### Auditory Filters: recent work connecting psychophysics, physiology, mechanics, and neuromorphic implementations

Richard F. Lyon Google, Inc. July, 2008 Telluride Neuromorphic Cognition Workshop

### My neuro history

(see http://www.dicklyon.com)

- 1978 "A Signal Processing Model of Hearing" paper for Stanford class on neural information processing
- 1981 Xerox optical mouse incorporates explicit model of lateral inhibition
- 1986–94 "Analog Electronic Cochlea" with Mead and his crew
- 1993–97 Neural-net–based handwriting recognition
- 1997–2006 Foveon color image sensors and cameras
- 2006–2008 Google "Machine-Hearing" project

#### What is the function of a rainbow?



### Pole–Zero Filter Cascade (PZFC) – preview

- Good fit to human masking data with simple parameters
- As with APFC, connection to traveling wave allows natural coupling effects, for masking, adaptation, etc.
- Like APGF & OZGF, unity-gain tail models lossless propagation of low-frequency energy; tail doesn't wag with Q or other parameters
- Easy to implement directly as standard second-order (pole-zero) filter sections, without further approximation
- Easy to vary parameters dynamically for "AGÇ" here response
- Low total order (complexity) for multi-channel filterbank, by sharing filter sections – total complexity just 2nd-order per channel, compared to 8th-order for gammatones



"Audito shapes a	ory filter" nd models	(86	0 -20 - 30 dB 70	P1312x2x dB
Resonance $N=1$ Schafer et al. 50 Patterson 70 Gaussian $N=\infty$	Roy Patterson Auditory Filter and their Pole O RoEx Patterson and Nimmo-Sm	Shapes rders (N)	-40 -60 -80 500 1000 1500 fr	2000 2500 3000 3500 equency (Hz)
Swets et al. 62 Patterson 76	RoEx(p), RoEx(p,r), RoE Patterson et al. 82	x(p,t,w)		APGF, OZGF N=8, 16, 32 Lyon 96
Symmetric approx N=2 Patterson 74	K. GTF	GTF N=4 Patterson 88		GCF N=6 Irino 95
GTF, OZGF, ++ N=3 Flanagan 60	Two-side RoEx Moore and Glas	sberg 87	Two-sided GTF Rosen and Baker	94
Except for	Flanagan 60, only psychophys	sical models ar	e included	
	It would be goo	od to conr	nect to phys	iology

#### Don't wag the tail when changing pole damping – Like APGF & OZGF, unlike Gammatone & Gammachirp



We can now plot the location of these poles and zeros assuming a gammatone filter with a bandwidth of 125Hz and a center frequency of 1000Hz (this approximates a 1000Hz cochlear channel.) In the graph below, the four zeros are shown along the real (horizontal) axis. There are four sets of poles at each of the dots indicated near  $\pm 6000j$  (which indicates a resonance near 1000Hz.)





# Jim Flanagan 1960–62 filter models of basilar membrane



a reasonable fit to Békésy's results is

 $F_{l}(s) = c_{1}\beta_{l}^{4} \left(\frac{2000 \pi \beta_{l}}{\beta_{l} + 2000 \pi}\right)^{0.8} \left(\frac{s + \varepsilon_{l}}{s + \beta_{l}}\right) \left[\frac{1}{(s + \alpha_{l})^{2} + \beta_{l}^{2}}\right]^{2} e^{\frac{-3\pi s}{4\beta_{l}}}$ 

The membrane response at any point is therefore approximated in terms of the poles and zeros of the rational function part of  $F_l(s)$ . As



Fig. 4.17 a. Pole-zero diagram for the approximating function  $F_l(s)$  (after FLANAGAN, 1962a) Fig. 4.17 b. Amplitude and phase response of the basilar membrane model  $F_{I}(s)$ . Frequency is normalized in terms of the characteristic frequency  $\beta_{I}$ 

order-3 "gamma-tone" filter

126

Techniques for Speech Analysis

good approximation to the displacement impulse response of the basilar membrane, at a point maximally responsive to radian frequency  $\beta$ , is

$$p(t) = (\beta t)^2 e^{-\beta t/2} \sin \beta t$$

$$= h_{bm}(t) \sin \beta t.$$
(5.9)

The time window for the basilar membrane, according to this modeling<sup>1</sup>, is the "surge" function plotted in Fig. 5.6. One notices that the time window has a duration inversely related to  $\beta$ . It has its maximum at  $t_{\max} = 4/\beta$ . If, as a crude estimate,  $2t_{\max}$  is taken as the effective dura-



tion D of the window, then for several membrane places:

