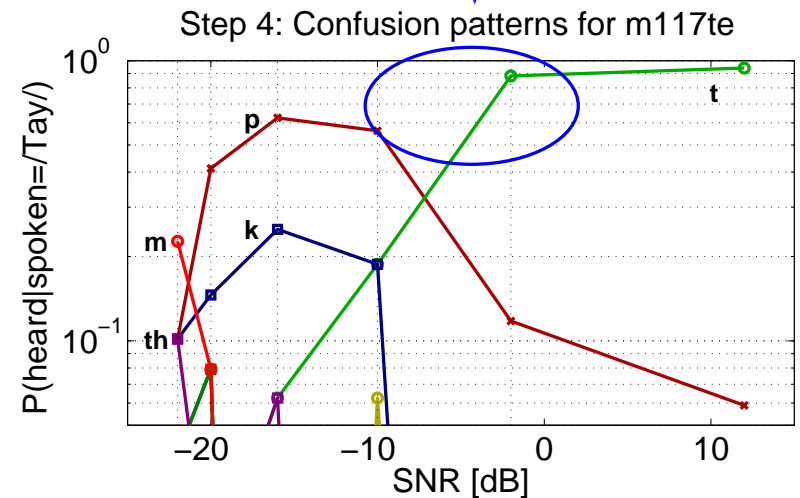
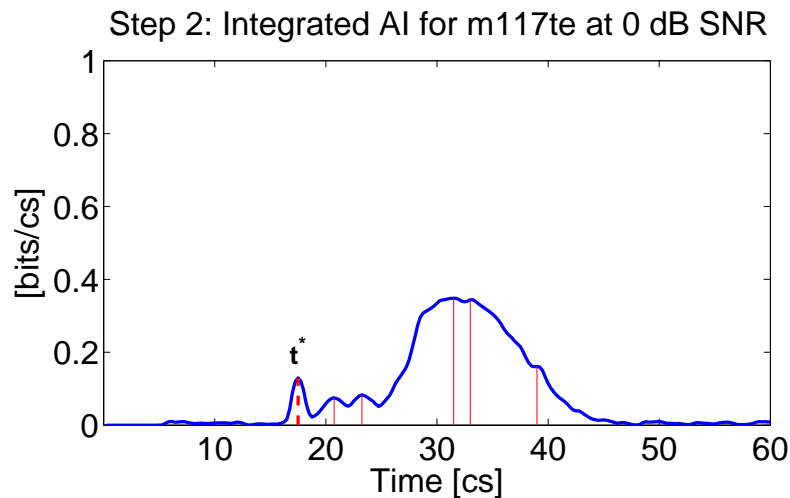
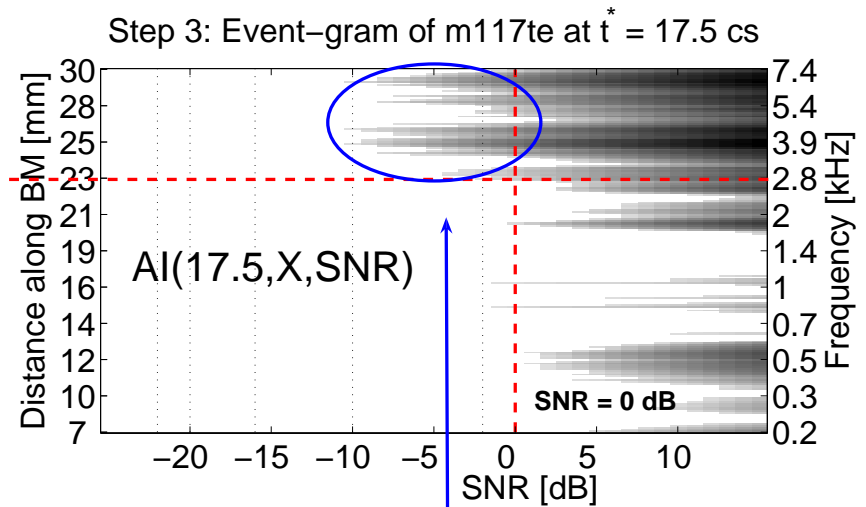
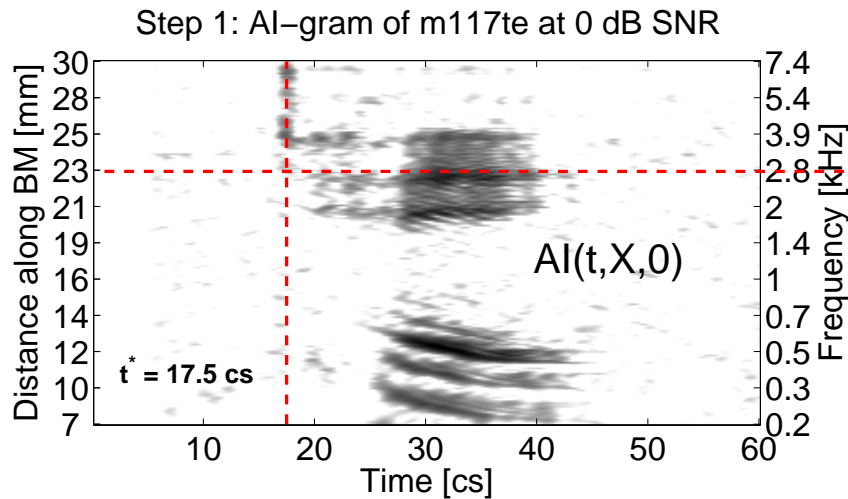
- For a Maximum Entropy (**MaxEnt**) speech source, the maximum information rate is determined by the SNR
- The **AI-gram** is a closely related measure:



