





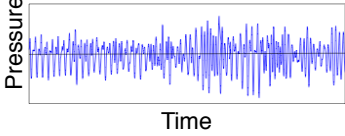
Sound, Ears, Brains and the World

Josh McDermott
Dept. of Brain and Cognitive Sciences, MIT
9.523

Consider some examples of typical auditory input:

- Scene from cafe: 
- Scene from sports bar: 
- Radio excerpt: 
- Barry White: 

The ear receives a pressure waveform.




AUDITION

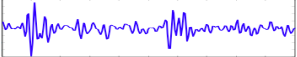
When objects in the world vibrate, they transmit acoustic energy through surrounding medium in the form of a wave.

The ears measure this sound energy and transmit it to the brain.

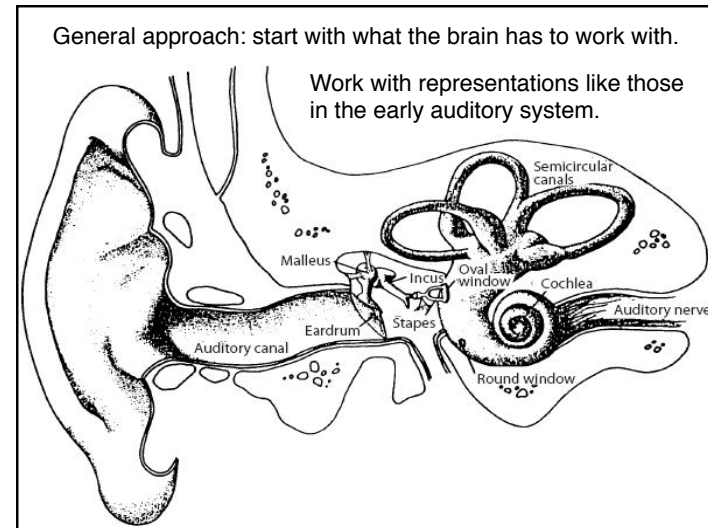
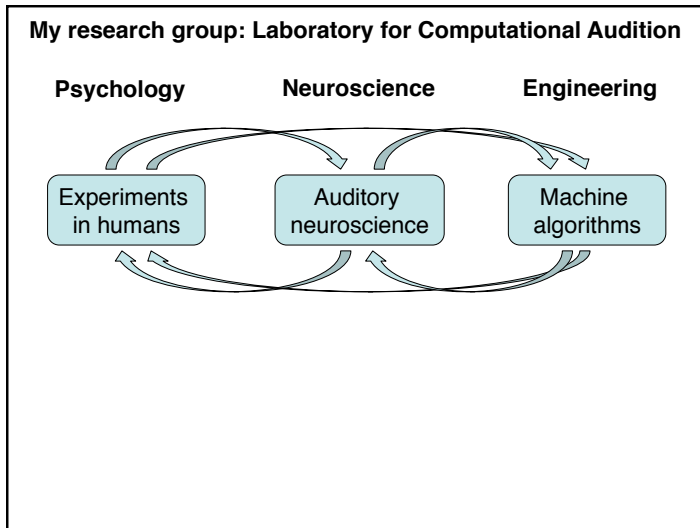
The task of the brain is to interpret this signal, and use it to figure out what is out there in the world.

The listener is interested in what happened in the world to cause the sound:



- Most properties of interest are not explicit in the waveform: 

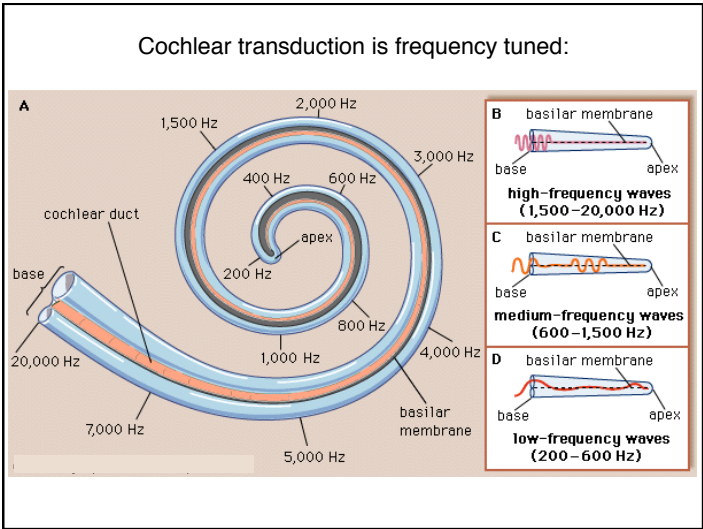
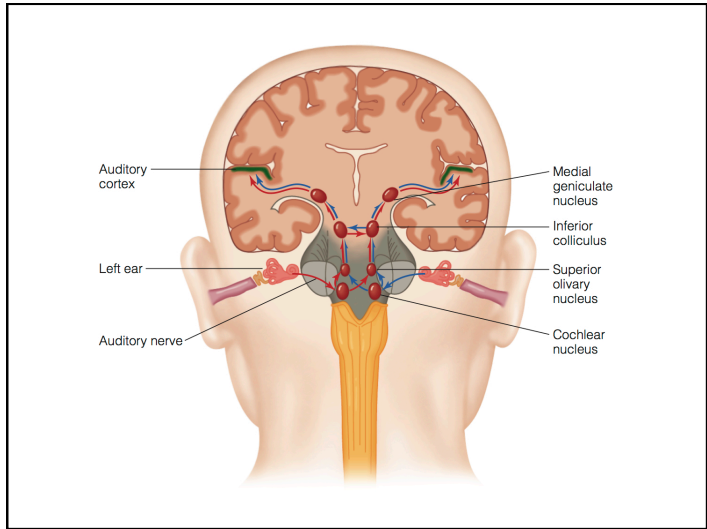
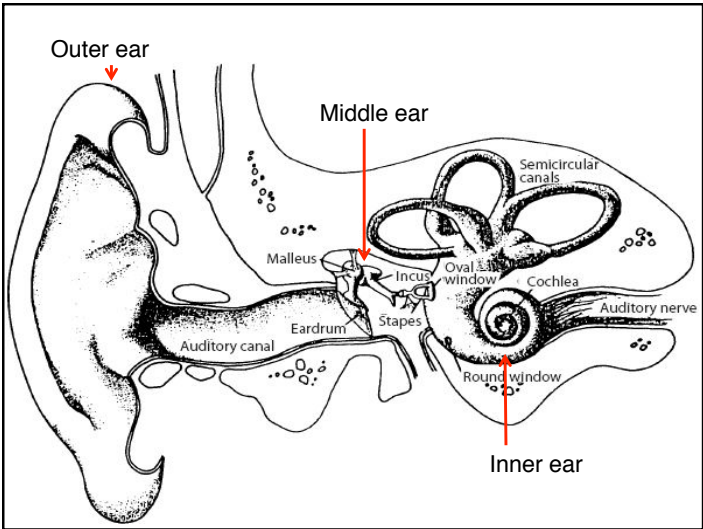
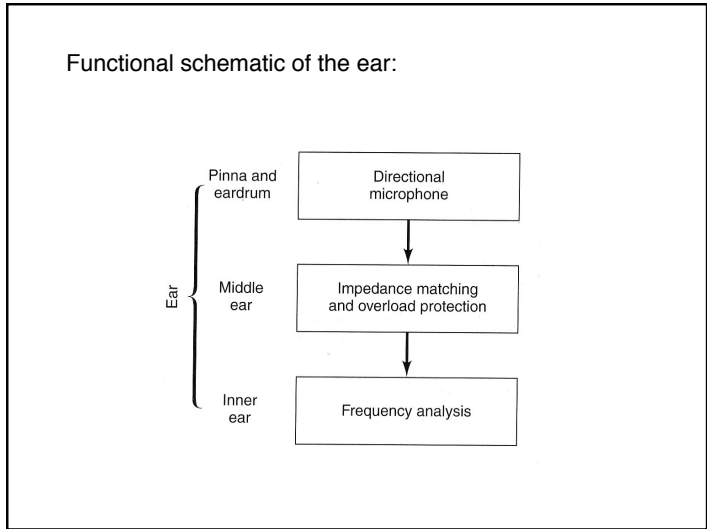
How do we derive information about the world from sound?

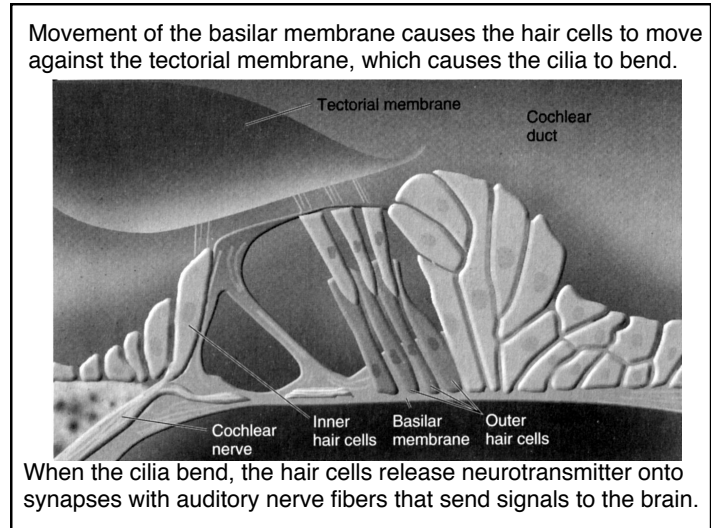
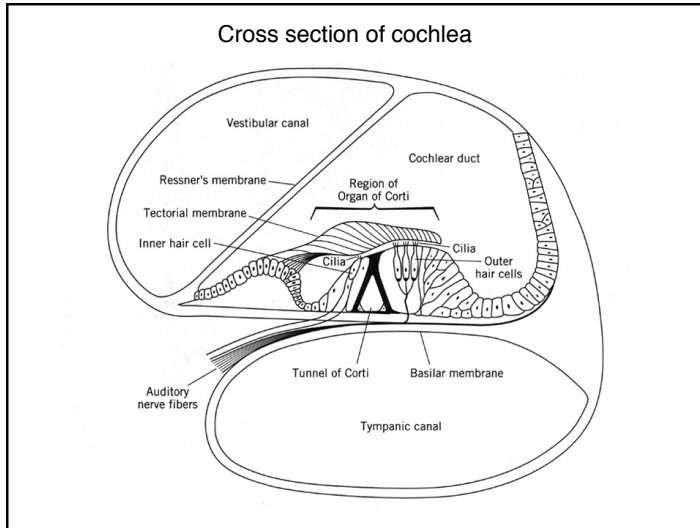


Plan for today:

1. Overview of Auditory System
2. Sound Texture Perception

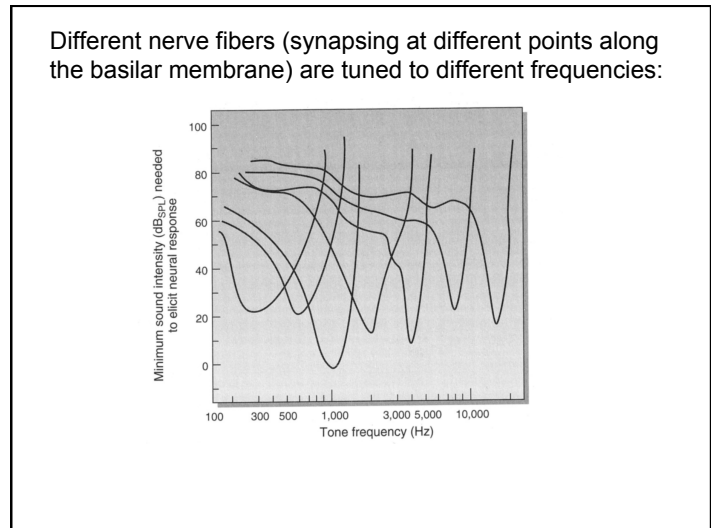
Part 1: Overview of Auditory System

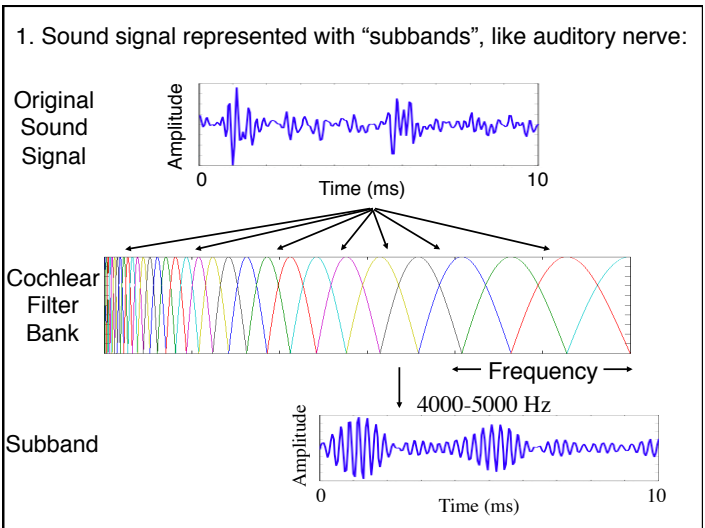
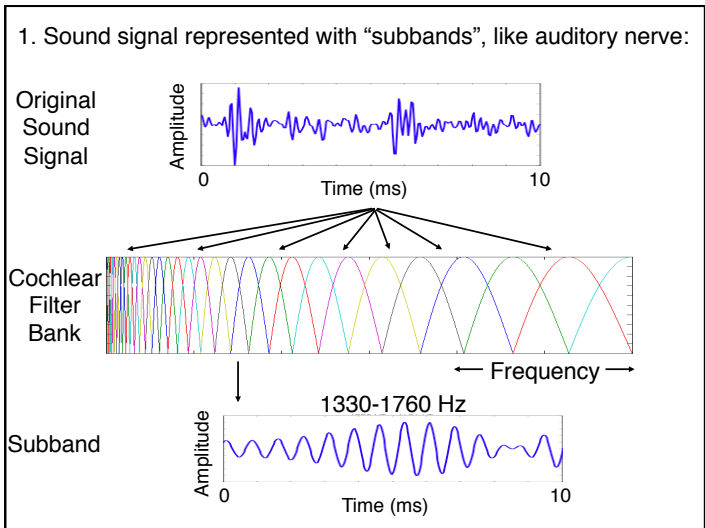
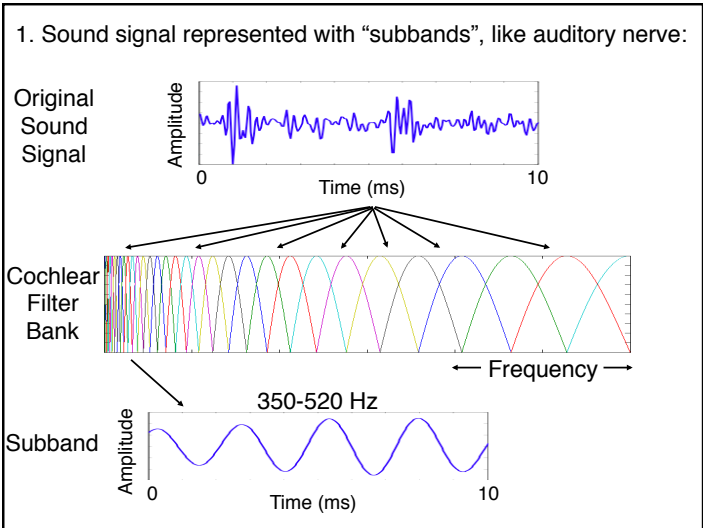
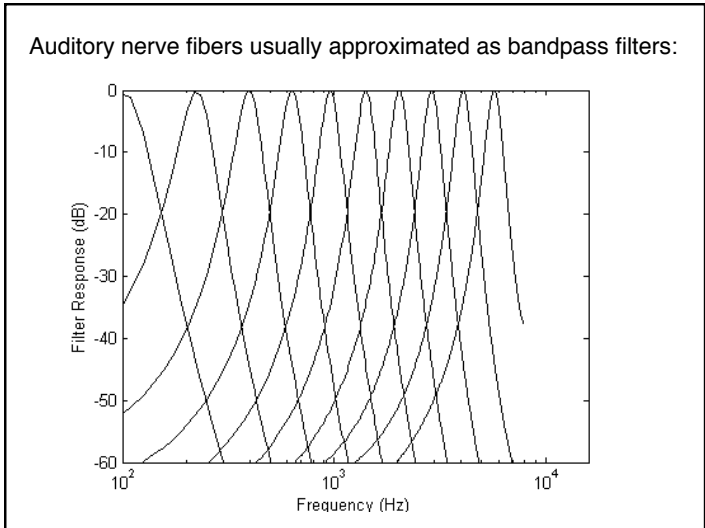


