

# Spatial Audio & The Vestibular System

Gordon Wetzstein  
Stanford University

EE 267 Virtual Reality

Lecture 13

[stanford.edu/class/ee267/](https://stanford.edu/class/ee267/)



# Updates

- lab this Friday will be released as a video
- TAs will be in lab on Friday, but as “extended office hours”

# Overview

- what is sound? how do we synthesize it?
- the human auditory system
- stereophonic sound
- spatial audio of point sound sources
- surround sound
- ambisonics
- brief overview of the vestibular system

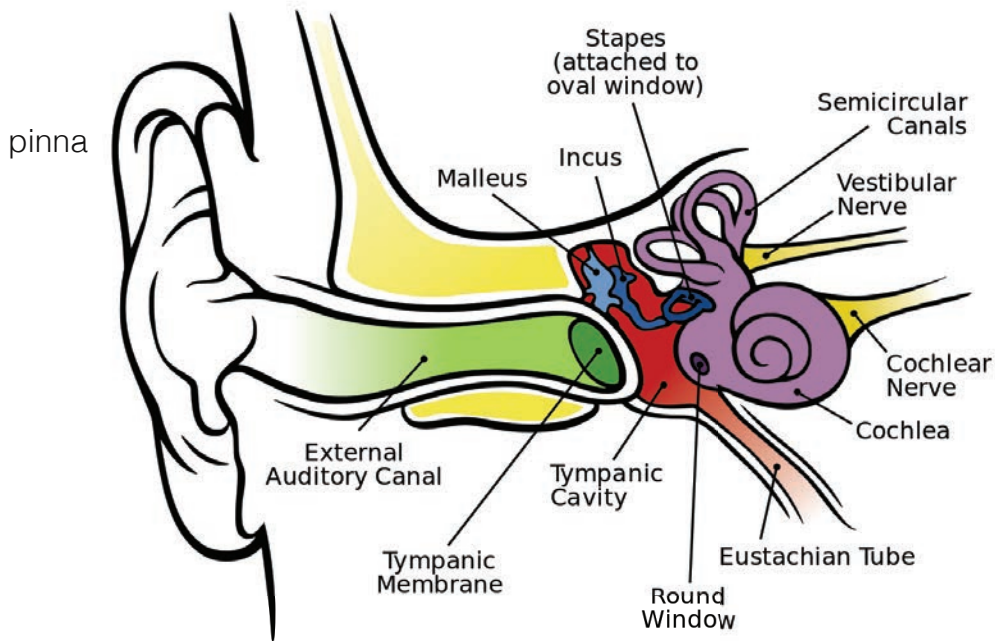
# What is Sound?

- “sound” is a pressure wave propagating in a medium
- speed of sound is  $c = \sqrt{K/\rho}$  where  $c$  is velocity,  $\rho$  is density of medium and  $K$  is elastic bulk modulus
- in air, speed of sound is 340 m/s
- in water, speed of sound is 1,483 m/s