β/2 π (cps)	$D = 2t_{max}$ (msec)
100	12.0
1000	1.2
5000	0.2

Fig. 5.6. The effective time window for short-time frequency analysis by the basilar membrane in the human ear. The weighting function is deduced from the ear model discussed in Chapter IV

For most speech signals, therefore, the mechanical analysis of the ear apparently provides better temporal resolu-

tion than spectral resolution. Generally, the only harmonic component resolved mechanically is the fundamental frequency of voiced segments. This result is borne out by observations on the models described in Chapter IV.

#### Computer Simulation of Membrane Motion

101

suitable for calculations in a digital computer. One such digital simulation represents the membrane motion at 40 points (FLANAGAN, 1962b).

As might be done in realizing the analog electrical circuit, the digital representation of the model can be constructed from sampled-implementations: data equivalents of the individual complex pole-pairs and the individual real poles and zeros. The sampled-data equivalents approximate the continuous functions over the frequency range of interest. The computer operations used to simulate the necessary poles and zeros are shown

 $Y_a(s)$ 

in Fig. 4.25. All of the square boxes labelled D are delays equal to the time between successive digital samples. The input sampling frequency, 1/D, in the present



2e<sup>-8</sup>cos¢

Xa(S) 0



Fig. 4.26. Functional block diagram for a digital computer simulation of basilar membrane displacement



Such filters have easy Flanagan 1962

# Waves in uniform media: sinusoidal functions of *x* and *t*

In a uniform medium, a wave propagating toward increasing values of the place dimension is given by

Complex wavenumber  $k(\omega)$  controls loss or gain

 $W(x) = A \exp(i\omega t - ikx)$ 

By examining the ratio of waves at two places separated by a distance  $\Delta x$ , we see that the wave at the farther place is equal to the wave at the nearer place multiplied by  $\exp(-ik\Delta x)$ , representing a frequency-dependent filter characterizing the stretch of length  $\Delta x$ .

- But in a non-uniform medium, the wavenumber k depends on place (x) as well as on frequency: k(ω, x)
- Net result is a **cascade of filters**  $exp(-ik(\omega, x)dx)$ .

# CMOS VLSI Cochlea: an all-pole filter cascade (APFC) – Lyon & Mead 1988

٧2 TAU TAIL X Figure 1: Second-order filter-section circuit. 2-pole stage is an order-1 GTF or APGF X



Figure 2: Floorplan of 100-stage cochlea chip.



## Some possible cascade stage frequency responses





### PZFC poles and zeros in the s plane



### PZFC stage frequency response: **pole** pair makes a bump (variable via pole Q),and **zero** pair makes a dip



### AGC (AQC) via feedback (automatic gain or Q control)

APGF (or OZGF) in feedback configuration



A Filter Cascade in feedback configuration uses all outputs to affect parameters of all stages, through a sort of diffusion spreading effect









## Response to sounds – no further compressive (log) nonlinearity needed



### PZFC Advantages – review

- Good fit to human masking data with simple params
- Connection to traveling wave allows natural coupling effects, for masking, adaptation, etc.
- Like APGF & OZGF, unity-gain tail models lossless propagation of low-frequency energy; tail doesn't wag with Q or other parameters
- Easy to implement directly as standard second-order filter sections, without further approximation
- Easy to vary parameters dynamically for "AGC"
- Low total order (complexity) for multi-channel filterbank, by sharing filter sections – total complexity just 2ndorder per channel, compared to 8th-order for gammatones

### Conclusion –

### continuous improvement path

- RoEx family good parameterized shapes, but not the best; no corresponding real filters
- Gammatone too symmetric, but otherwise a good filter, not hard to implement; tail problems in real case
- Gammachirp parameterized asymmetry, dynamically-varying peak gain; good improvement over the others, but the real filter still has tail issues
- AP/OZ GF a chirping asymmetric filter much like GC, but with rock-steady tail behavior as gain varies; exact easy implementation, even dynamic
- PZFC tied to traveling-wave concept; more efficient filterbank; shape fit at least as good as any others, with few parameters; can work on phase fits, too

# Fitting Nonlinear Auditory Filters: OZGF, PZFC, feedback versions, etc.

(Deriving the parameters for the humancalibrated versions of the OZGF and the PZFC from simultaneous masking data)