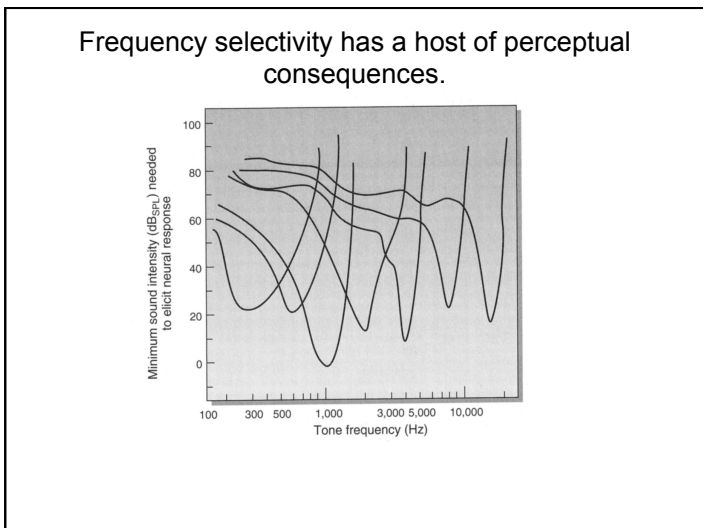
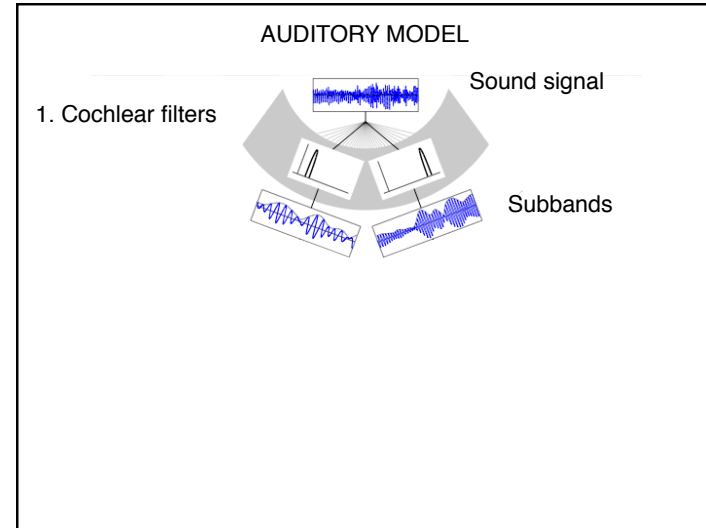
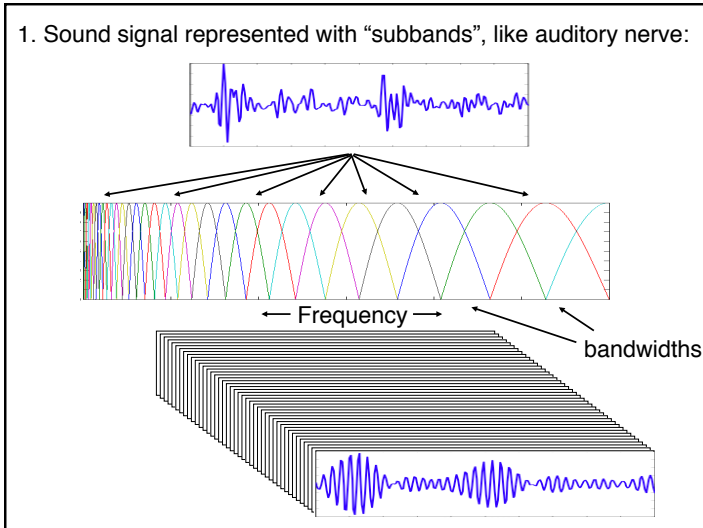


But because only part of the basilar membrane moves for a given frequency of sound, each hair cell and auditory nerve fiber signal only particular frequencies of sound.

One example:







Perception of beating constrained by freq. selectivity

Superposition of two pure tones waxes and wanes in amplitude.

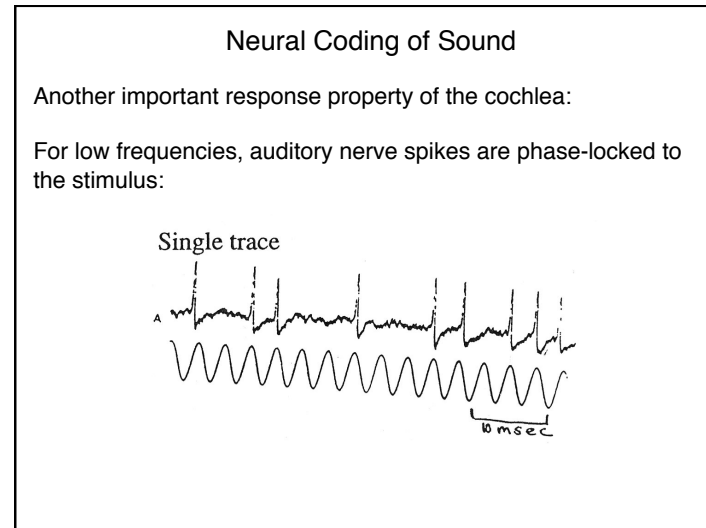
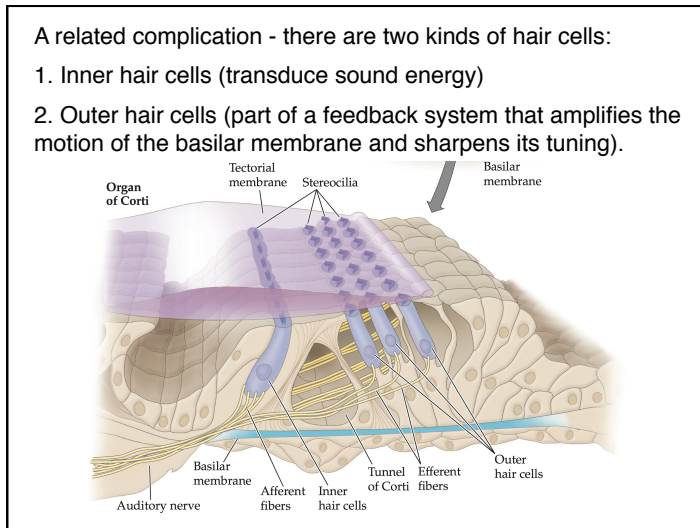
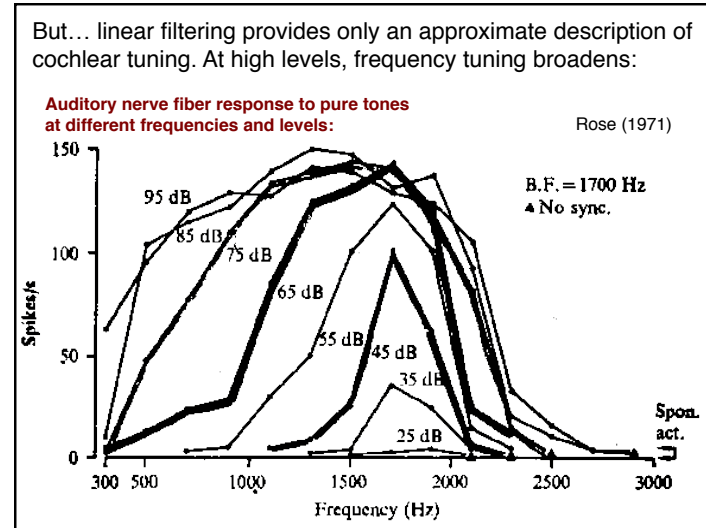
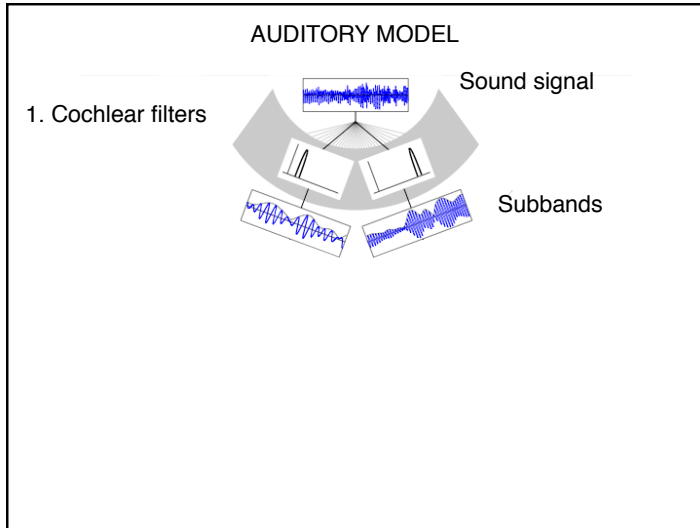
The perceptual correlate of rapid beating is known as roughness.

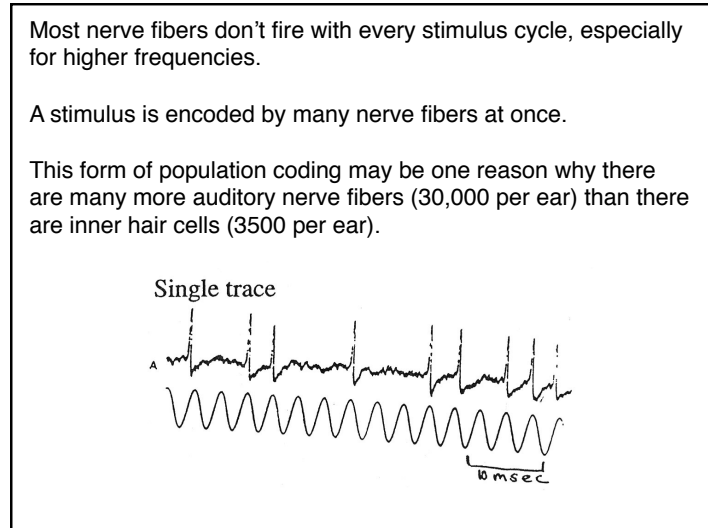
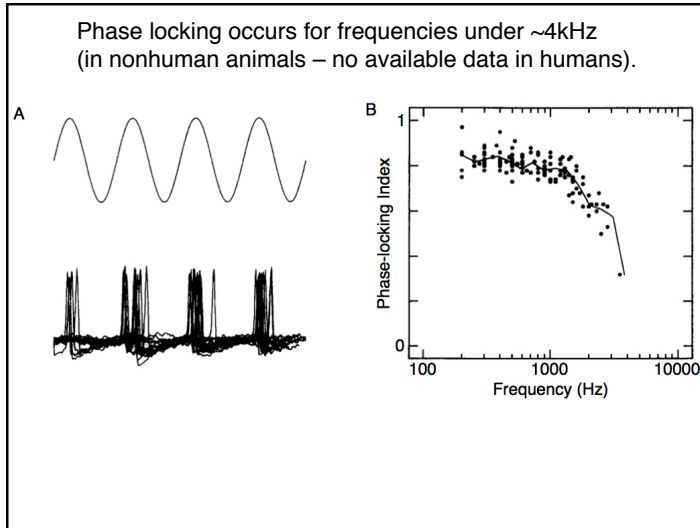
- Perception of beating is constrained by the cochlea:
 - Beats are only heard if two frequency components fall within the filter bandwidth of the cochlea:

1 semitone frequency difference:

3 semitones:

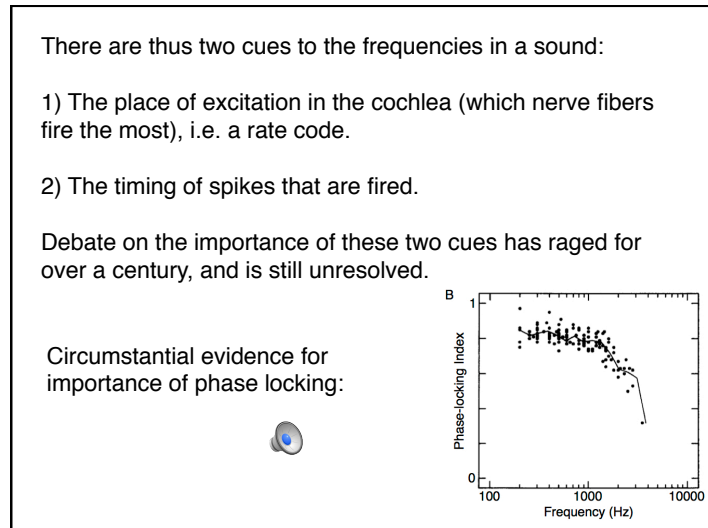
8 semitones:

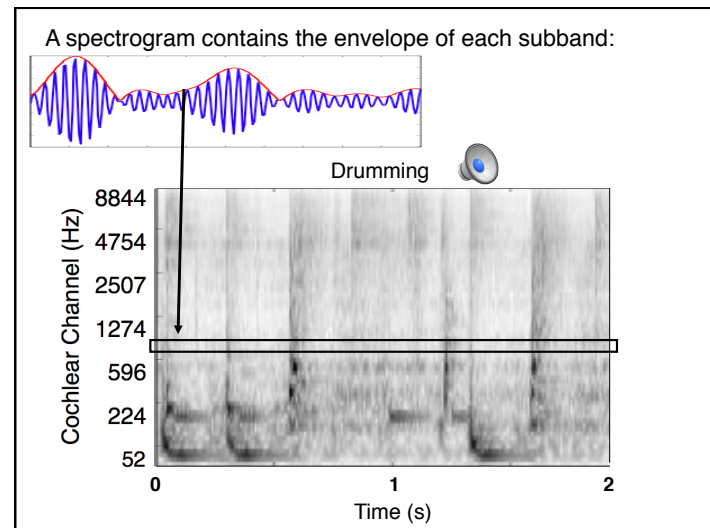
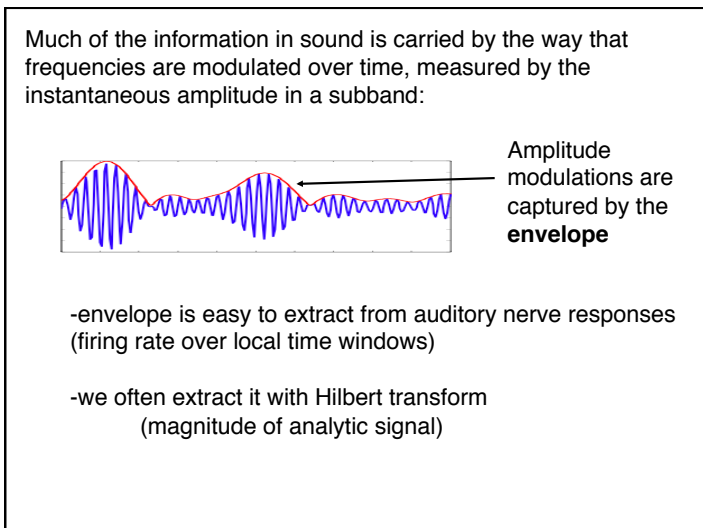
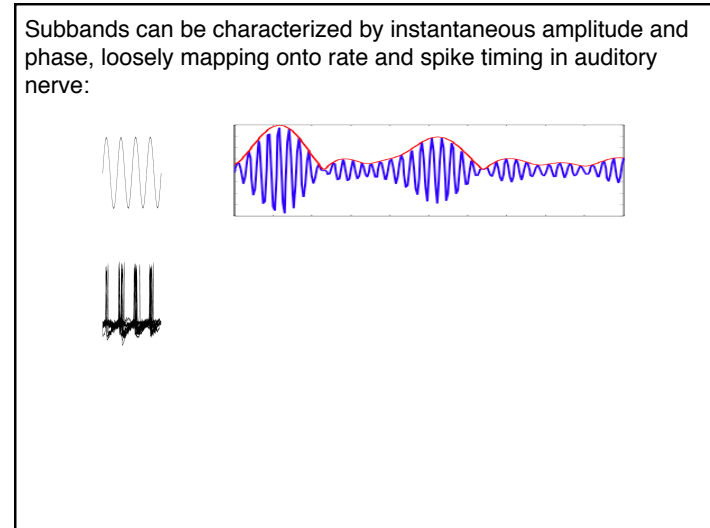
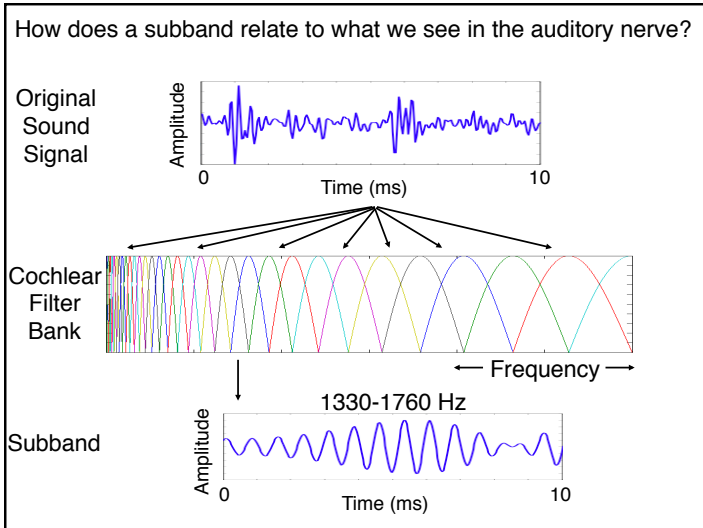




Some interesting numbers:

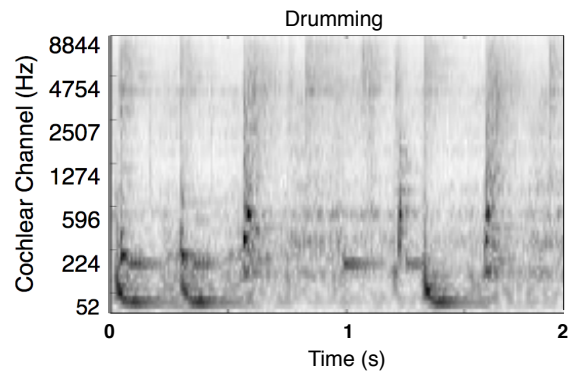
Per ear:	Per hemisphere:
3500 inner hair cells	60 million neurons in primary auditory cortex?
12,000 outer hair cells	
30,000 auditory nerve fibers	
Per eye:	
5 million cones	140 million neurons in primary visual cortex
100 million rods	
1.5 million optic nerve fibers	



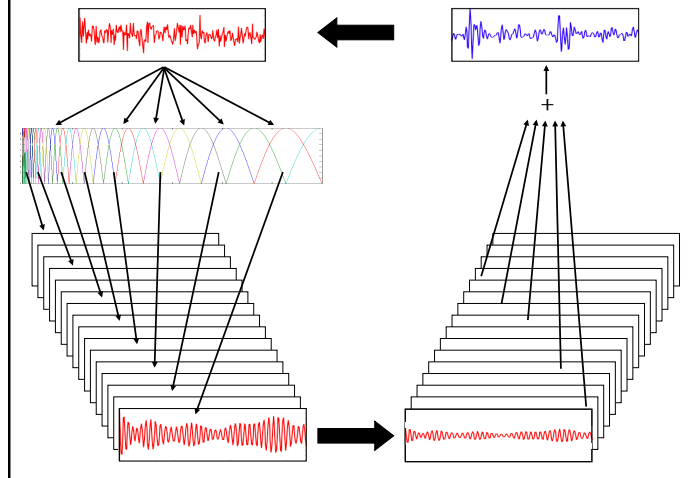


Envelopes often capture all the information that matters perceptually.

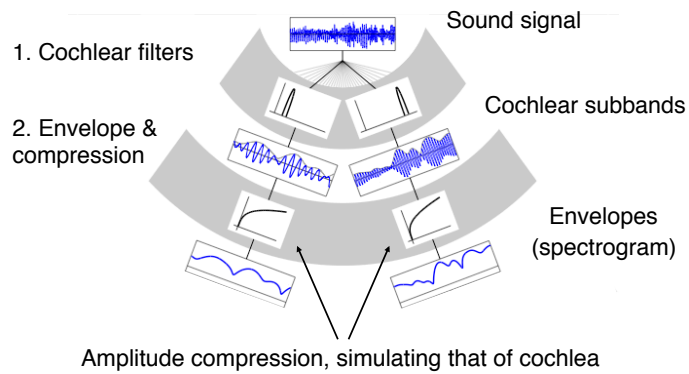
Sounds can be reconstructed just from the envelopes.



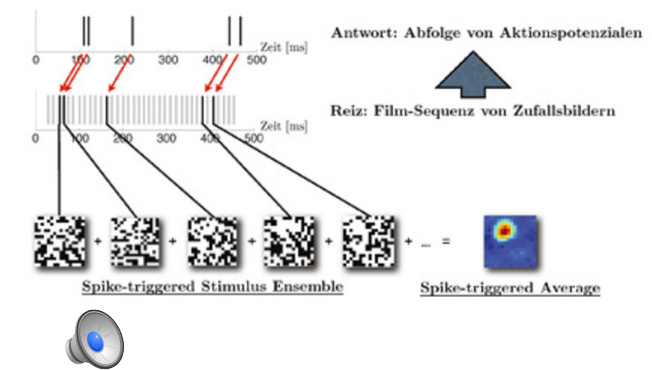
Start with noise, replace with envelopes, resynthesize, iterate:

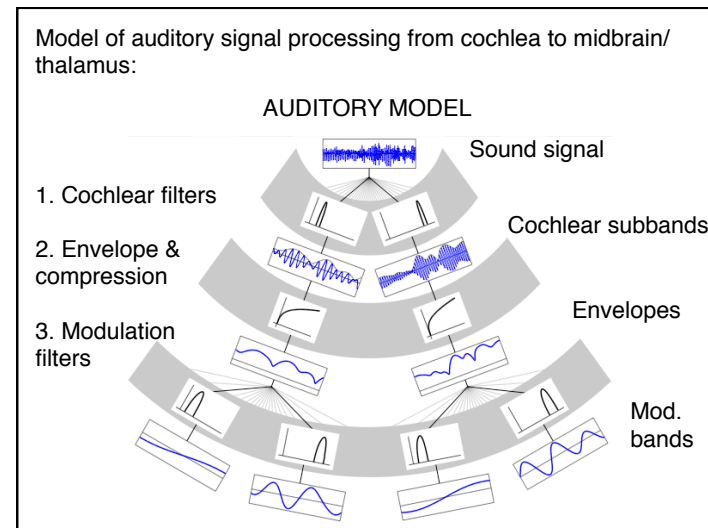
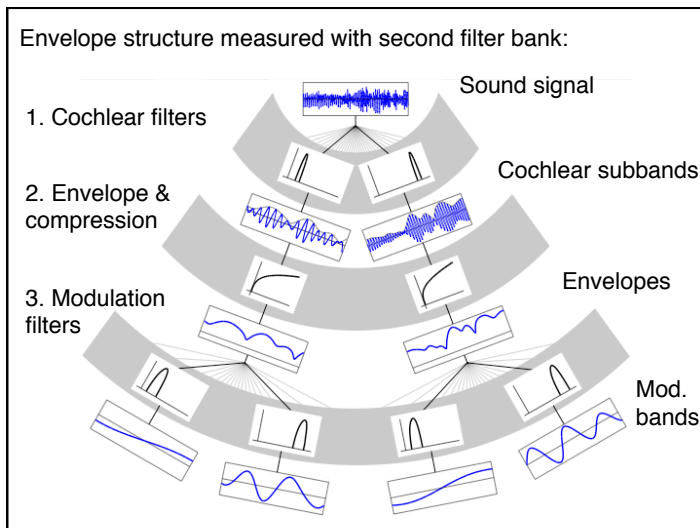
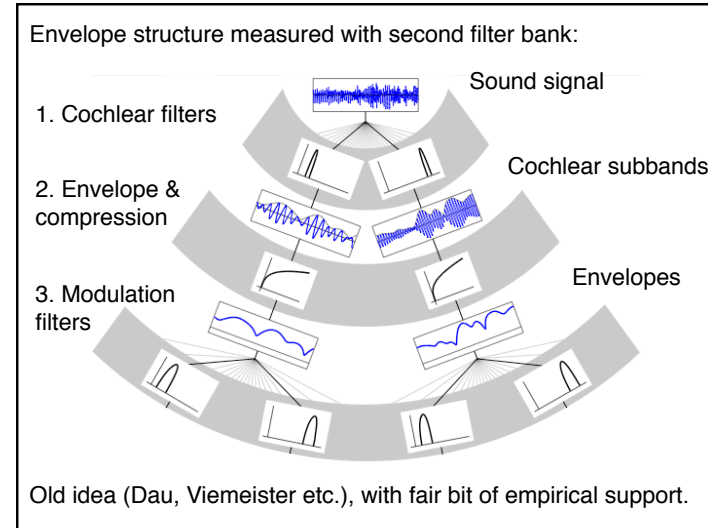
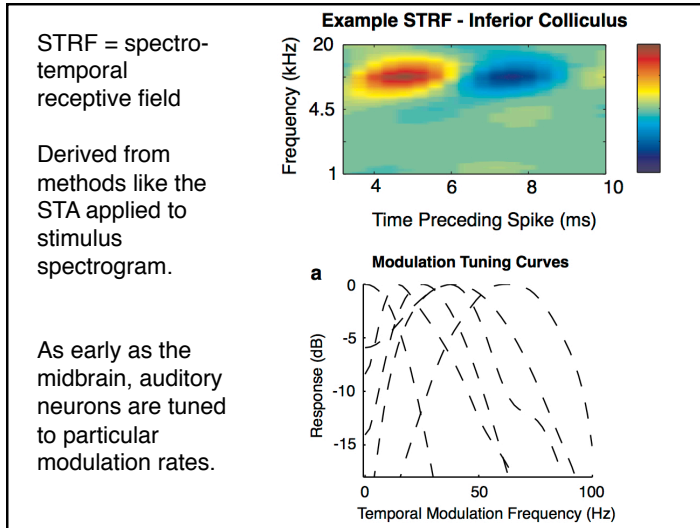


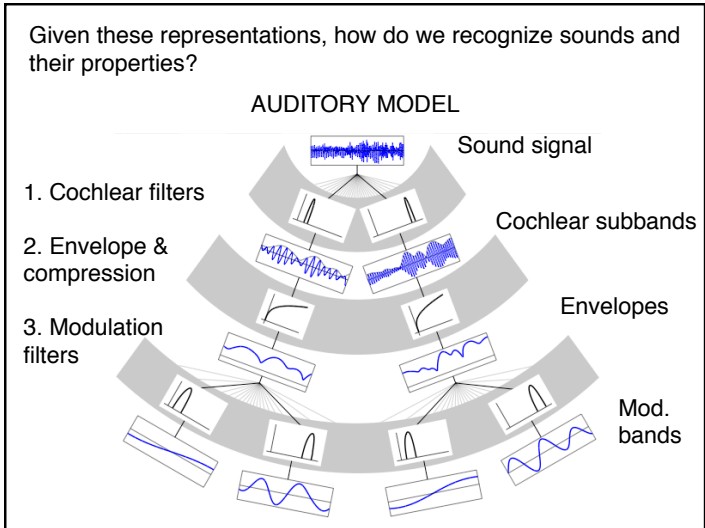
AUDITORY MODEL



Spike-triggered average: a method to characterize a neuron's receptive field.







Part 2: Sound Texture

SOUND TEXTURE

Textures result from large numbers of acoustic events.

- rain
- wind
- birds in a forest
- running water
- insects at night
- crowd noise
- applause
- fire

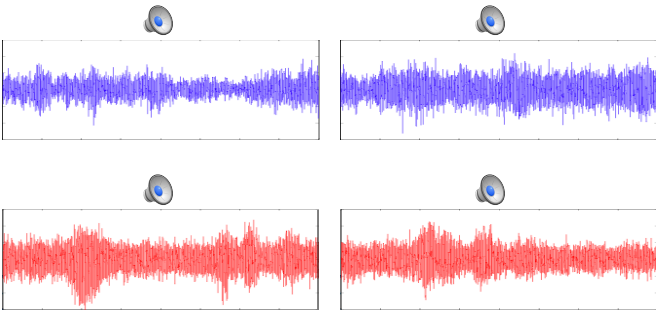
Common in the world, though historically neglected.

Much of hearing research is concerned with the sounds of individual events:

Unlike event sounds, textures are stationary - essential properties do not change over time.

- Stationarity makes textures a good starting point for understanding auditory representation.

How do people represent, recognize sound textures?