# How do we Synthesize Sound?



# The Human Auditory System

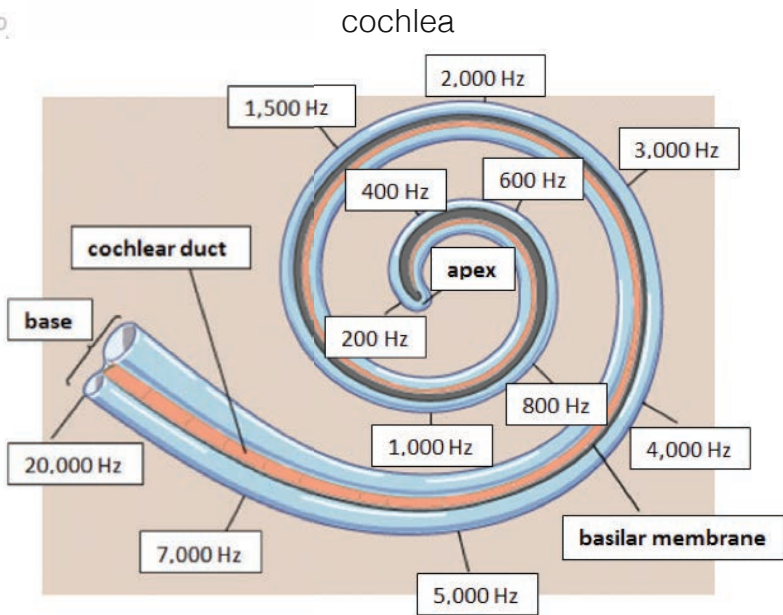
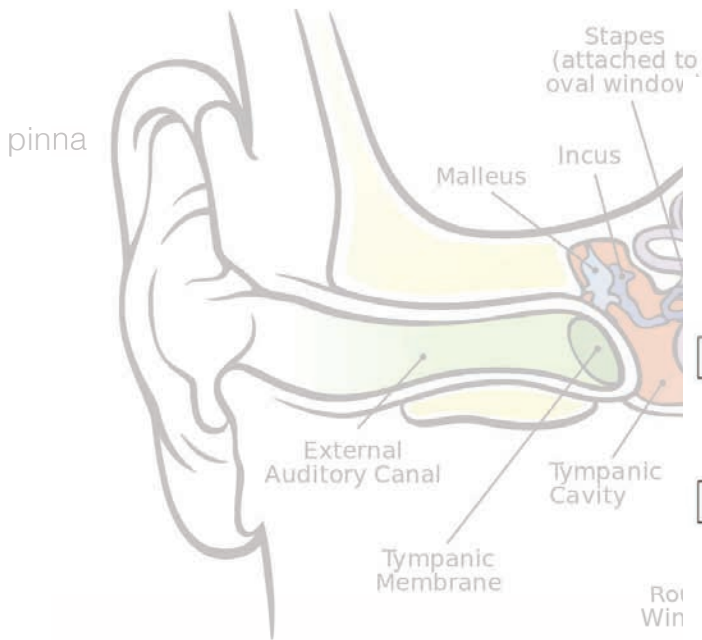


Primary auditory cortex



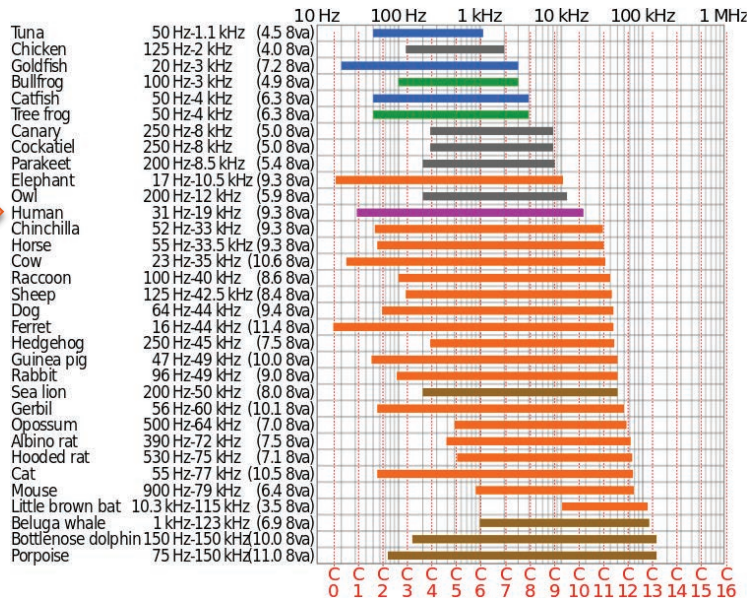
# The Human Auditory System

- hair receptor cells pick up vibrations



# The Human Auditory System

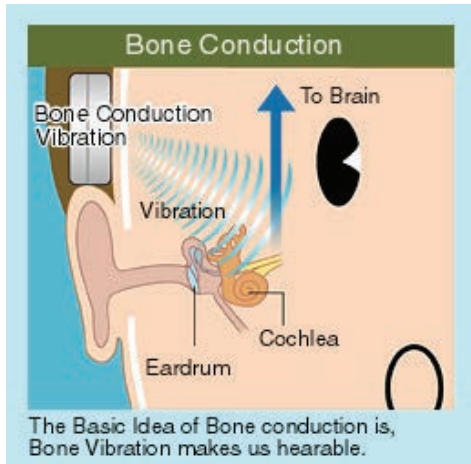
- human hearing range:  
~20 – 20,000 Hz
- variation between  
individuals and  
changes with age