# III–Results (30 mins)

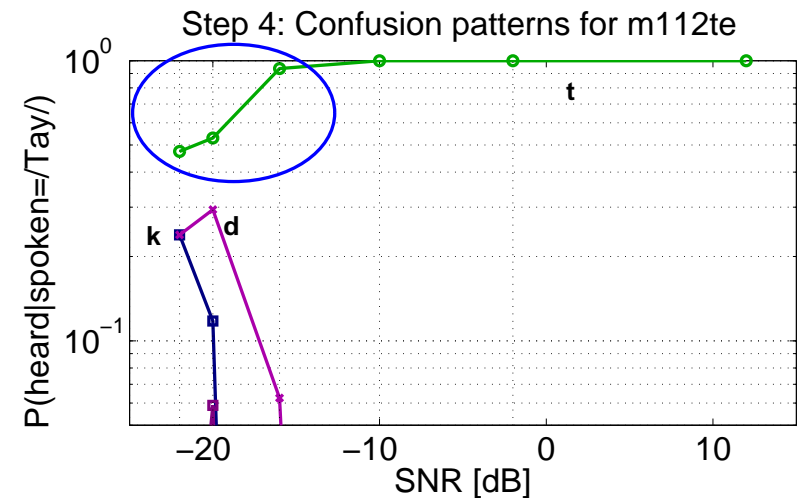
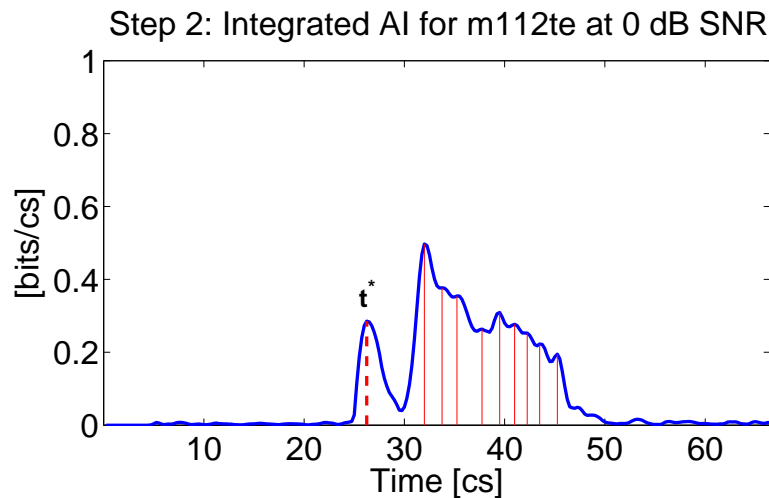
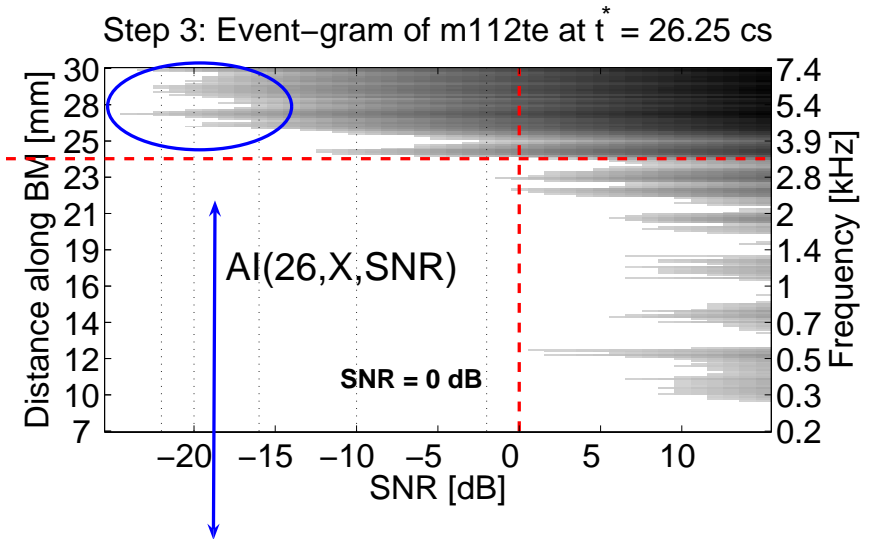
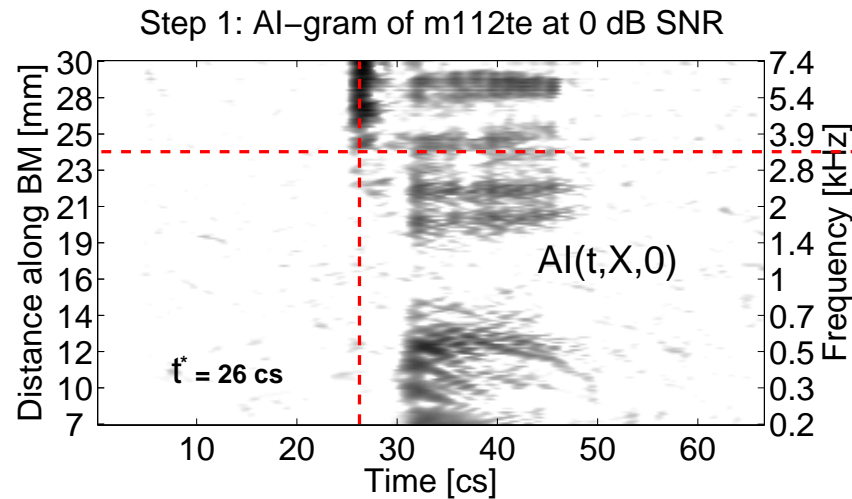
- Examples and Demos of events
  - Plosive CV events
  - Fricative CV events
- Conflicting cues
- **DEMOS:**
  - Event isolation
  - Consonant morphing
  - Consonant enhancement
  - Conflicting cues within consonants
  - Sentence meaning modification