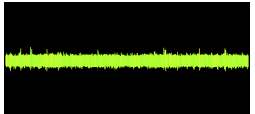
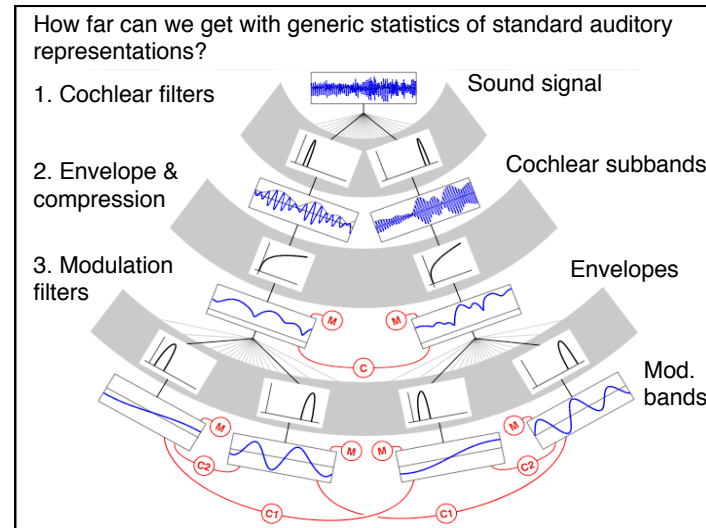
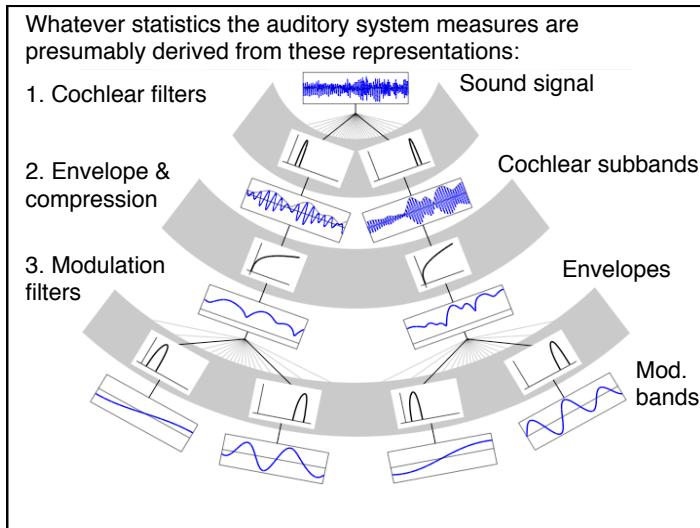


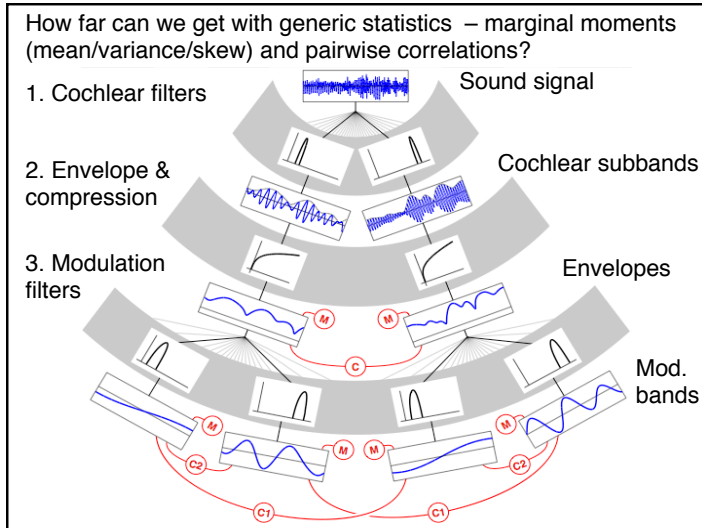
What do you extract and store about these waveforms to recognize that they are the same kind of thing?

Key Theoretical Proposal:

- Because they are stationary, textures can be captured by statistics that are *time-averages* of acoustic measurements.
- When you recognize the sound of fire or the sound of rain, you may be recognizing these statistics.

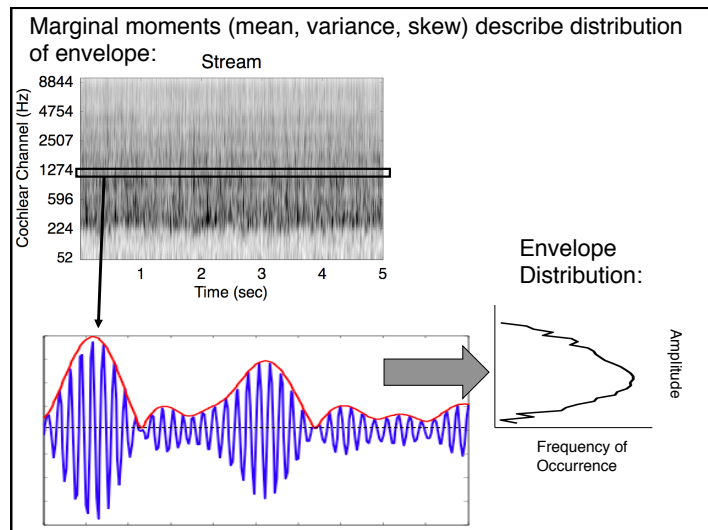
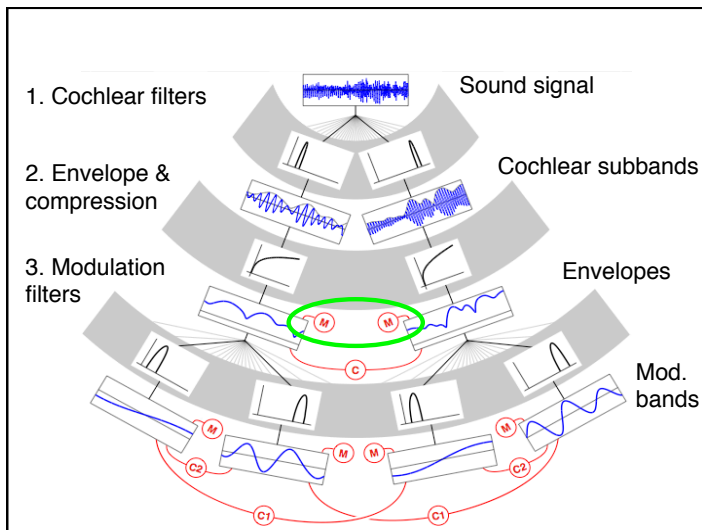
What kinds of statistics might we be measuring?

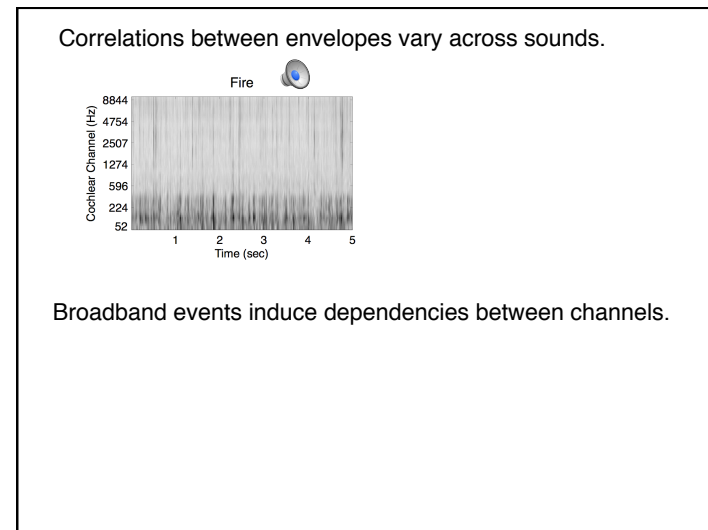
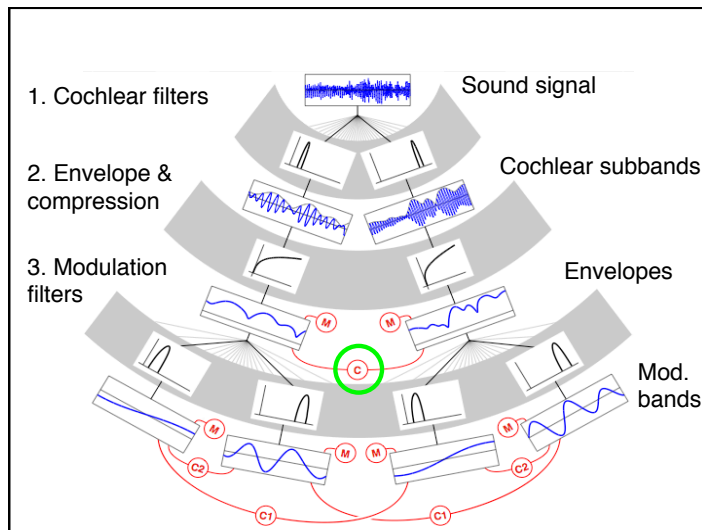
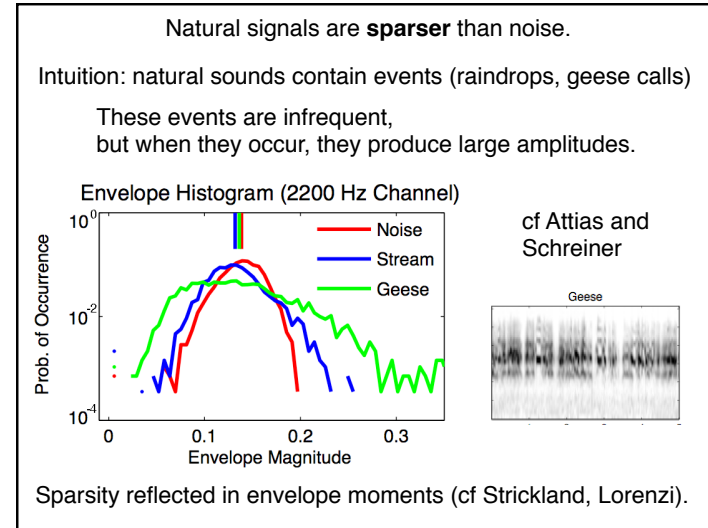
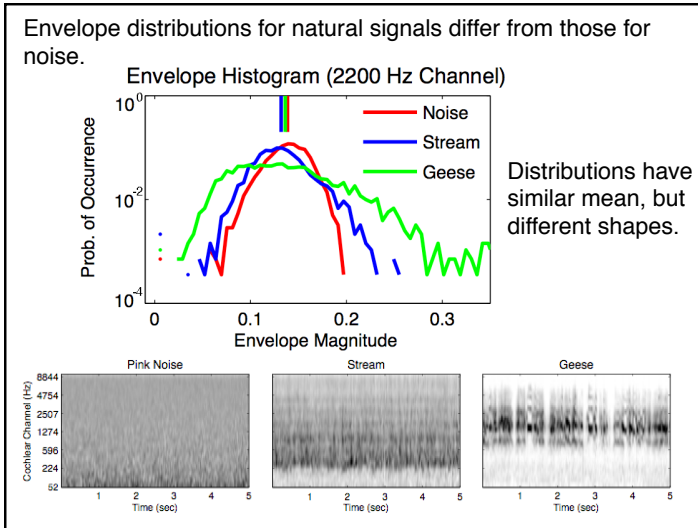



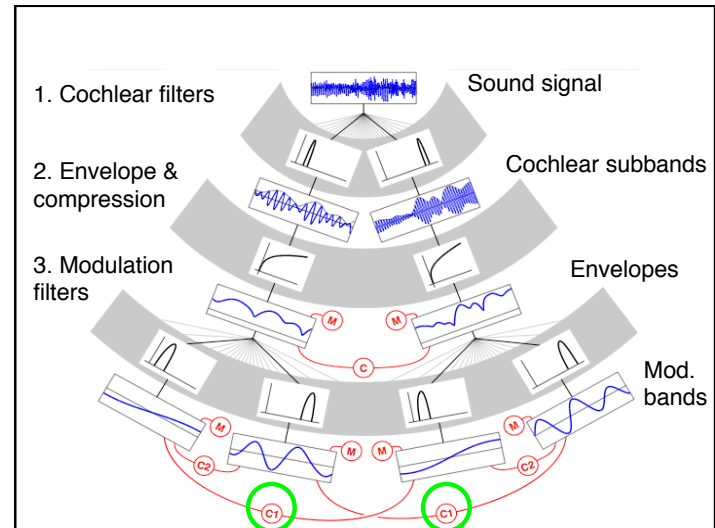
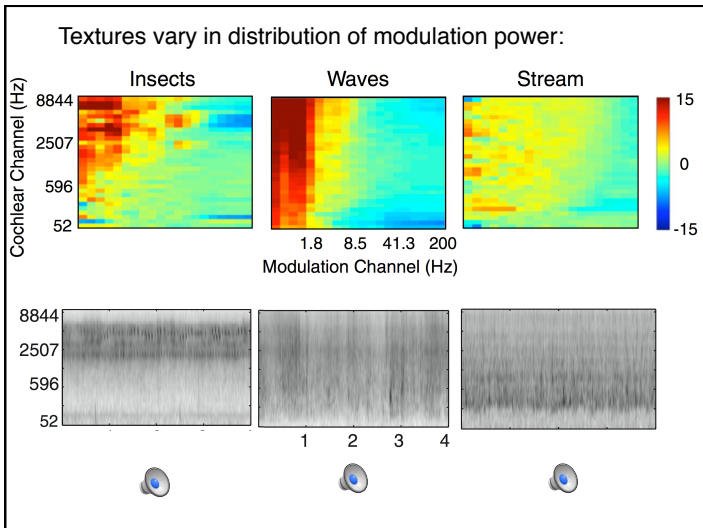
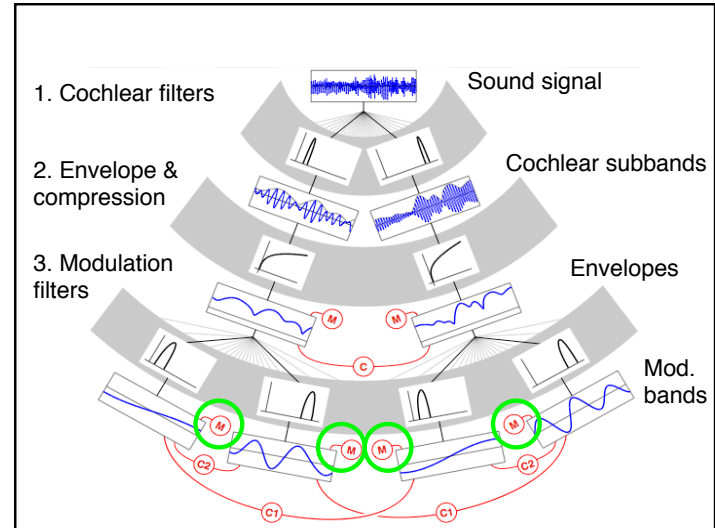
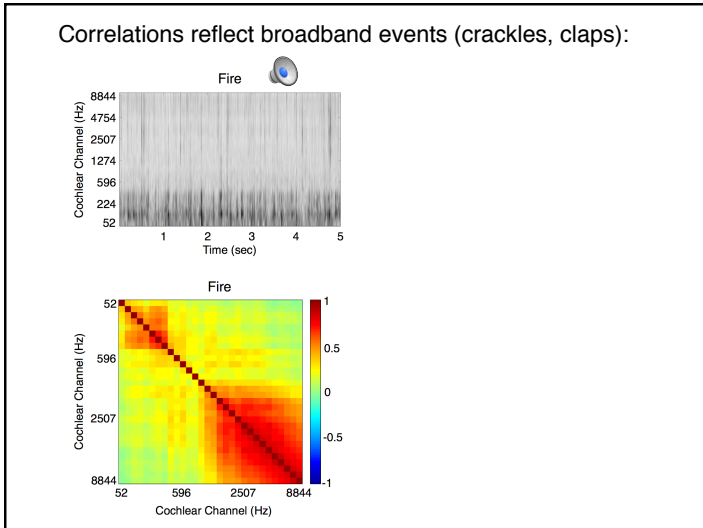


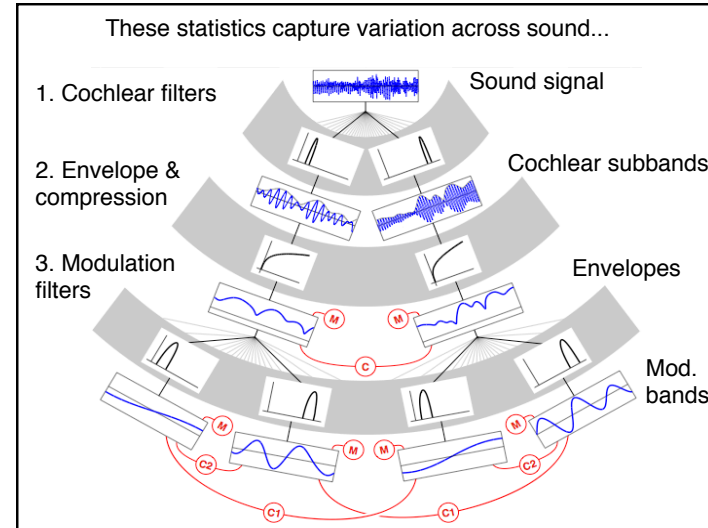
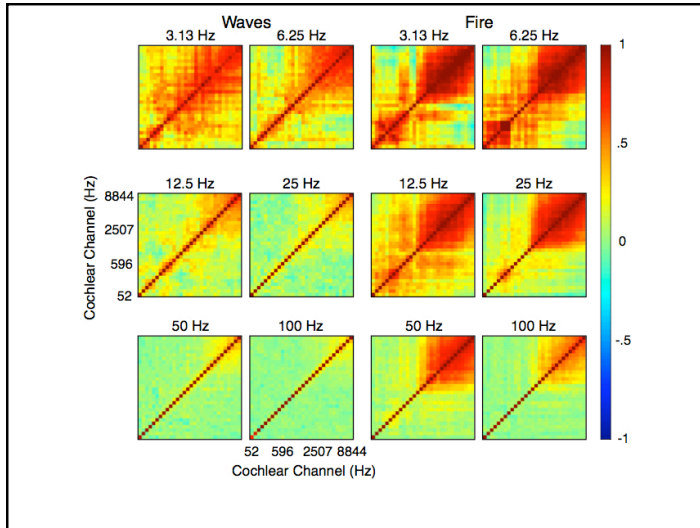
- Statistics are not specifically tailored to natural sounds
- Ultimately, would be nice to learn statistics from data (stay tuned...)
- But moments and correlations are interpretable, and might give insight.