# Bone Conduction

- can stimulate eardrum mechanically to create the illusion of audio, e.g. with bone conduction



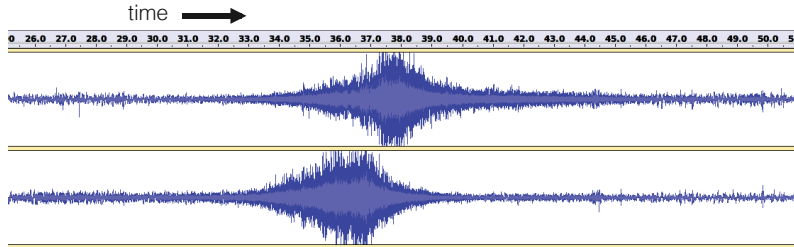
<http://www.goldendance.co.jp/English/boneconduct/01.html>



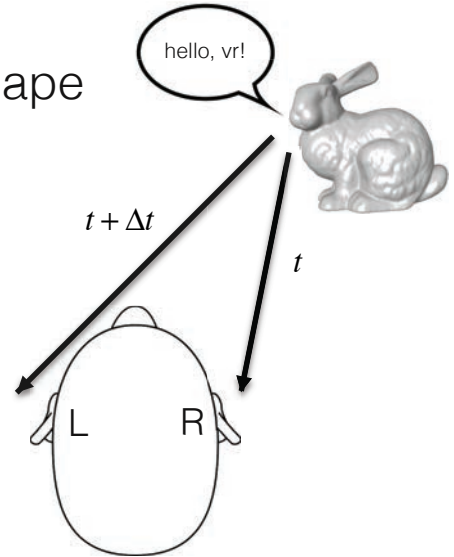
the verge

# Stereophonic Sound

- mainly captures differences between the ears:
  - interaural time difference
  - amplitude differences from body shape (nose, head, neck, shoulders, ...)

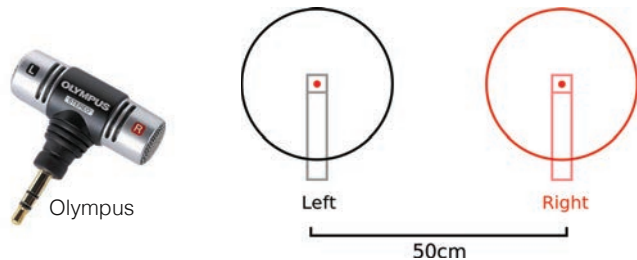


wikipedia



# Stereophonic Sound Recording

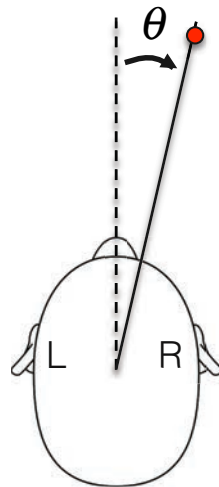
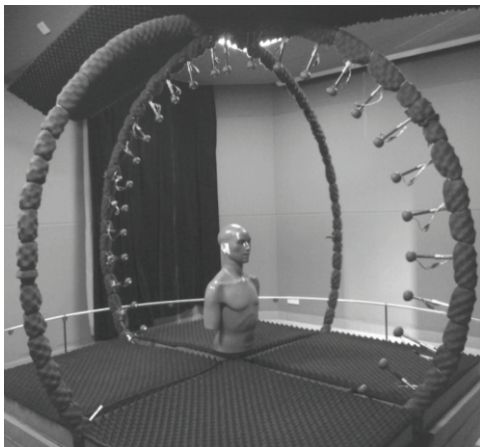
- use two microphones
- A-B techniques captures differences in time-of-arrival
- other configurations work too, capture differences in amplitude