# m117/tε/ in speech-weighted noise



- /t/ confusion threshold at  $P_c(SNR^* = -2) = 0.9$  correlated to Event-gram

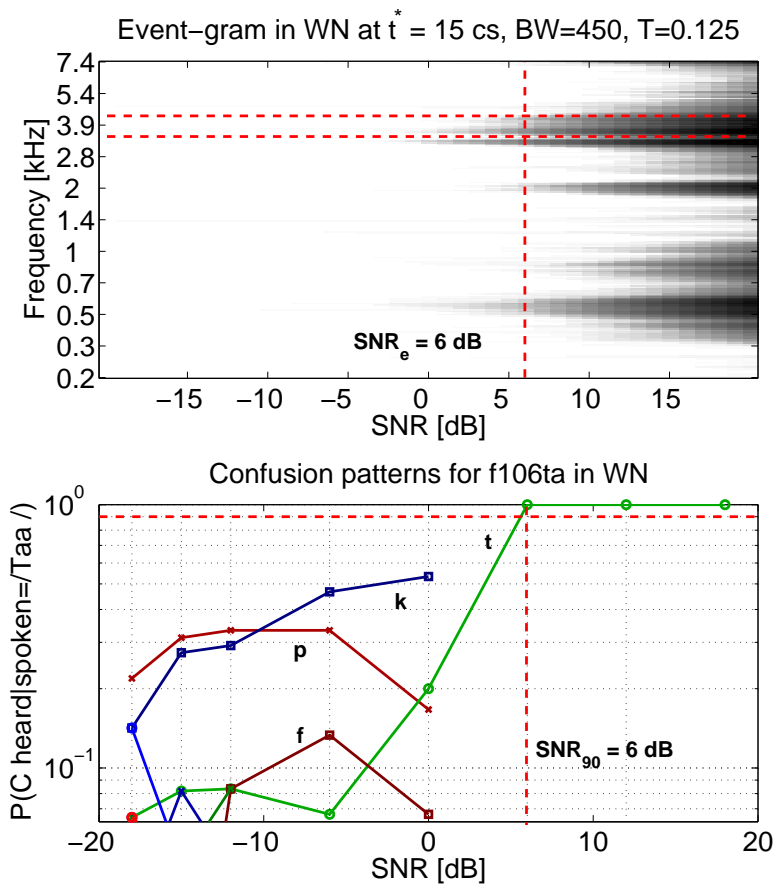
# m112/tε/ in speech-weighted noise



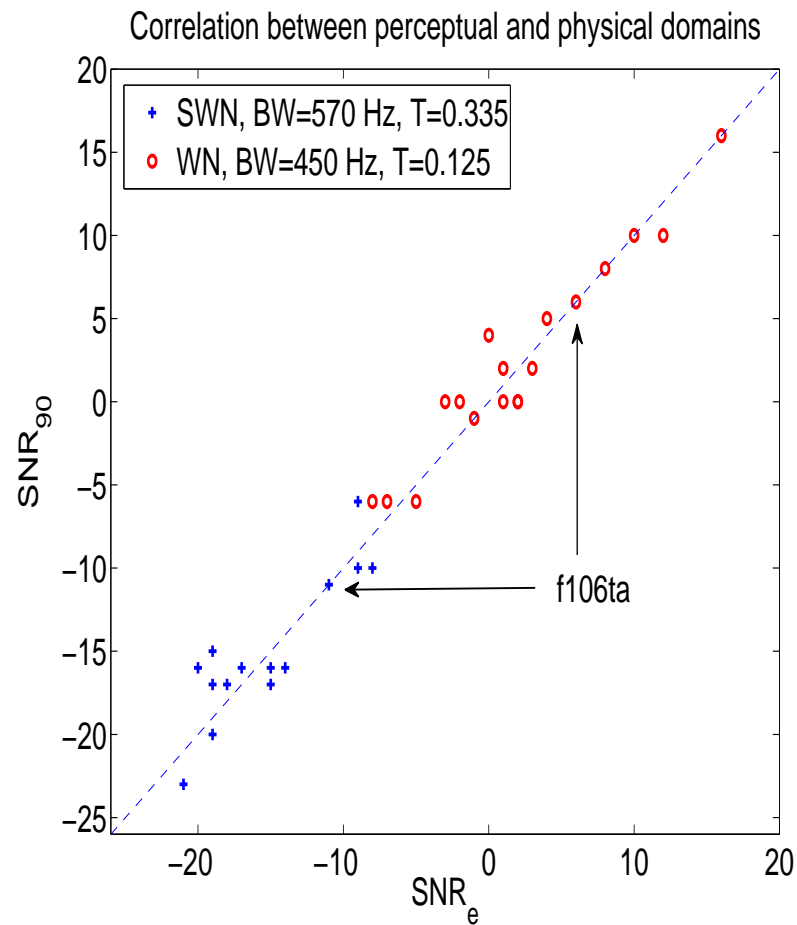
- /t/ confusion threshold at  $P_c(SNR^* = -16) = 0.9$  correlated to Event-gram

# Correlations of /t/ events

- High correlation across all /t/'s in the database

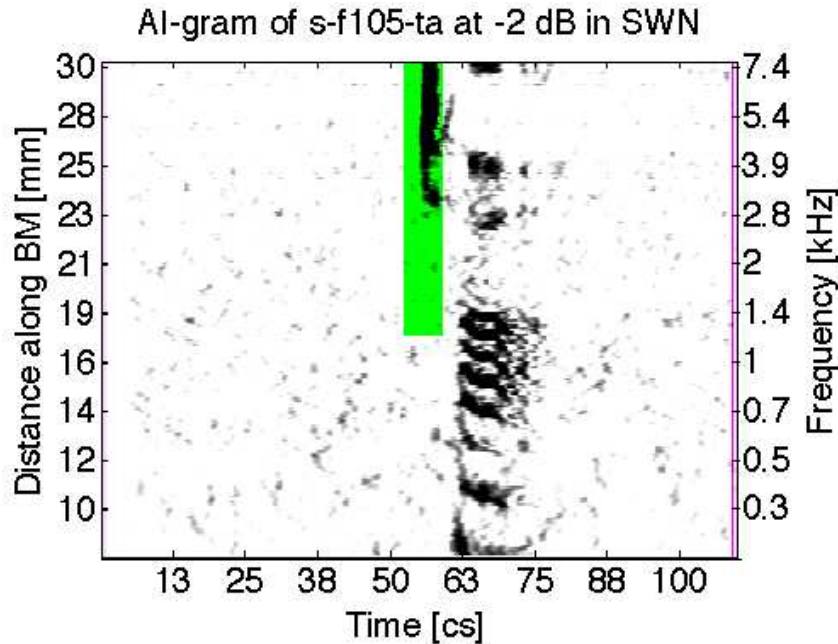


(g)

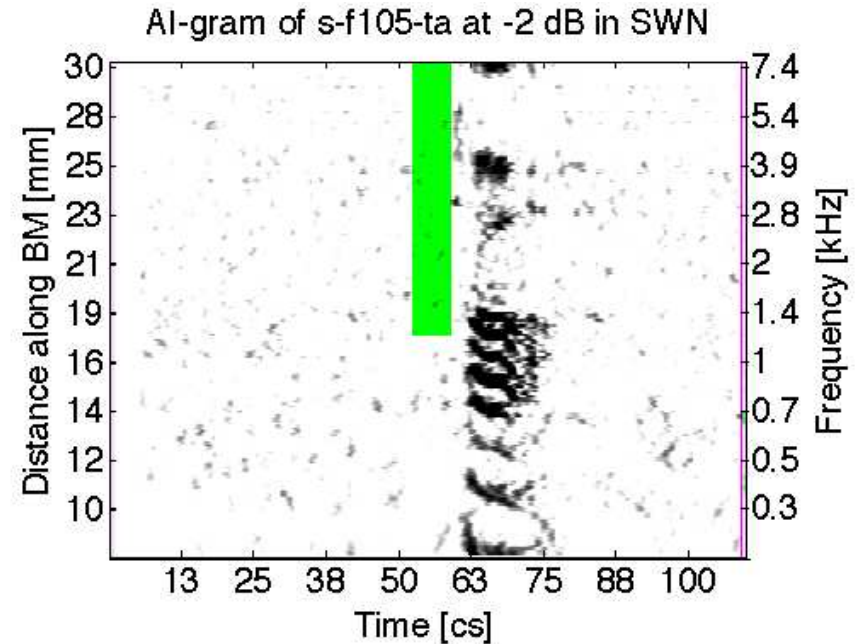


(h)