### Irino & Unoki's Framework

- Nonlinear optimization of parameters based on minimizing squared error in predicting masked thresholds based on tone SNR at filter output
- Also optimize over filter CF to get best SNR of the masked tone
- Level-dependent parameters depend on output of a "passive" filter with noise only (or with noise plus target tone)
- Nonlinear fit search also optimizes P0 and K
- Flexible frequency dependence of selected parameters: linear or quadratic on ERBrate scale

Level control via detection of output of "passive" filter

"passive" filter's output level controls the "active" filter's shape or gain, like this for parallel filters (PrIGC)



in "cascade" version (CasGC), the "passive" filter comes first, followed by the level-dependent part



Level control via detection at filter's own output

in "feedback" configuration, the main filter's output is used for level detection and control



Somewhat trickier, since you need a guaranteed stable way to let it settle to a consistent level

### **Framework Modifications**

- Better convergence: near-continuous CF search to minimize confusion of estimated gradients
- Robust small-signal behavior: include P0 as an inputreferred noise, accounted for in the SNR maximization; noise-only level parameters include P0
- Easily obtain optimal detection K given other parameters, instead of adding K to the search
- Feedback-type models: include a level-parameter convergence step when the main filter's output is used instead of a "passive" filter.
- Allow all models (GC, OZGF, PZFC, RoEx, etc.) to run in the same code, with minimal case switching
- Cache the level and CF between evaluations



### PZFC Fit#516

```
case 516 % order 2, 8 params
   FeedbackType = 1; % enable feedback iteration
   ModelName = 'PZFC' % pole-zero filter cascade
   ValParam = [ ...
% Final, Nfit = 516, 11-3 parameters, PZFC, cwt 0
                        0.00000 % SumSqrErr= 10226.95
     1.73848 0.00000
     0.62250 - 1.02349 \quad 0.94190 \% \text{ RMSErr} = 2.82994
     0.37208 0.00000 0.00000 % MeanErr = 0.00000
         Inf 0.00000 0.00000 % RMSCost =
                                                  NaN
     0.00000 0.00000 0.00000
     2.00000 0.00000 0.00000
     1.27403 -0.26291 0.21906
    11.30471 5.33017 0.33995
%
    -3.63143 -1.59230 4.68184 % Kv
     ];
   CtrlParam = [ ... % an 8-parameter fit
        0 0 % b1 zero BW relative to ERB
     1
     1 1 1 % B2
     1 0 0 % B21
     0 0 % c one extra zero maybe
     0 0
          0 % n1 unused
     0 0 % n2 order, stages per nominal ERB
     1 1 1 % frat Fzero:Fpole
1 1 1 % P0
     ];
```





case 360	
FeedbackType = 1; % end	able feedback iteration
$ModelName = 'CasGC_fb'$	% from 12-parameter fit
ValParam = [ %	
% Final, Nfit = 360, 14-3 pc	arameters, CasGC_fb, cwt 0
3.02522 1.15581 -1.	.72018 % SumSqrErr= 11201.52
-6.51804 2.64805 0.	.00000 % RMSErr = 2.96171
1.44194 -1.02186 0.	.00000 % MeanErr = -0.00000
0.01967 0.00292 0.	.00000 % RMSCost = NaN
1.77270 0.00000 0.	. 00000
3.81828 0.00000 0.	. 00000
0.00000 0.00000 0.	. 00000
0.00000 0.00000 0.	. 00000
4.00000 0.00000 0.	. 00000
9.69783 5.26747 3.	.93392
% -3.36401 -3.67324 4.	.80991 % Kv
];	
CtrlParam = [ %	
1 1 1 % b1	
1 1 0 % c1	
1 1 0 % Fr	
1 1 0 % Fr1	
1 0 0 % BZ	
エーエーズ PU つ・	
<b>」</b> ,	

### Cascade GC in feedback mode?

- Feedback can work on filter models for which it was not originally planned, sometimes.
- Notice:

```
case_360
```

```
FeedbackType = 1; % enable feedback iteration
ModelName = 'CasGC_fb' % from 12-parameter fit
```

 This works because the level-dependence affects the gain correctly (Fr: position of highpass active filter); attempting to use the bandwidth instead, as in CasGC Fit#316, doesn't work well, since the model equations hold the peak gain fixed independent of bandwidth