To be useful for recognition, statistics need to give different values for different sounds...









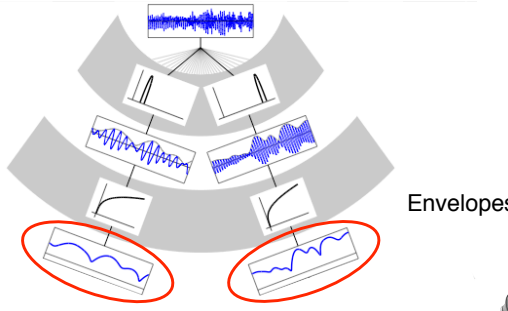
Can they account for perception of real-world textures?

Key Methodological Proposal:

- Synthesis is a powerful way to test a perceptual theory.
- If your brain represents sounds with a set of measurements, then:
 - Signals with the same values of those measurements should sound the same.
- Sounds synthesized to have the same measurements as a real-world recording should sound like it...

IF the measurements are what the brain is using to represent sound.

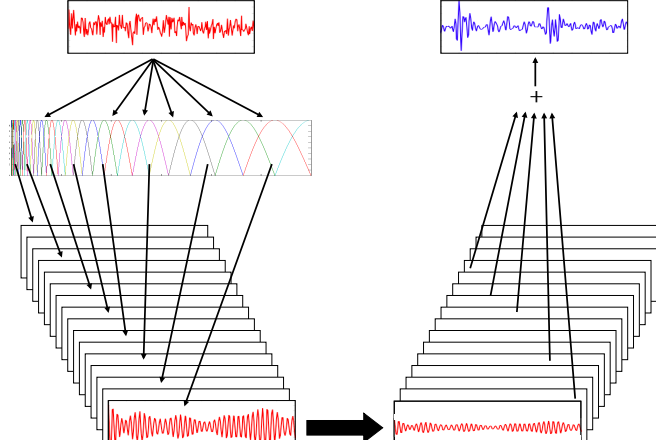
Simple example: test the role of the mean of each cochlear envelope (power spectrum)













Envelopes

- Measure average value of each envelope in real-world texture
- Then synthesize random signal with same envelope means.

Start with noise, rescale noise subbands, resynthesize:

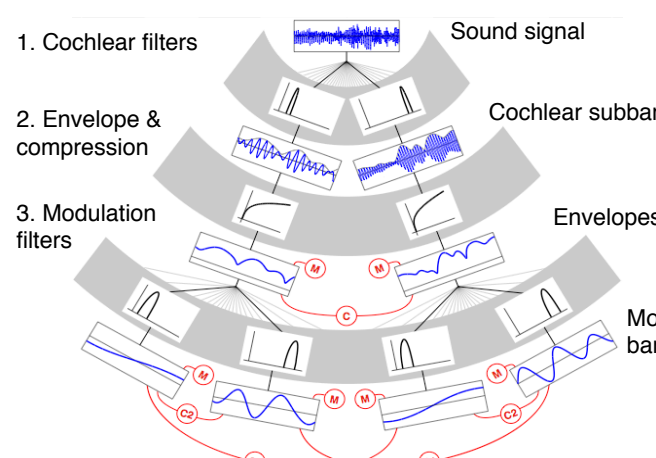


What do they sound like?

	Original	Power	
Rain			<ul style="list-style-type: none"> •Synthesis is not realistic (everything sounds like noise): •We aren't simply registering the spectrum (mean values of envelopes) when we recognize textures.
Stream			
Bubbles			
Fire			
Applause			

Will additional simple statistics do any better?

1. Cochlear filters
2. Envelope & compression
3. Modulation filters

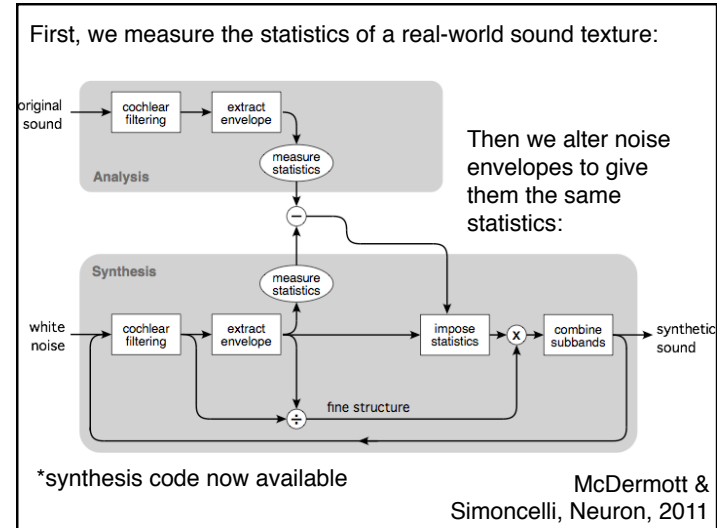
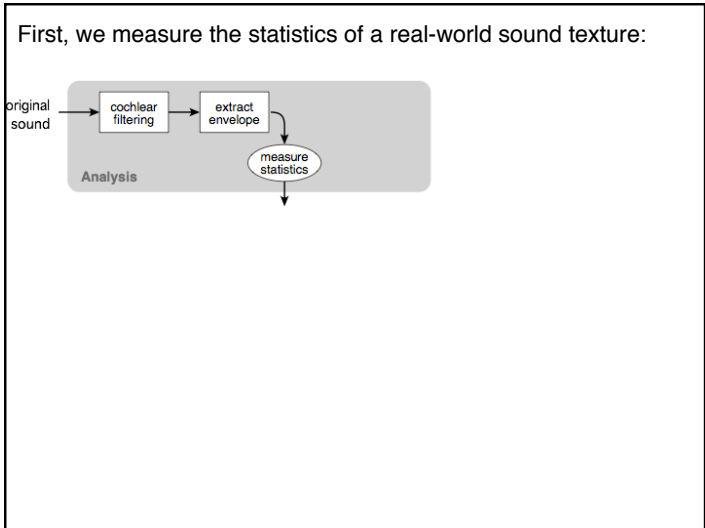


Sound signal

Cochlear subbands

Envelopes

Mod. bands



The result: a signal that shares the statistics of a real-world sound.







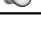
How do they sound?

If statistics account for texture perception, synthetic signals should sound like new examples of the real thing...

With marginal moments and pairwise correlations, synthesis is often compelling:

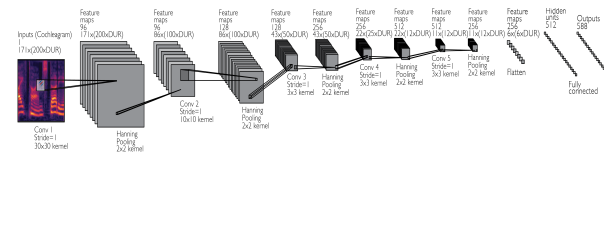
	Original	Power	All Stats
Rain			
Stream			
Bubbles			
Fire			
Applause			
Wind			
Insects			
Birds			
Crowd			

Also works for many “unnatural” sounds:

	Original	Power	All Stats
Rustling Paper			
Train			
Helicopter			
Jackhammer			

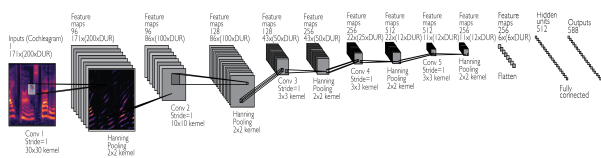
Success of synthesis suggests these statistics could underlie representation and recognition of textures.

Can we learn good representations of texture?



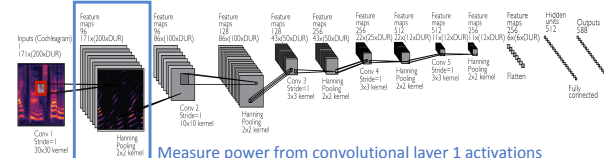
Jenelle Feather

Can we learn good representations of texture?

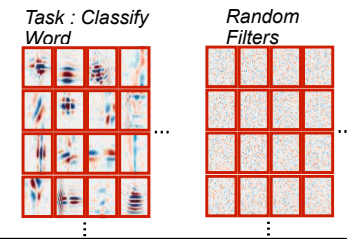


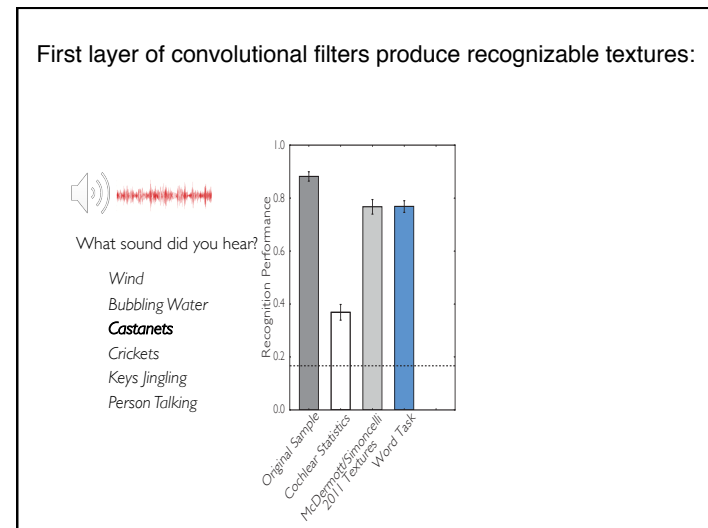
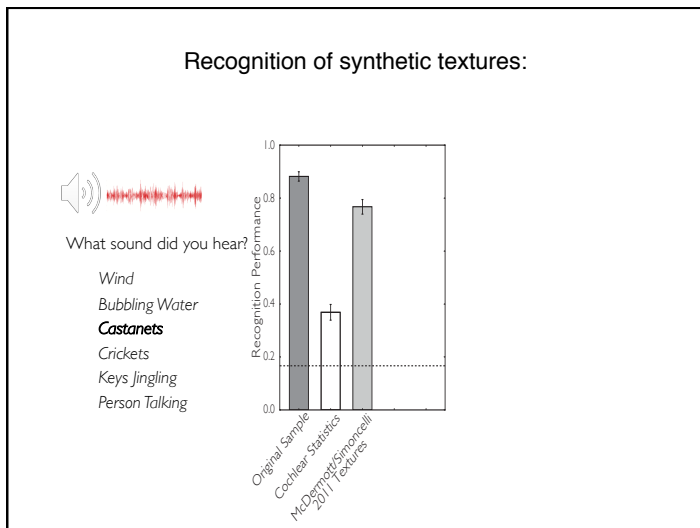
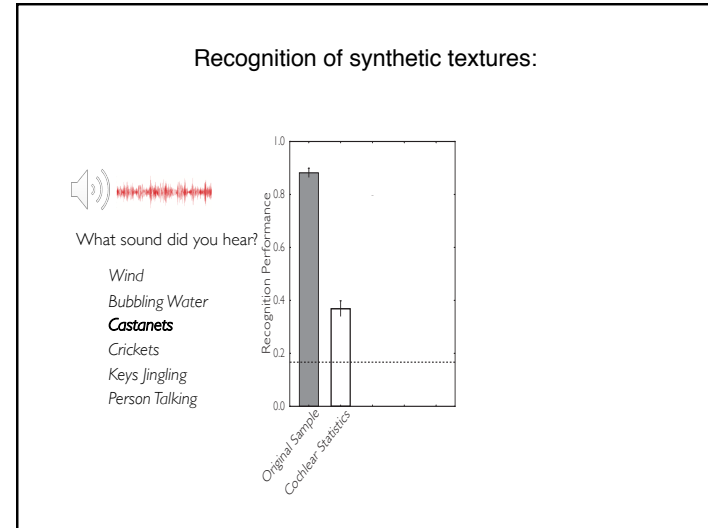
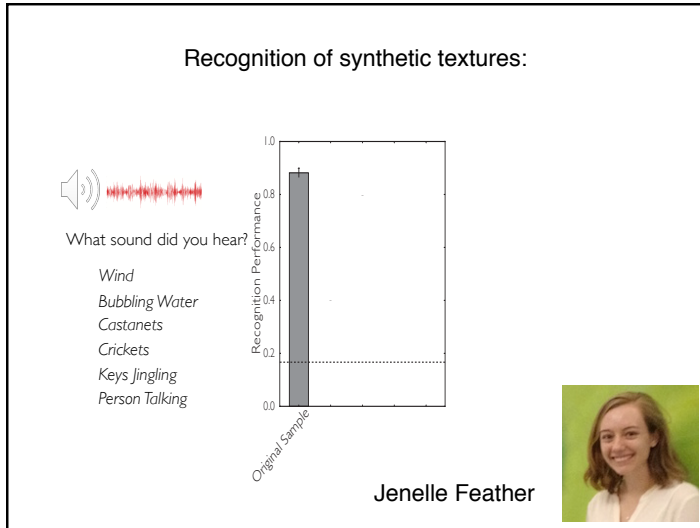
Jenelle Feather