# Masking of /ta/ timing cue



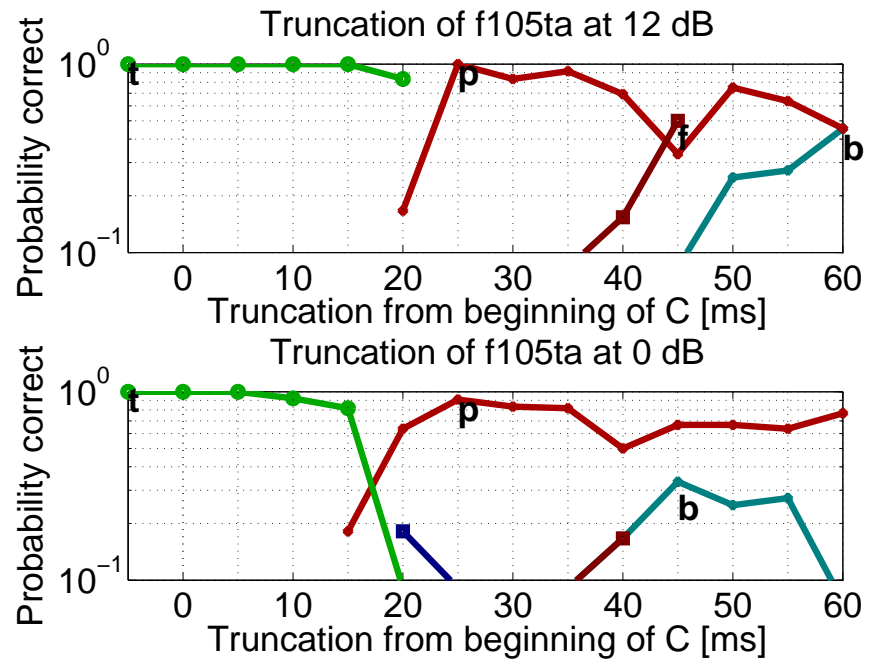
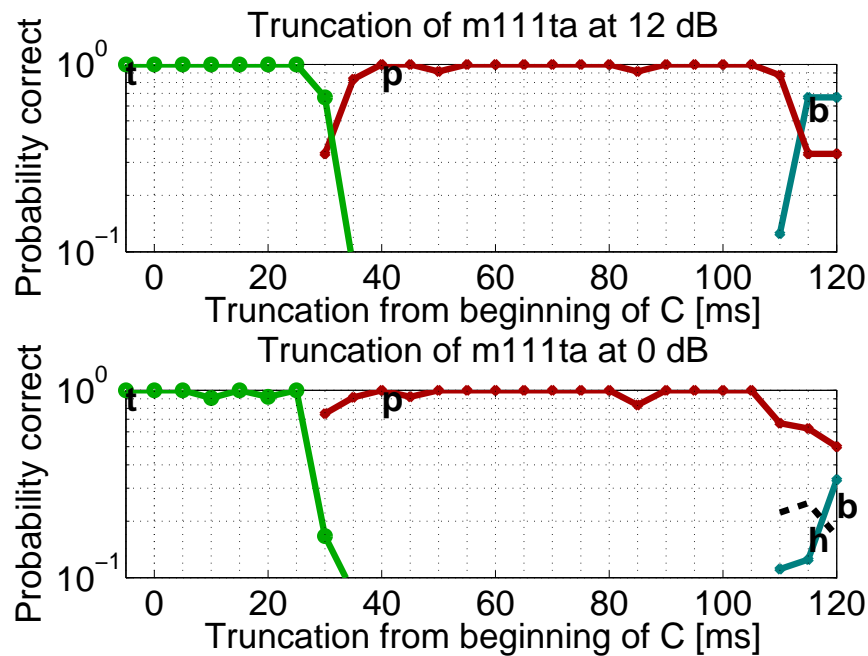
(i) Original /ta/



(j) Modified /ta/

- When the /t/ burst is masked by noise, the perception morphs to /p/
- DEMO 4

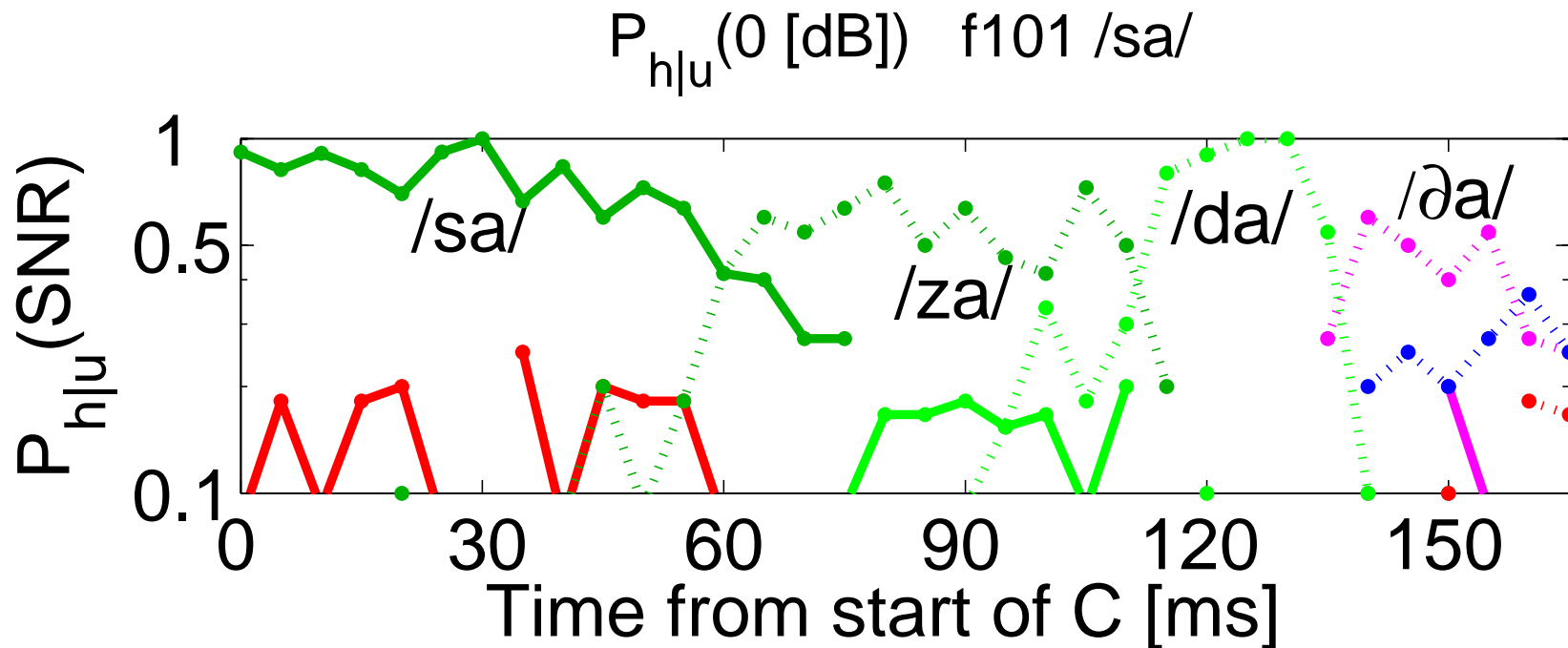
# Truncation of /ta/



- This represents the normal hearing responses to a truncated /ta/, from the start of the consonant
- Morphing from /ta/ to /pa/ to /ba/ at 0 and 12 dB SNR
- Similar to Furui 1986, and our extensive results