X-Y technique

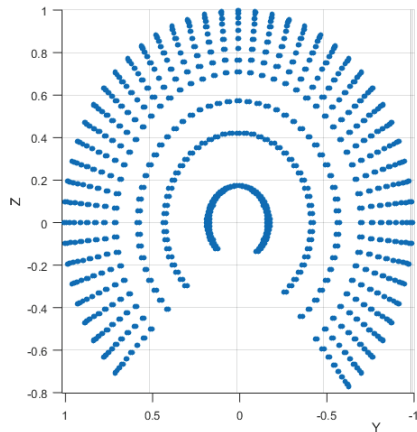
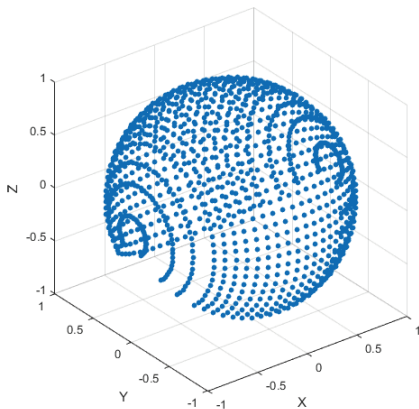
# Head-related Impulse Response (HRIR)

- models phase and amplitude differences for all possible sound directions parameterized by azimuth  $\theta$  and elevation  $\phi$
- can be measured with two microphones in ears of mannequin & speakers all around



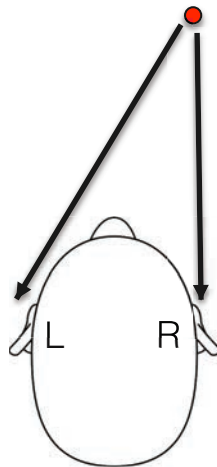
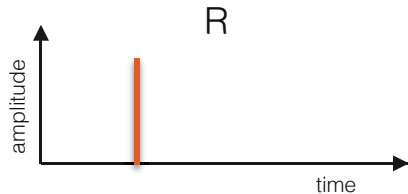
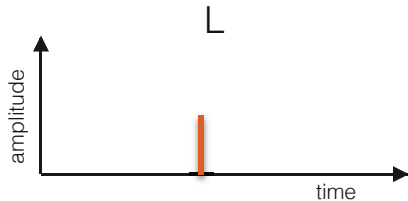
# Head-related Impulse Response (HRIR)

- CIPIC HRTF database: <http://interface.cipic.ucdavis.edu/sound/hrtf.html>
- elevation:  $-45^{\circ}$  to  $230.625^{\circ}$ , azimuth:  $-80^{\circ}$  to  $80^{\circ}$
- need to interpolate between discretely sampled directions



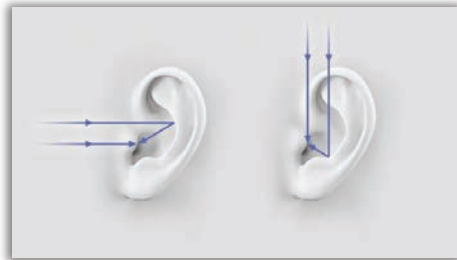
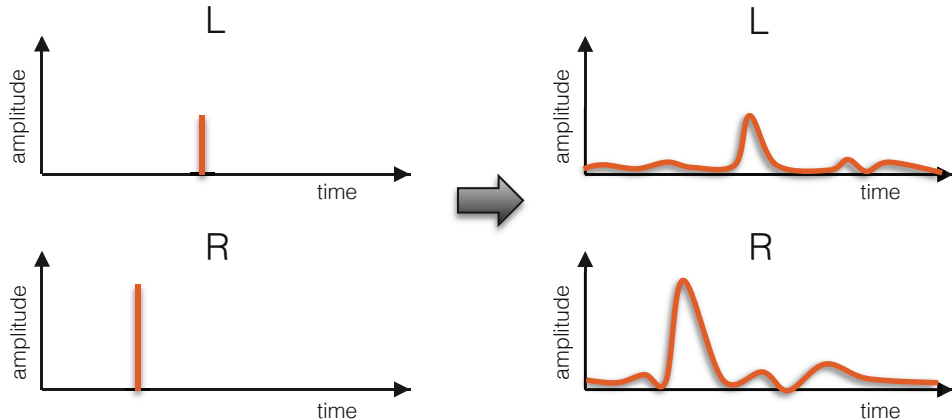
# Head-related Impulse Response (HRIR)

- measuring the HRIR
  - ideal case: scaled & shifted Dirac peaks



# Head-related Impulse Response (HRIR)

- measuring the HRIR
  - ideal case: scaled & shifted Dirac peaks
  - in practice: more complicated, includes scattering in the ear, shoulders etc.