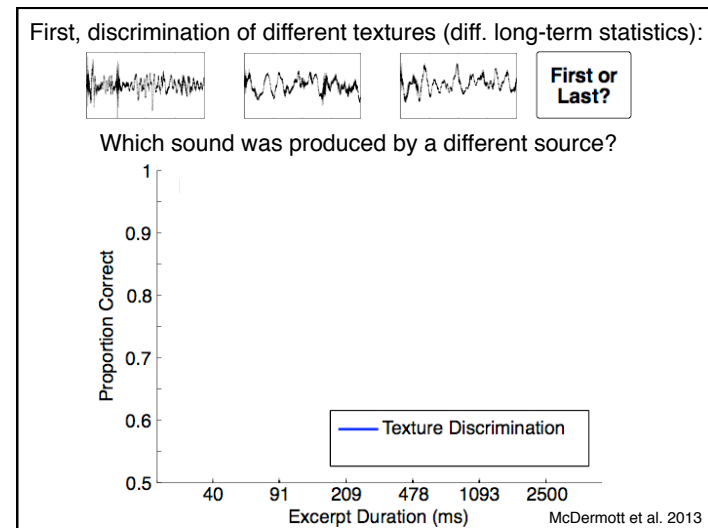
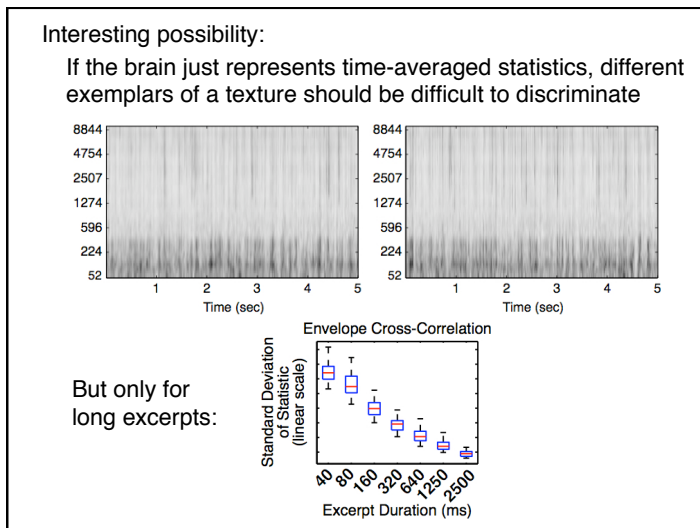
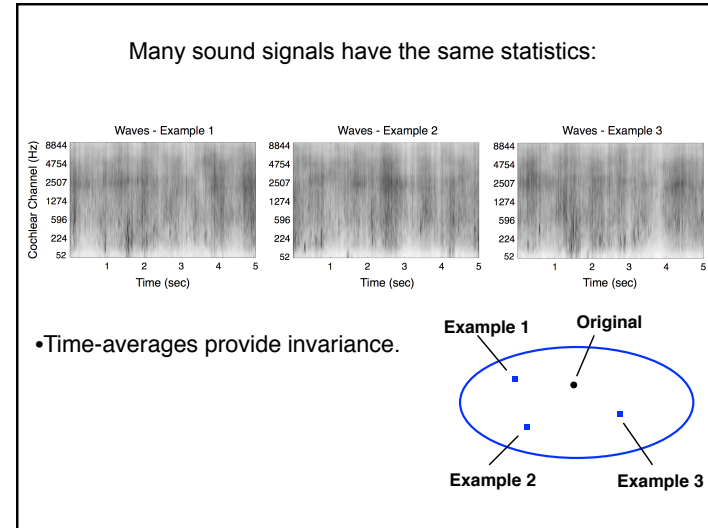
Can we learn good representations of texture?

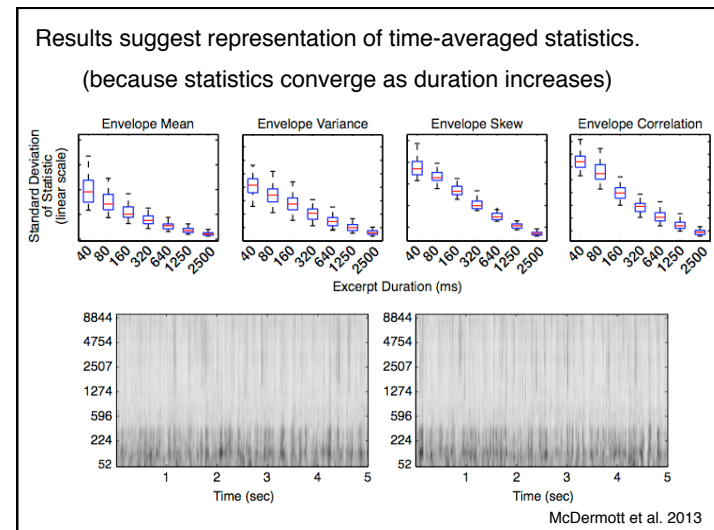
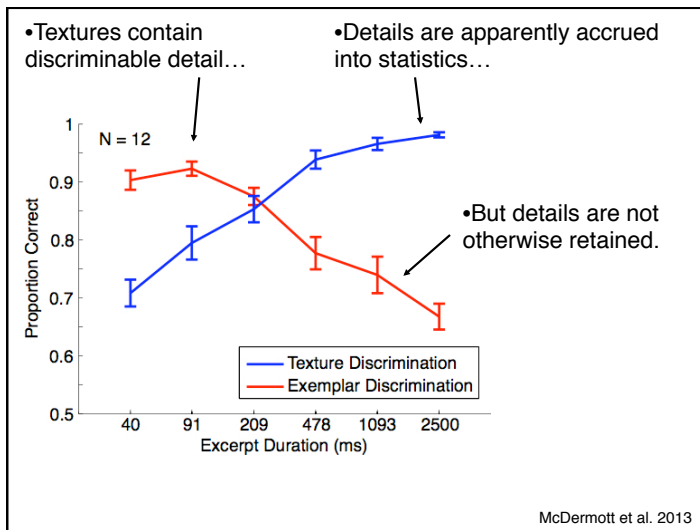
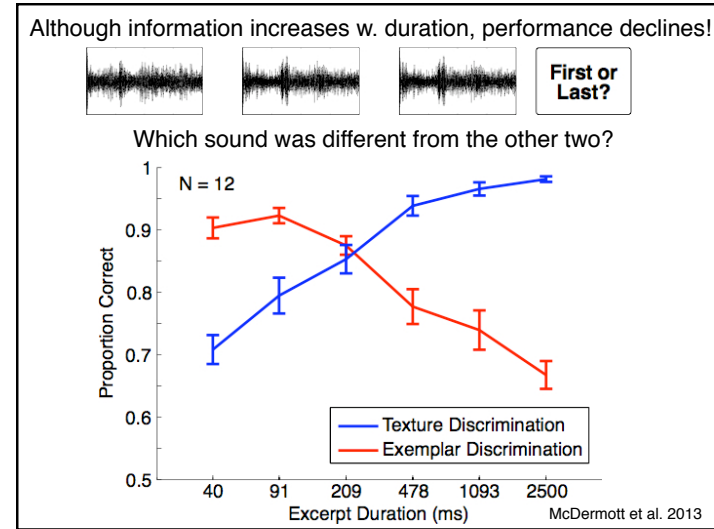
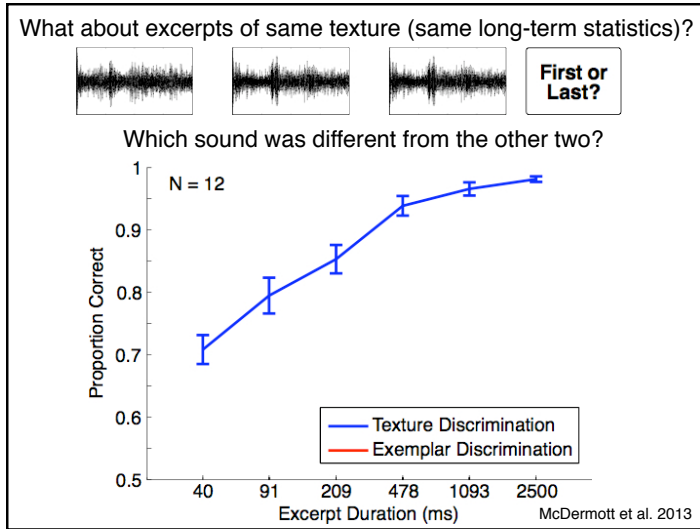


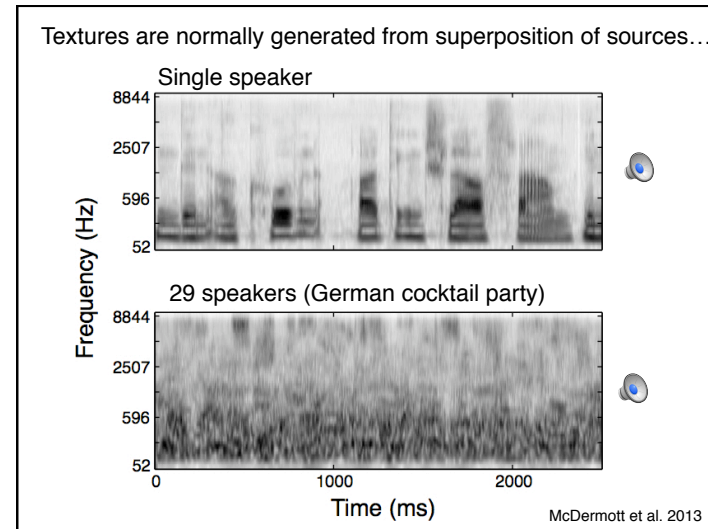
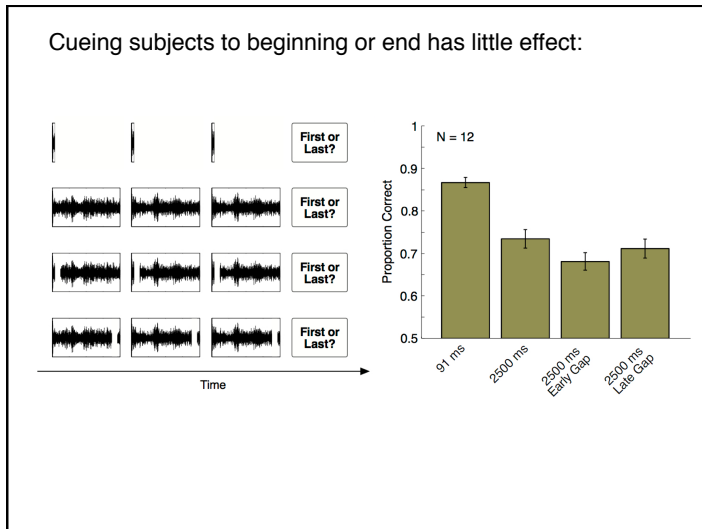
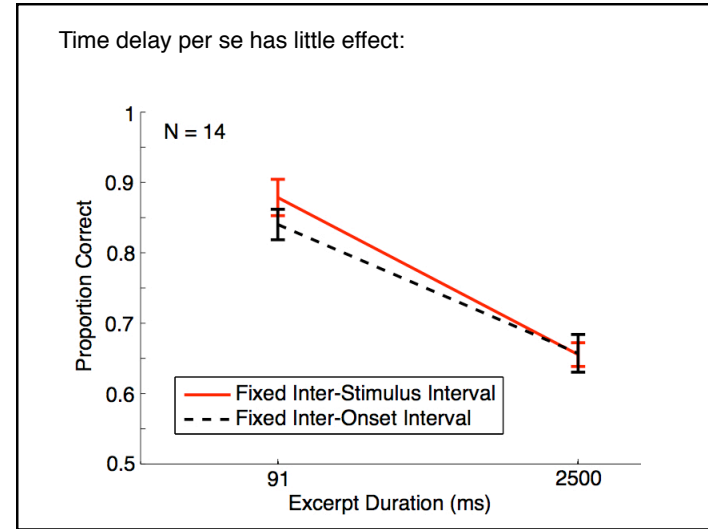
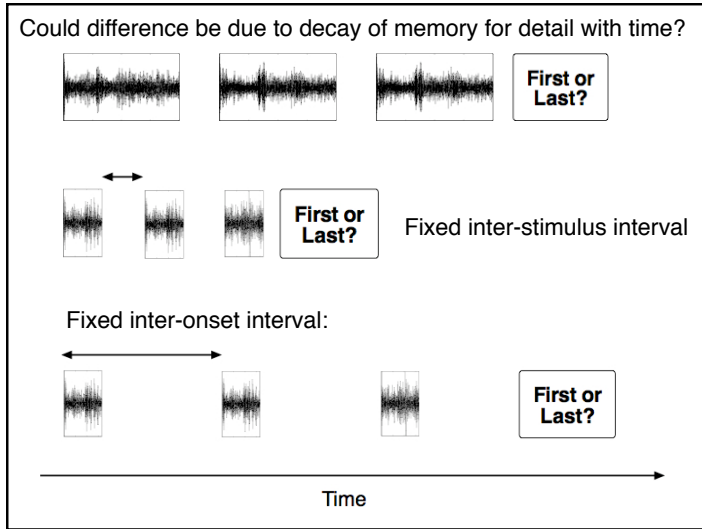
Measure power from convolutional layer 1 activations

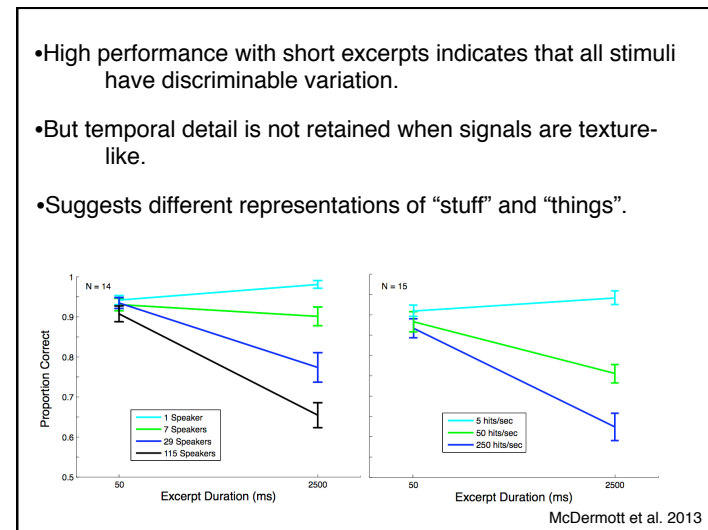
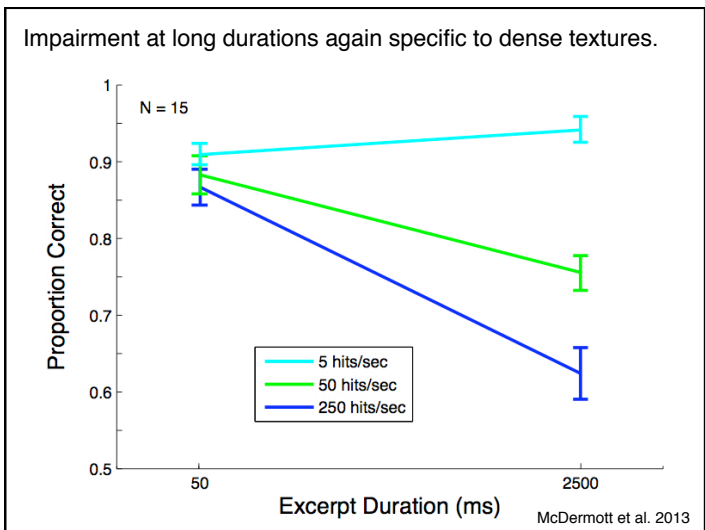
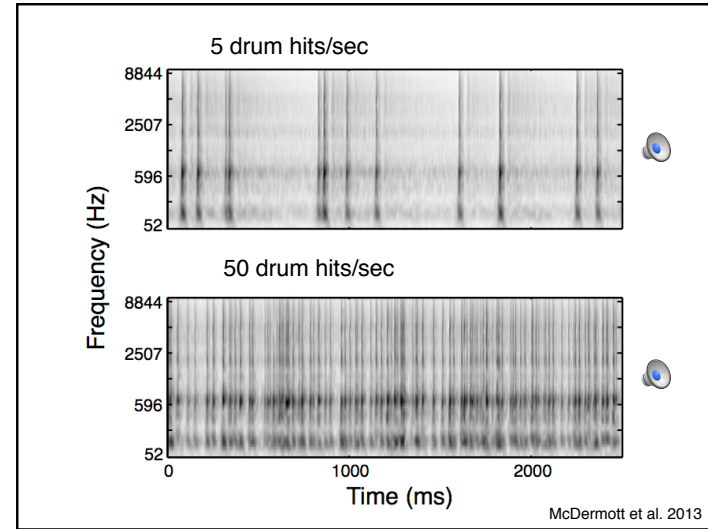
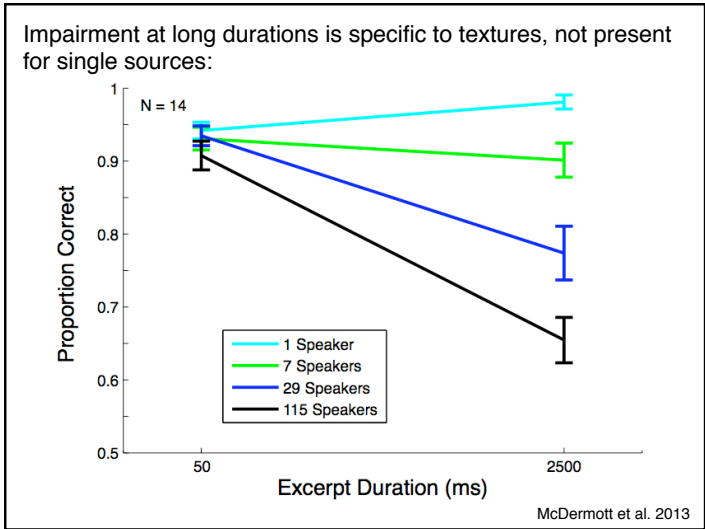












Over what time window are statistics averaged?