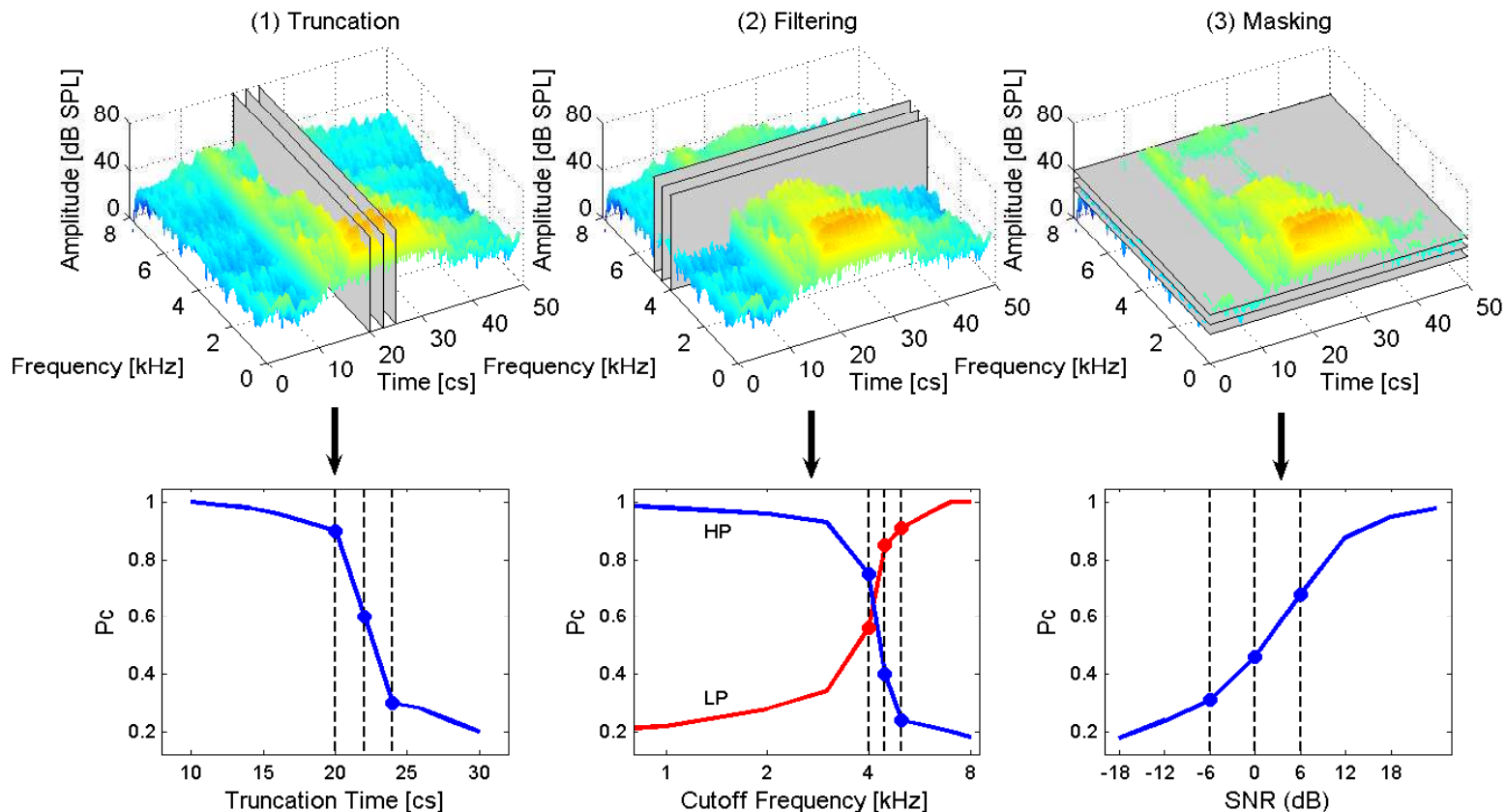
# Truncation of f101 /sa/



- This represents the normal hearing responses to a truncated /sa/, from the start of the consonant
- Morphing from /sa/ to /za/ to /da/ to /ðə/
- **Duration** seems to be a fricatives event

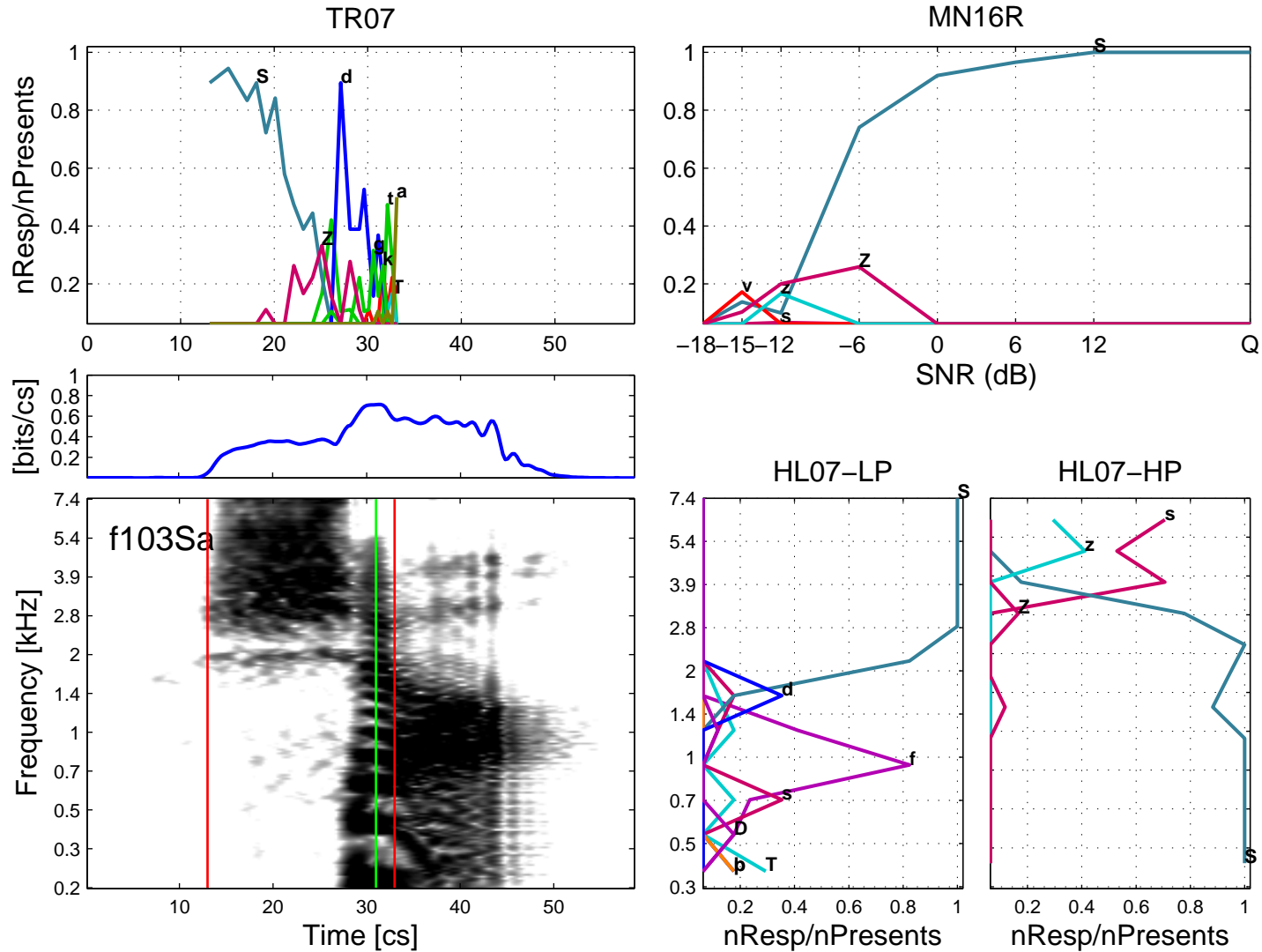
# Methods: $3^d$ Deep Search ( $3^d$ -DS)