# Head-related Impulse Response (HRIR)

- measuring the HRIR
  - need one temporally-varying function for each angle
  - total of  $2 \cdot N_\theta \cdot N_\phi \cdot N_t$  samples, where  $N_{\theta,\phi,t}$  is the number of samples for azimuth, elevation, and time, respectively

$$hrir\_l(\theta, \phi, t)$$

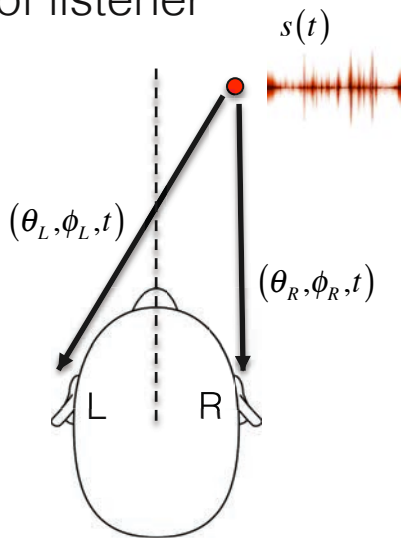
$$hrir\_r(\theta, \phi, t)$$



# Head-related Impulse Response (HRIR)

applying the HRIR:

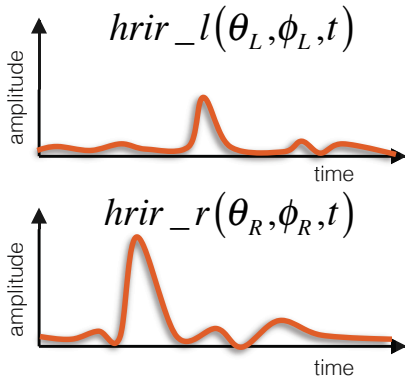
- given a mono sound source  $s(t)$  and it's 3D position
1. compute  $(\theta_L, \phi_L)$  and  $(\theta_R, \phi_R)$  relative to center of listener



# Head-related Impulse Response (HRIR)

applying the HRIR:

- given a mono sound source  $s(t)$  and it's 3D position
  1. compute  $(\theta_L, \phi_L)$  and  $(\theta_R, \phi_R)$  relative to center of listener
  2. look up measured HRIR for left and right ear at these angles



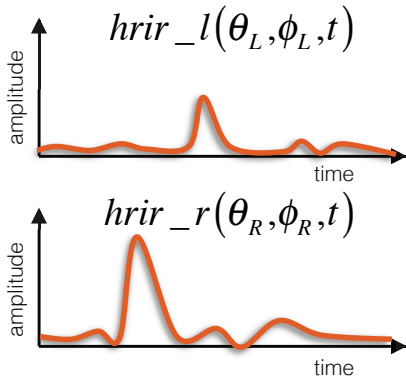
# Head-related Impulse Response (HRIR)

applying the HRIR:

- given a mono sound source  $s(t)$  and it's 3D position
  1. compute  $(\theta_L, \phi_L)$  and  $(\theta_R, \phi_R)$  relative to center of listener
  2. look up measured HRIR for left and right ear at these angles
  3. convolve signal with HRIRs to get response for each ear as

$$s_L(t) = hrir\_l(\theta_L, \phi_L, t) * s(t)$$

$$s_R(t) = hrir\_r(\theta_R, \phi_R, t) * s(t)$$



# Head-related Transfer Function (HRTF)

- HRTF is Fourier transform of HRIR! (you'll find the term HRTF more often than HRIR)

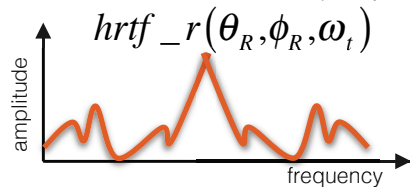
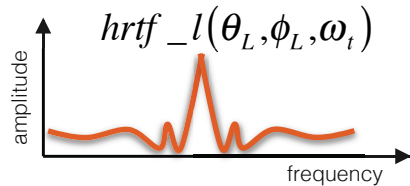
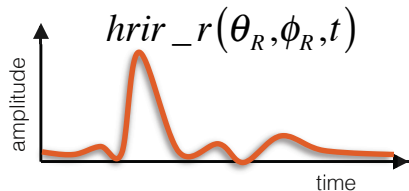
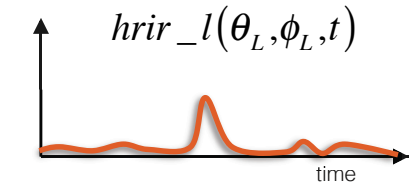
$$s_L(t) = \text{hrir}_l(\theta_L, \phi_L, t) * s(t)$$