### OZGF\_fb Fit#100 (APGF)

#### **case** 100

```
FeedbackType = 1; % enable feedback iteration
ModelName = 'OZGF_fb' % APGF
ValParam = [ ...
 % Final, Nfit = 100, 7-3 parameters, OZGF_fb, cwt 0
 0.00000 0.00000 0.00000 % SumSqrErr= 14684.11
 0.55936 -0.97985 0.89312 % RMSErr = 3.39101
 0.66293 0.00000 0.00000 % MeanErr = 0.00000
 Inf 0.00000 0.00000 % RMSCost =
                                          NaN
 4.00000 0.00000 0.00000
 4.00000 0.00000 0.00000
 0.00000 0.00000 0.00000
 13.59306 8.04349 -1.35988
 % -3.04024 -0.96252 3.35127 % Kv
 ];
CtrlParam = [ ... % a 4-parameter fit
 0 0 % b1 unused in feedback version
 1 1 1 % B2
 1 0 0 % B21
                                  4-parameter
 0 0 0 % c
 0 0 0 % n1
                               4th-order APGF,
 0 0
      0 % n2
                               feedback version
 0 0 0 % frat unused
    1 1 % P0
 1
 ];
```





### Conclusion

- The filter fitting framework is a useful tool for seeing how proposed auditory filters relate to others
- The modified framework allows more types of models, including feedback configurations
- Being able to specify different model types in one framework will make this approach more accessible and useful to others
- The APGF, OZGF and PZFC are good auditory filter shapes

## Machine Hearing Research Agenda Why?

- Help "Machine Hearing" become a firstclass academic and commercial field like "Machine Vision"
- Motivate, plan, and promote my project activities
- Do something useful with all the uninterpretable audio media out there
- Motivate the recording of more...

### Three areas to get right:

- Leveraging techniques already developed in the machine-vision and machine-learning fields
- 2. Productive interaction with the wider field of hearing research, to keep models honest and motivate better experiments
- Focus on applications for which the challenge has to do with what things sound like, as opposed to specialized domain knowledge ("non-speech non-music audio")

### How we proceed...

- Good auditory front end, based on good hearing research, leads to representation of what things "sound like"
- Content-based retrieval of sound tracks (or other audio content) based on what it "sounds like is going on" is a good hard application
- Noisy data is good data
- Analogy to content-based image retrieval, based on features that encode what things "look like," leads to workable system structure



Stabilized Auditory Image (Correlogram)

副人



Auditory front end based on stable models from long ago, with new feature extraction ideas.

PAMIR multi-label retrieval (MLR) for the trainable back-end retrieval.

What about sound segmentation or separation?

Sound-retrieval precision in top k, for multi-label retrieval versus support-vector-machine classifier



Large query vocabulary (about 2000 words) makes it difficult. Results are not great yet with this baseline front end; auditory front end has been shown to help a lot, in smaller experiments so far

# Other ways to leverage machine vision

- In combination:
  - Audio/visual robots and 3D perception
  - Content-base retrieval, content classification, etc., based on joint sound/image features
- By analogy:
  - Object tracking -> sound source tracking
  - Key-point features –> key sound features?

### Other good applications

- Indexing, retrieval, summarization of personal audio diaries, movie soundtracks, etc.
- Real-time and retrospective analysis of audio security/surveillance recordings
- Front end to speech transcription systems
- •

Avoid humans are electrical in nature and so are the signals handled by chips. Already universities have produced silicon chips that can be narrow implanted in the human body to replace damaged nerves to restore views. the flow of signals from the nervous system to the brain.

And, to an extent, the five senses and the brain can already be SECTECY replicated in silicon. Limited machine vision and touch are available to industrial robots. Machine hearing is something on which a lot of research money is being spent around the world and a great deal of secrecy attaches to progress. Machines can recognize individual words, but are not so hot when it comes to sentences. None the less, the best people in the field think that machines capable of recognizing human speech will be on the market in the 1997-2000 time frame.

That opens a vast range of product possibilities from dictation machines that print out speech, to portable translators that speak out in a different language from the one spoken into them, and to a whole bunch of voice-operated machines, from telephones to cars to computers.

Smell and taste have had less research money than sight and hearing, probably because there are fewer commercial applications

1995

### Bake-offs

- As in speech and vision, good shared datasets with defined tasks can lead to good competitive/cooperative progress
- Quantitative performance testing requires lots of labeled training and testing data. How will we get it?

# Conclusion (memo to self)

- Machine hearing field is full of unrealized potential; needs some focus on strategies for progress
- Cooperation among groups and fields is key to rapid progress; don't be insular or re-invent the wheel