- $3^d$  Deep-Search ( $3^d$ -DS) via truncation:
  - SNR truncation (i.e., masking)
  - Frequency truncation (High/Low-pass filtering)
  - Time truncation (Furui 1986)



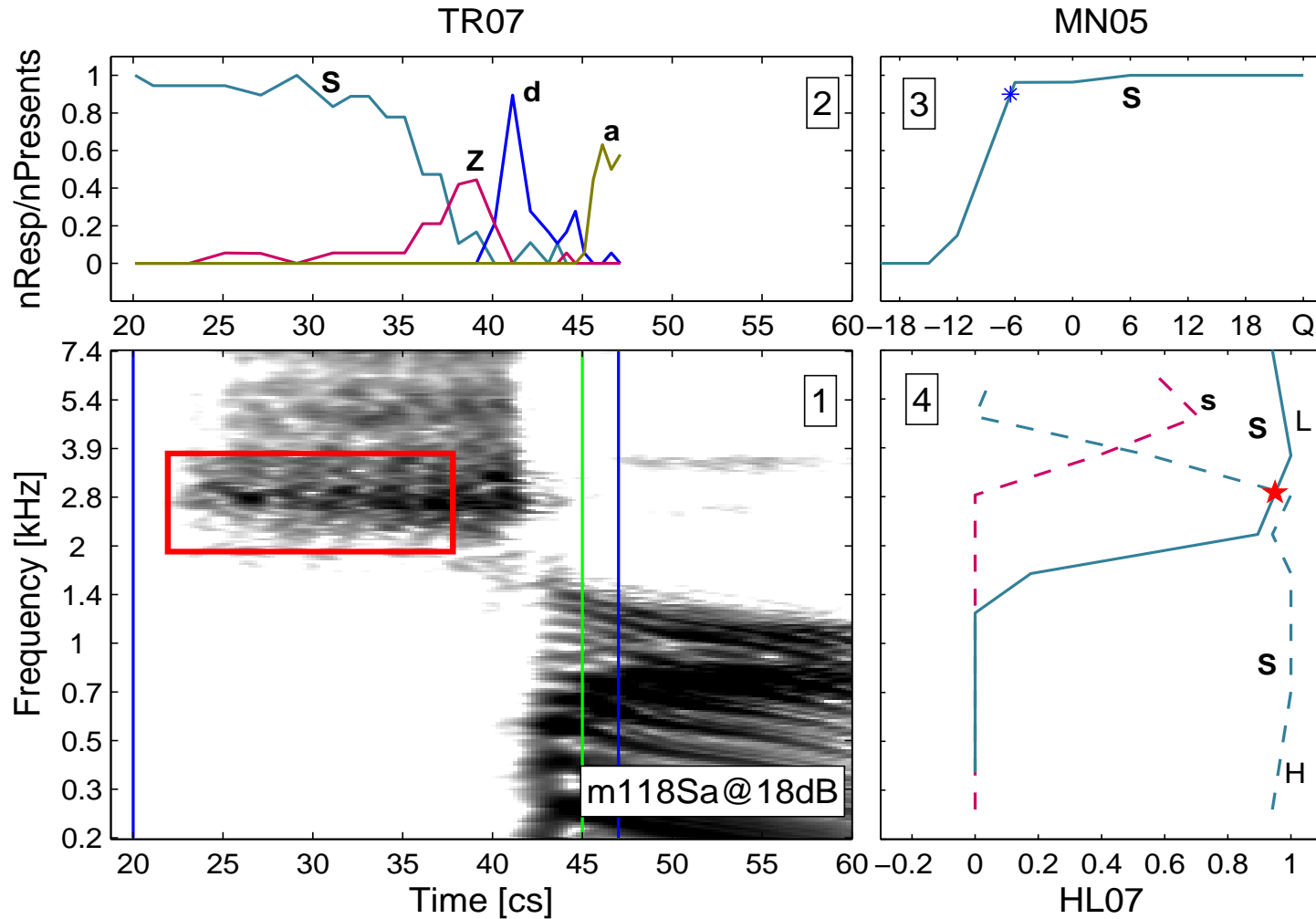
# 3<sup>d</sup>-DS Method /ja/

## ● Truncation in Time, Intensity and Frequency