$$s_R(t) = \text{hrir}_r(\theta_R, \phi_R, t) * s(t)$$



$$s_L(t) = F^{-1} \{ \text{hrtf}_l(\theta_L, \phi_L, \omega_t) \cdot F \{ s(t) \} \}$$

$$s_R(t) = F^{-1} \{ \text{hrtf}_r(\theta_R, \phi_R, \omega_t) \cdot F \{ s(t) \} \}$$



# Head-related Transfer Function (HRTF)

- HRTF is Fourier transform of HRIR! (you'll find the term HRTF more often than HRIR)

$$s_L(t) = hrir\_l(\theta_L, \phi_L, t) * s(t)$$

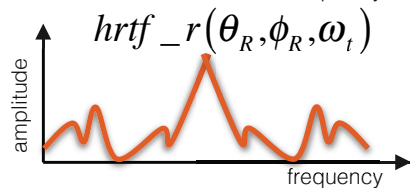
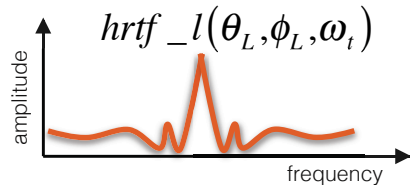
$$s_R(t) = hrir\_r(\theta_R, \phi_R, t) * s(t)$$



convolution theorem

$$s_L(t) = F^{-1}\{hrtf\_l(\theta_L, \phi_L, \omega_t) \cdot F\{s(t)\}\}$$

$$s_R(t) = F^{-1}\{hrtf\_r(\theta_R, \phi_R, \omega_t) \cdot F\{s(t)\}\}$$



# Head-related Transfer Function (HRTF)

- HRTF is Fourier transform of HRIR! (you'll find the term HRTF more often than HRIR)

$$s_L(t) = \text{hrir}_l(\theta_L, \phi_L, t) * s(t)$$

$$s_R(t) = \text{hrir}_r(\theta_R, \phi_R, t) * s(t)$$

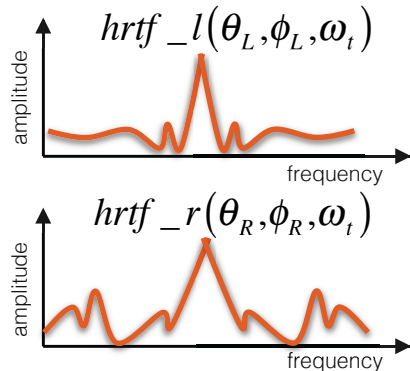


$$s_L(t) = F^{-1} \{ \text{hrtf}_l(\theta_L, \phi_L, \omega_t) \cdot F \{ s(t) \} \}$$

$$s_R(t) = F^{-1} \{ \text{hrtf}_r(\theta_R, \phi_R, \omega_t) \cdot F \{ s(t) \} \}$$

- properties of HRTF:

- complex-valued
- symmetric (because HRIR is real-valued)



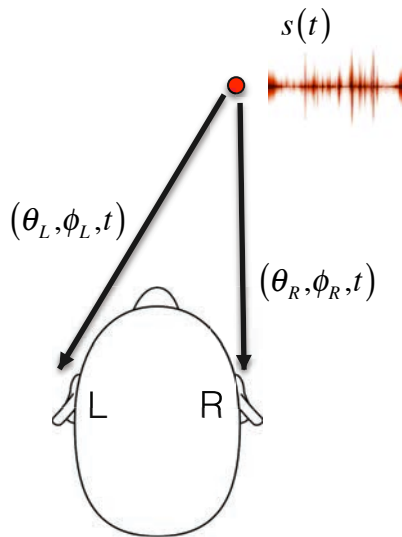
# Head-related Transfer Function (HRTF)

$$s_L(t) = F^{-1} \left\{ hrtf\_l(\theta_L, \phi_L, \omega_t) \cdot F \{ s(t) \} \right\}$$

$$s_R(t) = F^{-1} \left\{ hrtf\_r(\theta_R, \phi_R, \omega_t) \cdot F \{ s(t) \} \right\}$$

