# The role of the cochlea in Human speech recognition

Where is the speech information lurking?

Jont Allen

Sandeep Phatak

Marion Regnier

Feipeng Li

Univ. of IL, Beckman Inst., Urbana IL







# Objective



#### ACOUSTIC FEATURES

EVENTS

- Develop rigorous procedures for analyzing and modifying speech in noise, to:
  - identify perceptual features, denoted events
  - Develop a theory of human speech recognition (HSR) based on two basic measures:
    - 1. Al-Gram (speech audibility measure)
    - 2. Confusion matrix (speech discrimination measure)
- Show that across-frequency timing cues are events

## Human listeners as a Shannon Channel

- My approach is inspired by information theory using a classic 3-pronged approach: Simplify, simplify, simplify
  - 1. The Channel capacity theorem gives the maximum information rate as:

$$\mathcal{C} \equiv \int \log_2 \left( 1 + \operatorname{snr}^2(f) \right) df \tag{1}$$

- The basic idea is to use a Maximum entropy (MaxEnt) speech source, and reduce the maximum information rate for by increasing the noise.
- Take full advantage of Articulation Index predictions of the average phone score  $s = P_c(AI)$

# **Model of human speech recognition (HSR**

- The research goal is to identify elemental HSR events
  - An event is defined as a perceptual feature
  - Event errors are measured by band errors  $e_k$



# Articulation Matrices and elemental event

Miller-Nicely's 1955 articulation matrix A measured at [-18, -12, -6 shown, 0, 6, 12] dB SNR

k ð Ъ đ Þ g  $\boldsymbol{z}$ m n р l STIMULUS 3 θ S S Ь d U δ m n VOICED NASAL UNVOICED RESPONSE

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

Confusion groups imply underlying elemental events

# Case of /pa/, /ta/, /ka/ with /ta/ spoken

Phone groups imply sub-phonemic units (i.e., events)



- How many events, and of what form?
- Plot of  $A_{i,j}(snr)$  for row i = 2

• Solid red curve is total error  $e_2 \equiv 1 - A_{2,2} = \sum_{j \neq 2} A_{2,j}$ 

# The case of /ma/ vs. /na/



- This 2-group of sounds is closed since  $\mathcal{A}_{/ma/,/ma/}(SNR) + \mathcal{A}_{/ma/,/na/}(SNR) \approx 1$ 
  - There can be only 1 event
  - Solid red curve is the total error:  $e_i \equiv 1 - A_{i,i} = \sum_{j \neq i} A_{i,j}(SNR)$

## Fletcher's Lopass/Hipass result

#### The AI is based on the band error product formula

 $1 - P_c(SNR) \equiv e_{total}(SNR, f_c) = e_{lp}(SNR, f_c)e_{hp}(SNR, f_c)$ (2)



# **Probabilistic measures of recognition**

- $k^{th}$  band articulation index:  $AI_k = \frac{10}{30}\log(1 + c^2 snr_k^2)$ 
  - c = 2,  $k = 1 \cdots K$  with K = 20
- **• Band (event) error:**  $e_k = e_{\min}^{Al_k/K} = \sqrt[20]{0.02}^{Al_k} = 0.822^{Al_k}$
- The  $AI \equiv \frac{1}{K} \sum_k AI_k$ , ?
- MaxEnt phone score:  $s = 1 e_1 e_2 \dots e_K = 1 e_{\min}^{AI}$



#### How can we find events?

- A 4-Step analysis relates confusions to an audibility measure (???):
- Modification of speech sounds
  - We developed a tool based on the Short-Time Fourier Transform (STFT) (?) that allows us to selectively:
  - Mask with noise specific time and frequency regions
    - so that this specific part of the speech becomes inaudible
    - selectively amplify specific regions
    - to increase intelligibility
    - We will present audio examples of original and modified sounds

# m117/te/ in speech-weighted noise



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

# m112/te/ 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



# Masking of /ta/ timing cue



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 previous studies ?, and our more 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 /ða/
- Duration seems to be a fricatives event

# /ma/- /na/ discrimination

- Ind/recognition from /md/ relies on a  $\approx 50$  ms delay formed from the F<sub>1</sub> and F<sub>2</sub> collision
- When we edit the speech so that the onset is simultaneous above 0.6 kHz, the /na/ is robustly and naturally heard as /ma/
- METHODS: 9 listeners evaluated these sounds in open response random trial experiment.

# **Deletion of /na/ timing cue**



# **Creation of /ng/ timing cue**



## **Enhancement of /tɛ/ event**



- 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



## **Enhancement of /ta/ event**



- 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



## Conclusion

- We have shown that normal listeners use across-frequency timing coincidences to discriminate consonants in noise
- We have developed a tool to modify speech sounds
  - Morph sounds. Ex: /ma/ /na/
  - Decrease or increase intelligibility. Ex: /tα/, /tε/
- This could well lead to the design of new hearing aids