Idea: present signals whose statistics undergo a change:

Judgments of end state of texture should be biased depending on length of averaging window.

Changes in statistics enabled by morphing:

Richard McWalter

Experimental Paradigm

Which signal was closer to reference texture?
(Base judgments on end of standard)

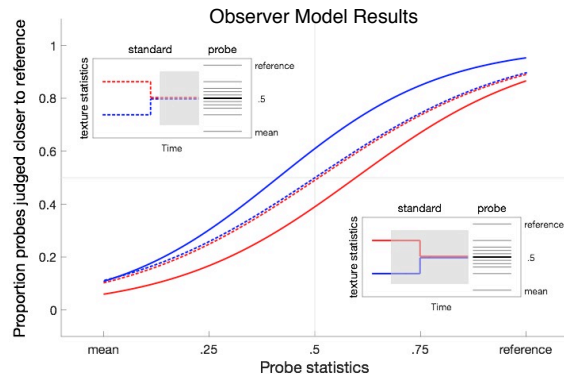
McWalter & McDermott, 2018

Experimental Paradigm

Which signal was closer to reference texture?
(Base judgments on end of standard)

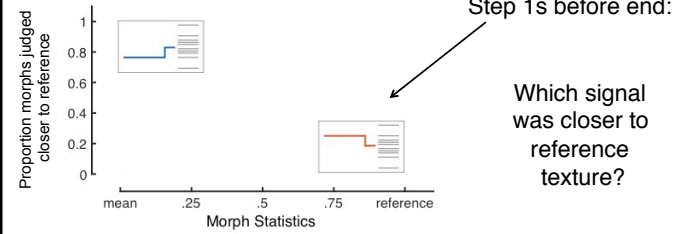
McWalter & McDermott, 2018

Bias in judgments should depend on averaging window extent:



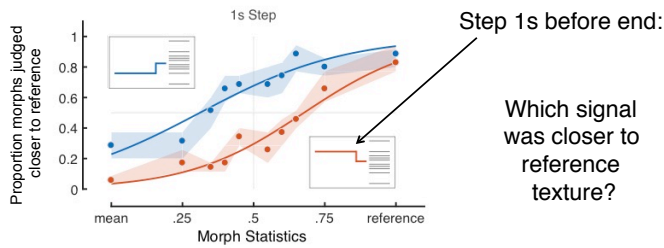
McWalter & McDermott, 2018

Bias in judgments indicates integration of stimulus history:



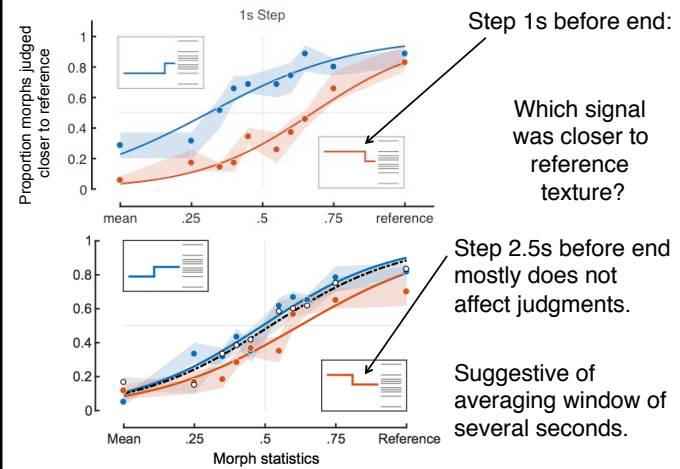
McWalter & McDermott, 2018

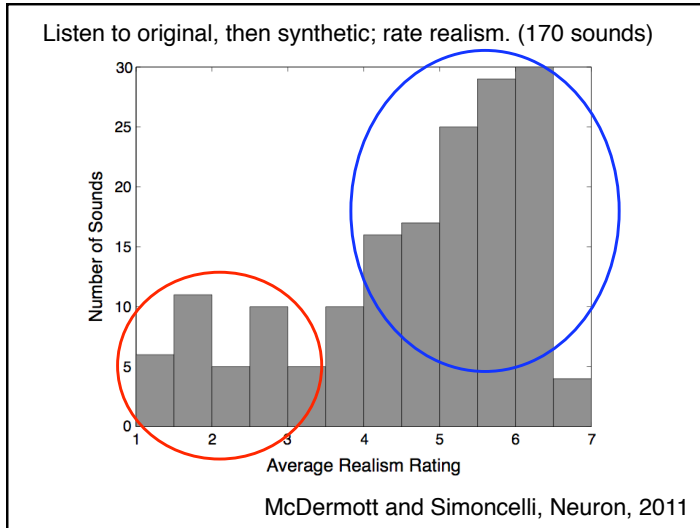
Bias in judgments indicates integration of stimulus history:



McWalter & McDermott, 2018

Bias in judgments indicates integration of stimulus history:





Lowest rated sounds are among most interesting, as they imply brain is measuring something model is not:

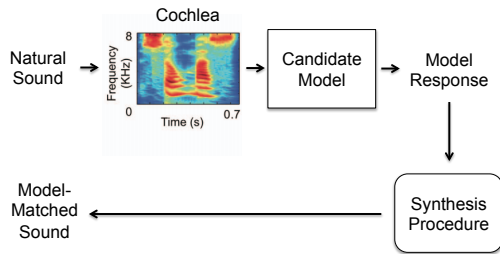
Pitch	1.93	Railroad crossing	
Rhythm	1.90	Tapping rhythm - quarter note pairs	
Pitch	1.77	Wind chimes	
Reverb	1.77	Running up stairs	
Rhythm	1.70	Tapping rhythm - quarter note triplets	
Reverb	1.67	Snare drum beats	
	1.63	Walking on gravel	
Reverb	1.60	Snare drum rimshot sequence	
Rhythm	1.60	Music - drum break	
Pitch	1.50	Music - mambo	
Rhythm	1.50	Bongo drum loop	
Reverb	1.47	Firecracker explosions	
Pitch	1.40	Person speaking French	
Pitch	1.37	Church bells	
Pitch	1.20	Person speaking English	

- ### TAKE-HOME MESSAGES
- Sound synthesis can help us test/explore theories of audition.
 - variables that produce compelling synthesis could underlie perception.
 - synthesis failures point the way to new variables that might be important for the perceptual system.
 - Textures are a nice point of entry into real-world hearing
 - Many natural sounds may be recognized with relatively simple statistics of early auditory representations
 - simplest statistics (spectrum) are not that informative
 - slightly more complex statistics are quite powerful
 - for textures of moderate length, statistics may be all we retain

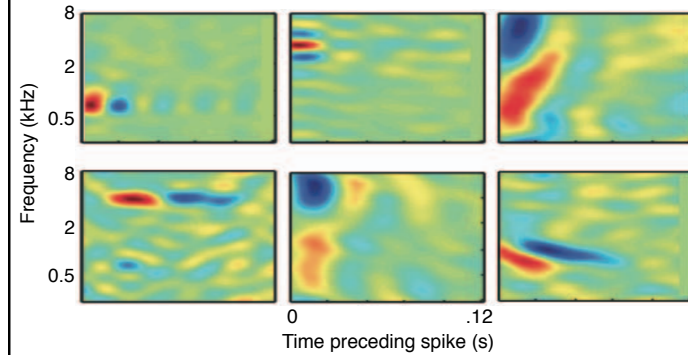
- ### OPEN QUESTIONS
- Locus of time-averaging?
 - Presumptive integration windows of several seconds are long relative to typical timescales in auditory system...
 - Relation to scene analysis?
 - What happens when foreground sounds are superimposed on a texture?
 - What statistics are needed to account for synthesis failures?

Another application of model-based synthesis:

- Idea: present a natural sound, and a synthetic signal that produces same response in a model
- If model is good description of neural response, responses to natural and model-matched sounds should be similar.

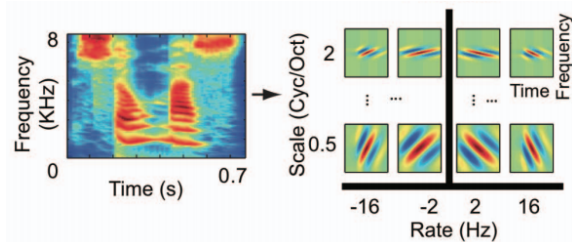


Linear filter characterizations in auditory cortex of animals (w.r.t. spectrogram):



Mesgarani et al. 2007

Common model of auditory cortex: linear spectrotemporal filtering

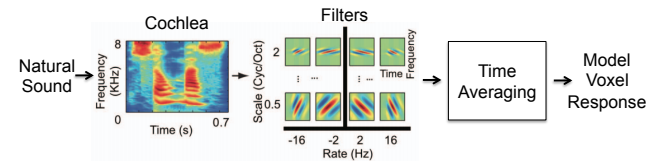


Very useful (and widely used), but clearly limited.


Shamma and colleagues

Another application of model-based synthesis:

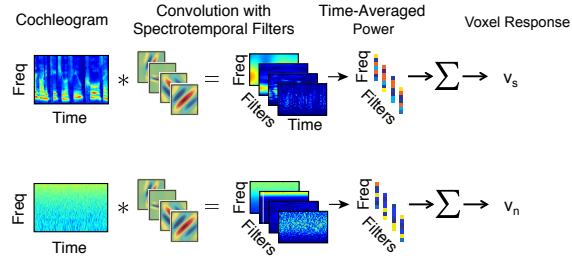
- Idea: present a natural sound, and a synthetic signal that produces same response in a model
- If model is good description of neural response, responses to natural and model-matched sounds should be similar.
- To test models with fMRI, we model the BOLD signal from a voxel as a sum of time-averaged responses of each unit (filter) in model.



- Work by Sam Norman-Haignere

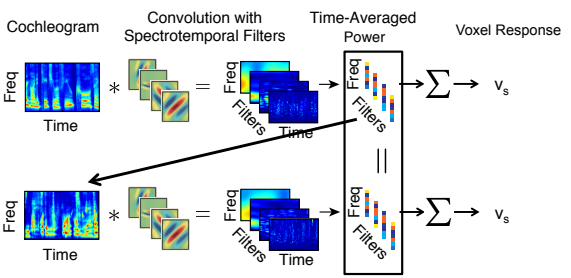


Model-Matched Synthetic Stimuli



Measure time-averaged statistics from a natural sound
 Measure same statistics for a noise sound

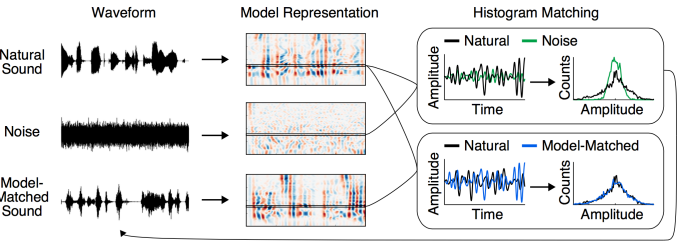
Model-Matched Synthetic Stimuli



Measure time-averaged statistics from a natural sound
 Measure same statistics for a noise sound
 Coerce a noise signal to have the matched statistics
 => via texture-synthesis procedure (e.g. McDermott & Simoncelli, 2011)

Synthesis of model-matched sounds

Synthesis algorithm used relies on histogram matching (Heeger and Bergen, 1995)



Matches time-average of any instantaneous function of the filter response.
 (matches power, but also other moments)

$$\int g(u_k(t, s_{mm}))dt = \int g(u_k(t, s_{nat}))dt \quad \forall k$$

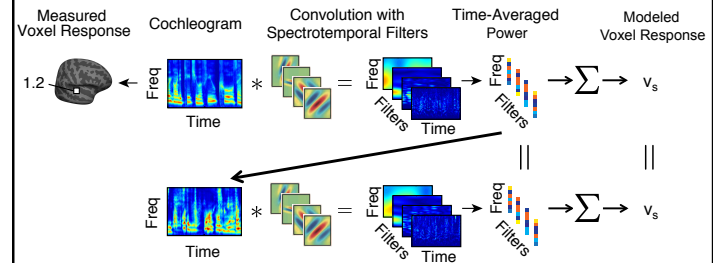
Synthesis of model-matched sounds

Synthesis algorithm used relies on histogram matching (Heeger and Bergen, 1995)

Matches time-average of any instantaneous function of the filter response.
(matches power, but also other moments)

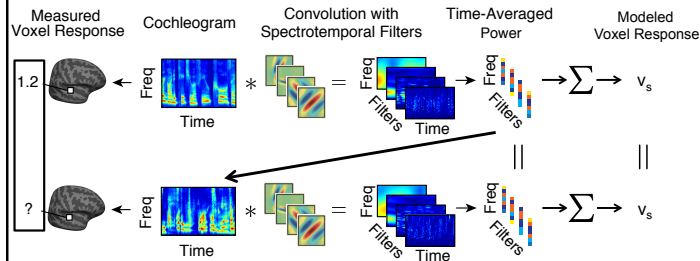
$$\int g(u_k(t, s_{mm}))dt = \int g(u_k(t, s_{nat}))dt \quad \forall k$$

Model-Matched Synthetic Stimuli



Natural and synthetic sound by definition give the same model response.