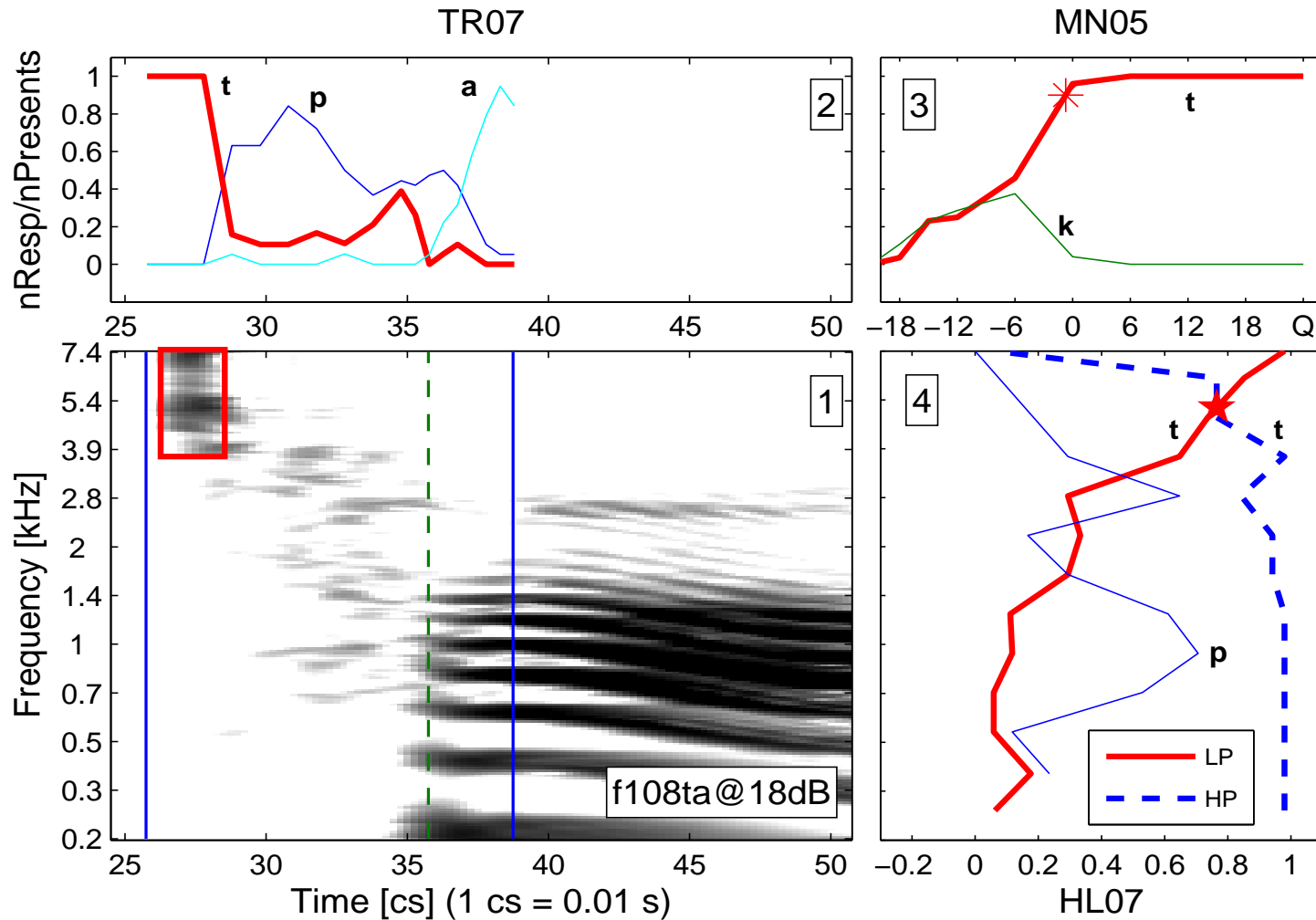
# 3<sup>d</sup>-DS Method /sa/

- Truncation in Intensity, time and frequency

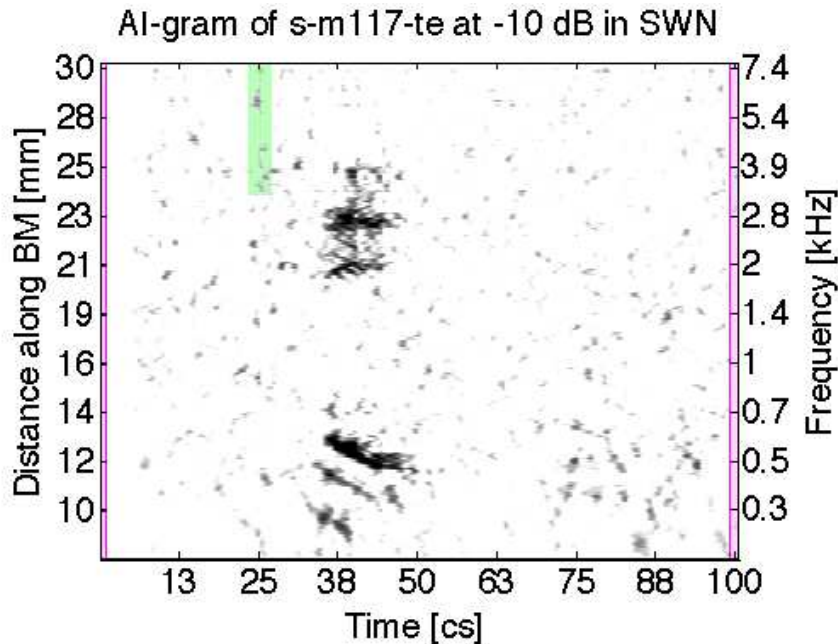


# 3<sup>d</sup>-DS Method /ta/

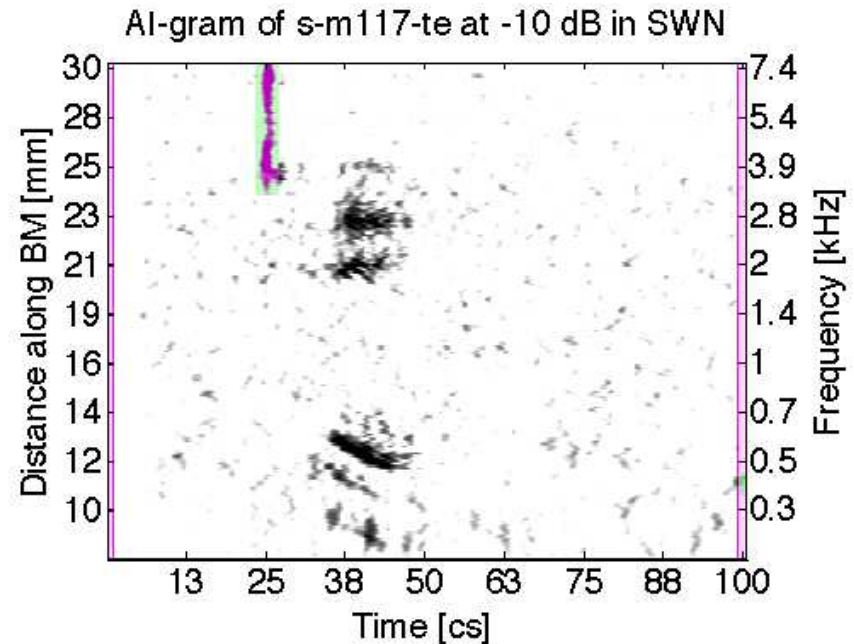
- Truncation in Intensity, time and frequency