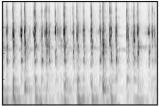
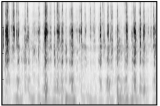

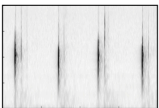
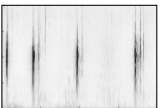

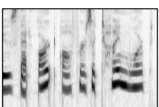
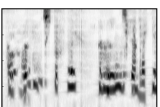

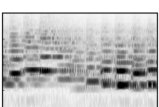
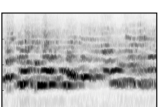
Model-Matched Synthetic Stimuli

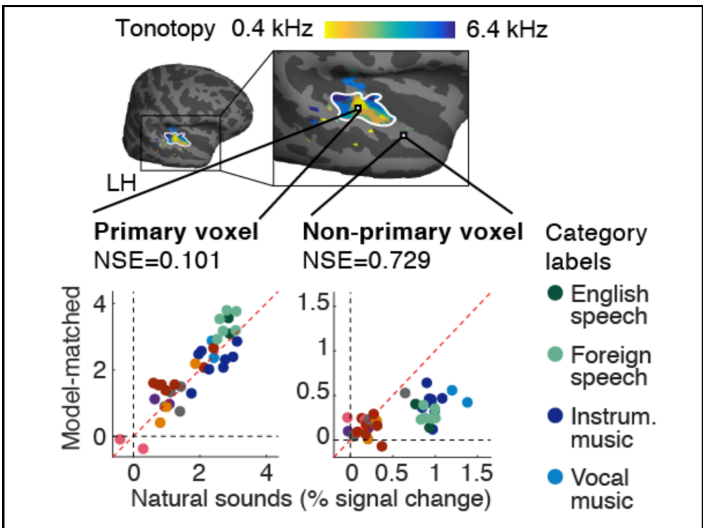
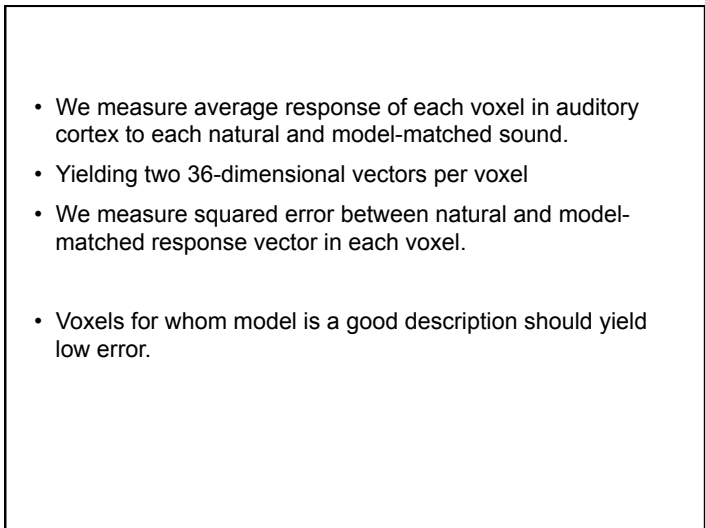
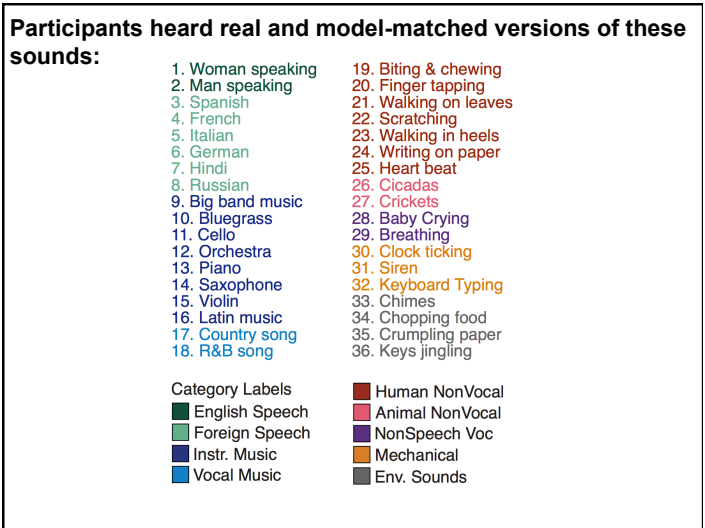
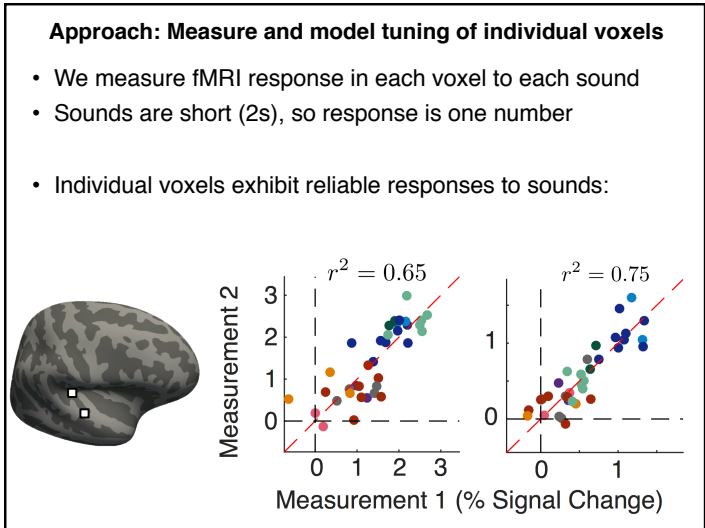


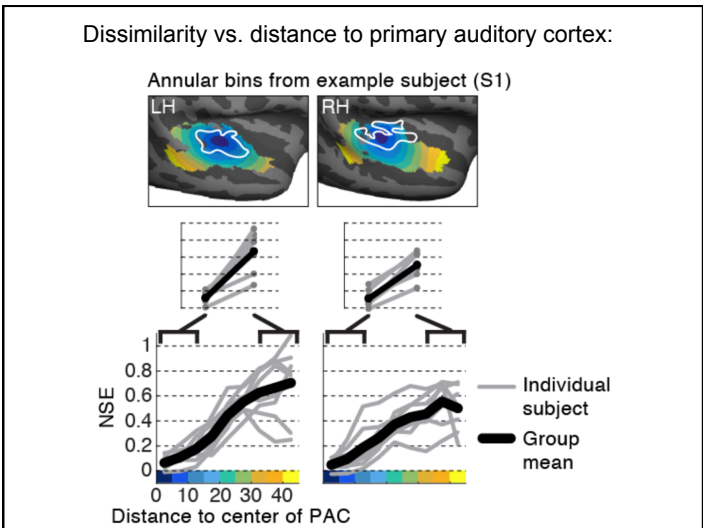
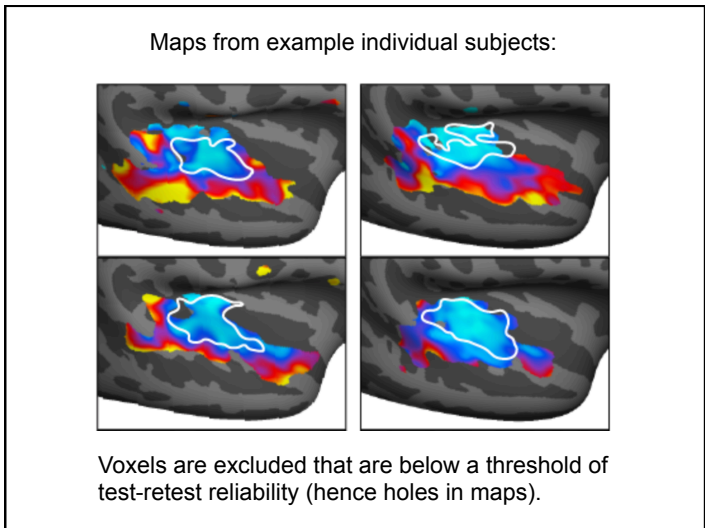
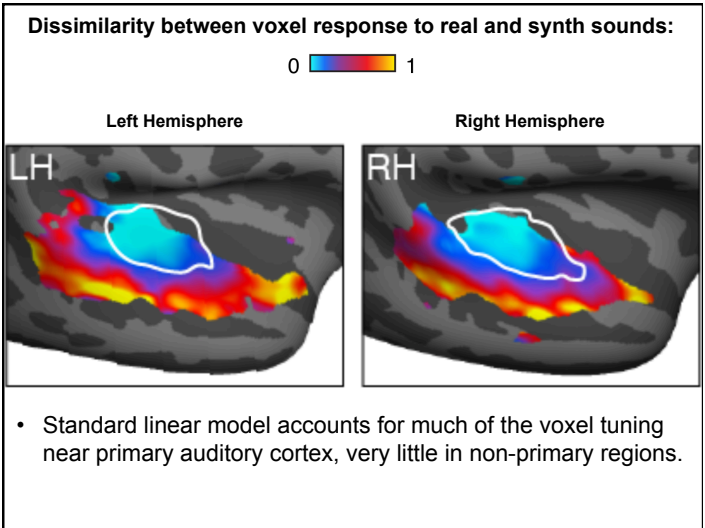
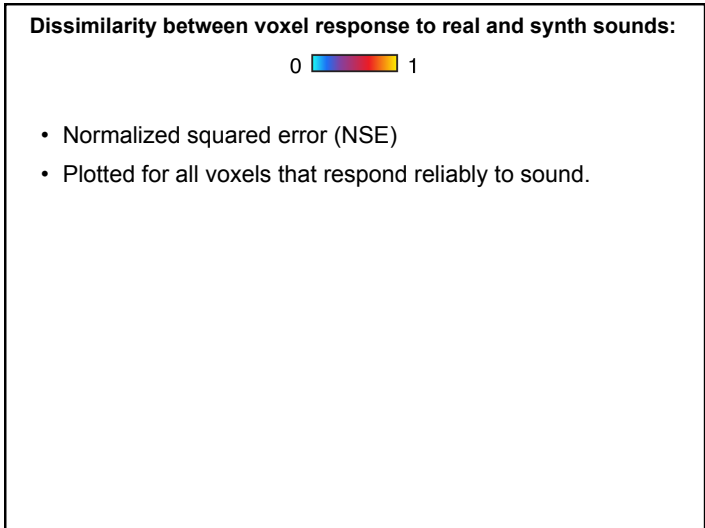
Natural and synthetic sound by definition give the same model response

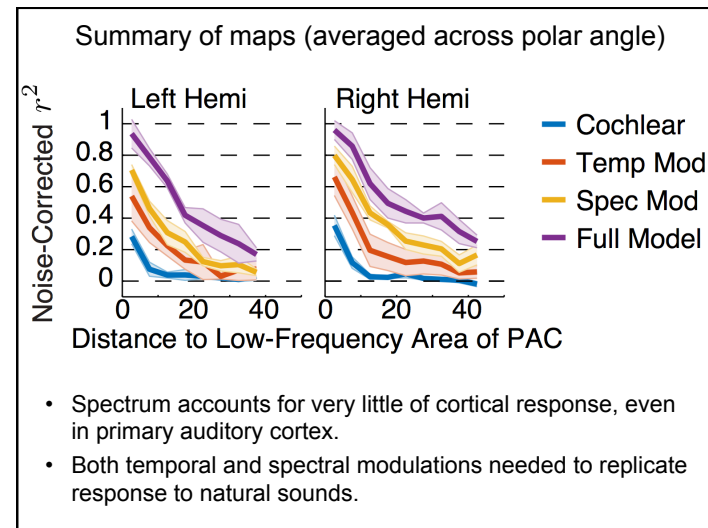
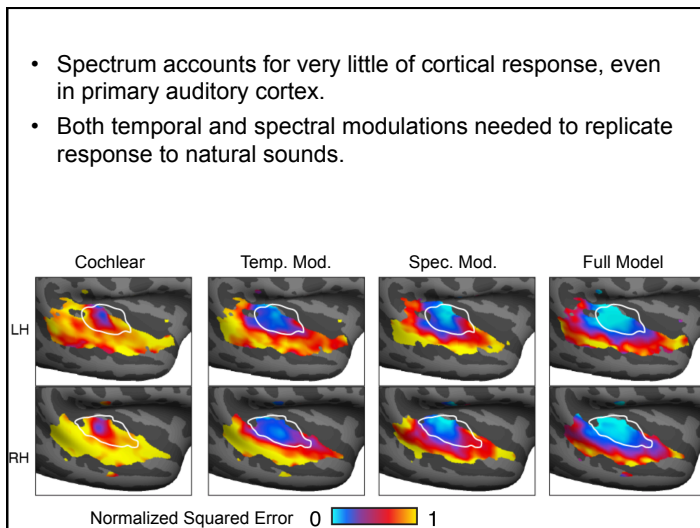
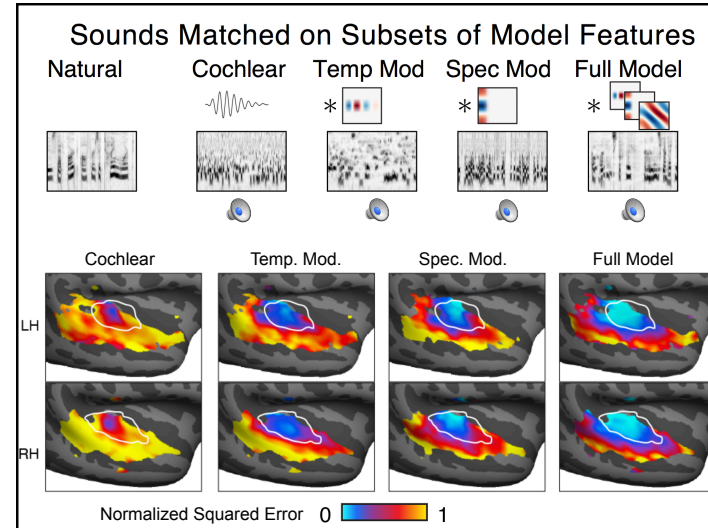
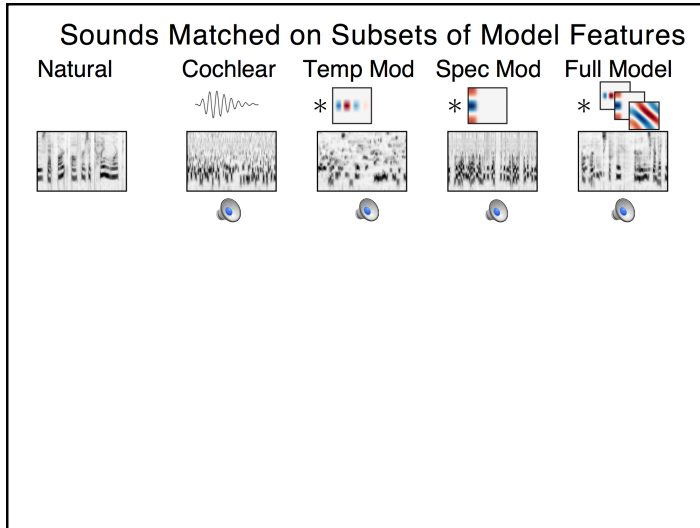
Similar neural response?

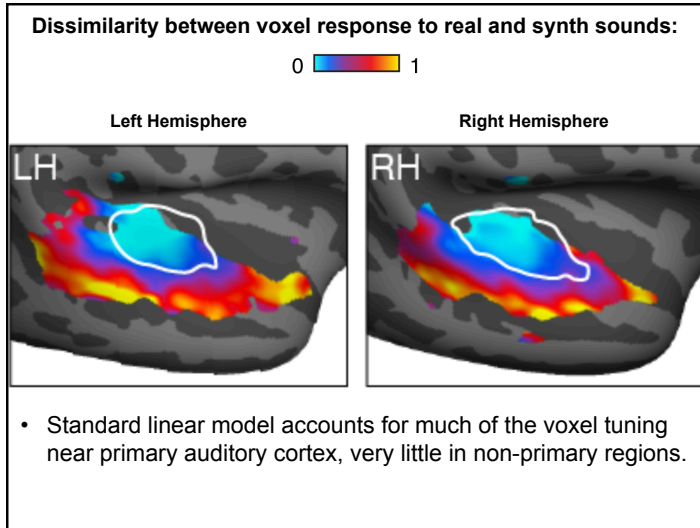
⇒ If model is accurate, neural response should also be matched

	Natural Sound	Model-Matched Sound
 Keyboard Typing		
 Walking in Heels		
 Woman Speaking		
 Violin		




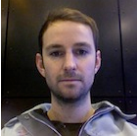








ACKNOWLEDGMENTS

Eero Simoncelli  Michael Schemitsch

Richard McWalter  Jenelle Feather 

Sam Norman-Haignere  Howard Hughes Medical Institute
James S McDonnell Foundation
National Science Foundation