# Enhancement of /tɛ/ event



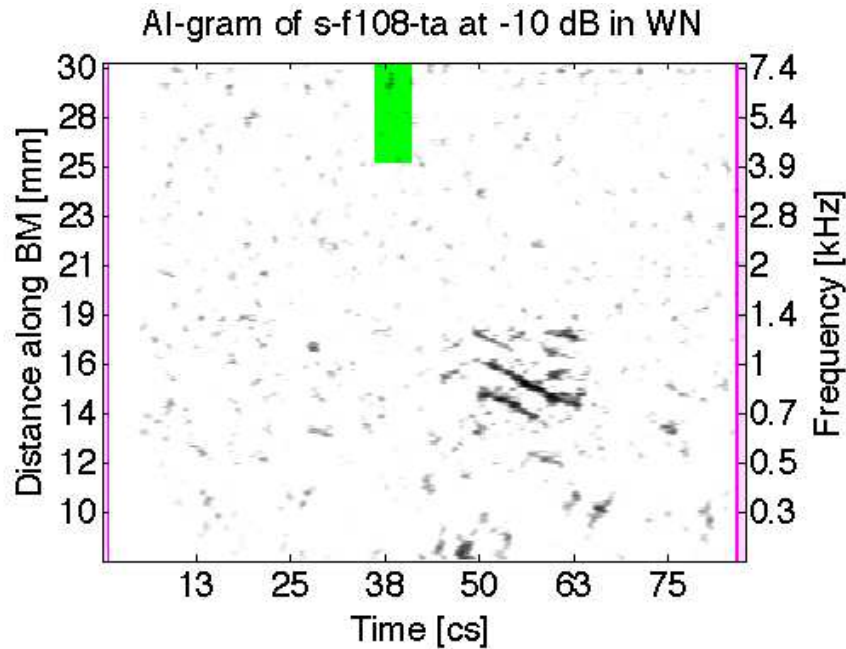
(k) Original /tɛ/



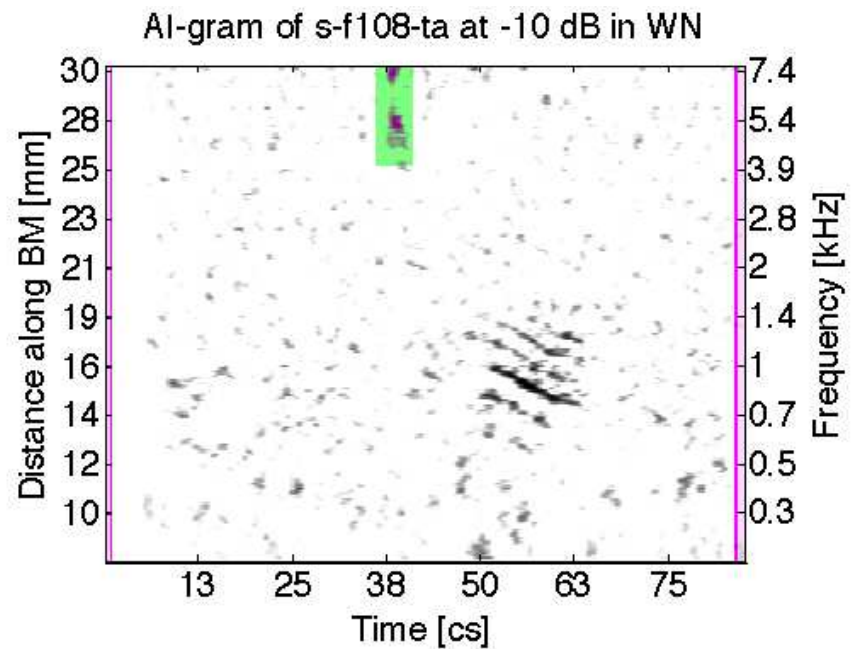
(l) Modified /tɛ/

- The sound is heard as /t/ again, we suppressed the morph (see confusion patterns of slide 4)
- METHODS: The /t/ burst is enhanced (14 dB) on the quiet sound, then noise is added
- DEMO

# Enhancement of /ta/ event



(m) Original /ta/



(n) Modified /ta/

- The sound is heard as /t/ again, we increase /t/ recognition
- METHODS: The /t/ burst is enhanced (14 dB) on the quiet sound, then noise is added
- DEMO

# Conclusion I

We have:

- 1 isolated events for CV: **Plosives** /p, t, k/ and /b, d, g/ and **Fricatives** /θ, ʃ, tʃ, s, h, f/ and /z, ʒ, v, ð/) + **Vowels** /o, ε, ɪ/
  - for many individual talkers
  - via new tools (AI-gram, Event-gram and 3<sup>d</sup>-DS)



# Conclusion I

We have:

- 1 isolated events for CV: **Plosives** /p, t, k/ and /b, d, g/ and **Fricatives** /θ, ʃ, tʃ, s, h, f/ and /z, ʒ, v, ð/) + **Vowels** /o, ε, ɪ/
  - for many individual talkers
  - via new tools (AI-gram, Event-gram and 3<sup>d</sup>-DS)
- 2 shown that normal listeners use
  - *across-frequency timing coincidences*
  - duration and bandwidthto discriminate consonants in noise

# Conclusion I

We have:

- 1 isolated events for CV: **Plosives** /p, t, k/ and /b, d, g/ and **Fricatives** /θ, ʃ, tʃ, s, h, f/ and /z, ʒ, v, ð/) + **Vowels** /o, ε, ɪ/
  - for many individual talkers
  - via new tools (AI-gram, Event-gram and 3<sup>d</sup>-DS)
- 2 shown that normal listeners use
  - *across-frequency timing coincidences*
  - duration and bandwidthto discriminate consonants in noise
- 3 developed tools to
  - Morphed speech sounds
  - Decrease or increase intelligibility. Ex: /tɑ/, /tɛ/

# Conclusion II

We have shown:

- 1 the existence of conflicting cues
  - Thus MaxEnt consonants are NOT redundant

# Conclusion II

We have shown:

- 1 the existence of conflicting cues
  - Thus MaxEnt consonants are NOT redundant
- 2 that the event threshold is abrupt (i.e., 6 dB)

# Conclusion II

We have shown:

- 1 the existence of conflicting cues
  - Thus MaxEnt consonants are NOT redundant
- 2 that the event threshold is abrupt (i.e., 6 dB)
- 3 proven the AI band-product formula (yet again)

# Conclusion II

We have shown:

- 1 the existence of conflicting cues
  - Thus MaxEnt consonants are NOT redundant
- 2 that the event threshold is abrupt (i.e., 6 dB)
- 3 proven the AI band-product formula (yet again)
- 4 why the AI works
  - Due to the frequency and SNR event distribution

# Conclusion II

We have shown:

- 1 the existence of conflicting cues
  - Thus MaxEnt consonants are NOT redundant
- 2 that the event threshold is abrupt (i.e., 6 dB)
- 3 proven the AI band-product formula (yet again)
- 4 why the AI works
  - Due to the frequency and SNR event distribution
- 5 the role of **forward** and **upward** masking spread

# Conclusion III

This could lead to:

- 1 Improved automatic speech recognition front-ends



# Conclusion III

This could lead to:

- 1 Improved automatic speech recognition front-ends
- 2 The design of new hearing aids

**Question your basic assumptions**

**Thanks for your attention**  
**<http://hear.ai.uiuc.edu>**