# Spatial Sound of 1 Point Sound Source

- given  $s(t)$  and 3D position, follow instructions from last slides by convolving Fourier transform of  $s$  with HRTFs for each ear



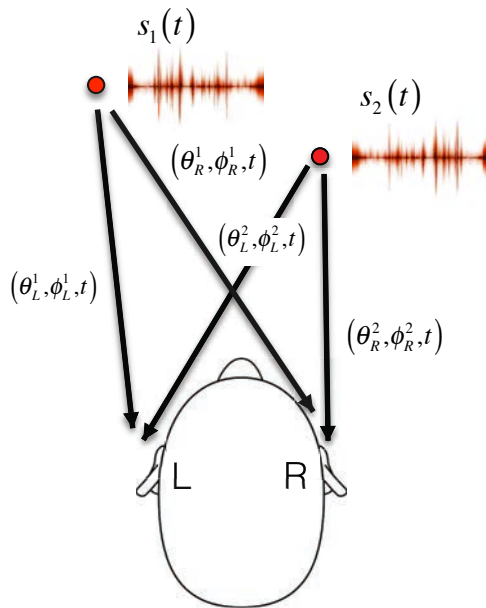


# Spatial Sound of N Point Sound Sources

- superposition principle holds, so just sum the contributions of each

$$s_L(t) = \sum_{i=1}^N F^{-1} \left\{ hrtf\_l(\theta_L^i, \phi_L^i, \omega_t) \cdot F \{ s_i(t) \} \right\}$$

$$s_R(t) = \sum_{i=1}^N F^{-1} \left\{ hrtf\_r(\theta_R^i, \phi_R^i, \omega_t) \cdot F \{ s_i(t) \} \right\}$$



# Surround Sound

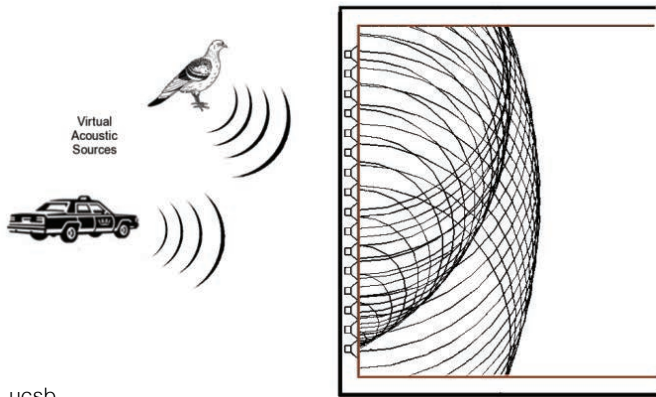
- approximate continuous wave field with discrete set of speakers



- most common:  
5.1 surround sound =  
5 (channels) . 1 (bass)  
→ 6 channels total

# Surround Sound

- approximate continuous wave field with discrete set of speakers
- can also use more speakers for “wave field synthesis” (i.e. audio hologram)



# Surround Sound

- approximate continuous wave field with discrete set of speakers
- can also use more speakers for “wave field synthesis” (i.e. audio hologram)
- for wave field synthesis, phase of speakers needs to be synchronized, i.e. a phased array!

# Surround Sound & HRTF

- for all speaker-based (surround) sound, we don't need an HRTF because the ears of the listener will apply them!
- speaker setup usually needs to be calibrated

# Spatial Audio for VR

- VR/AR requires us to re-think audio, especially spatial audio!
- could use 5.1 surround sound and set up “virtual speakers” in the virtual environment – can use existing content, but not super easy to capture new content; also doesn’t capture directionality from above/below

# Spatial Audio for VR

Two primary approaches:

## 1. Real-time sound engine

- render 3D sound sources via HRTF in real-time, just as discussed in the previous slides
- used for games and synthetic virtual environments
- a lot of libraries available: FMOD, OpenAL, ...

# Spatial Audio for VR

Two primary approaches:

## 2. Spatial sound recorded from real environments

- most widely used format now: ambisonics
- simple microphones exist
- relatively easy mathematical model
- only need 4 channels for starters
- used in YouTube VR and many other platforms



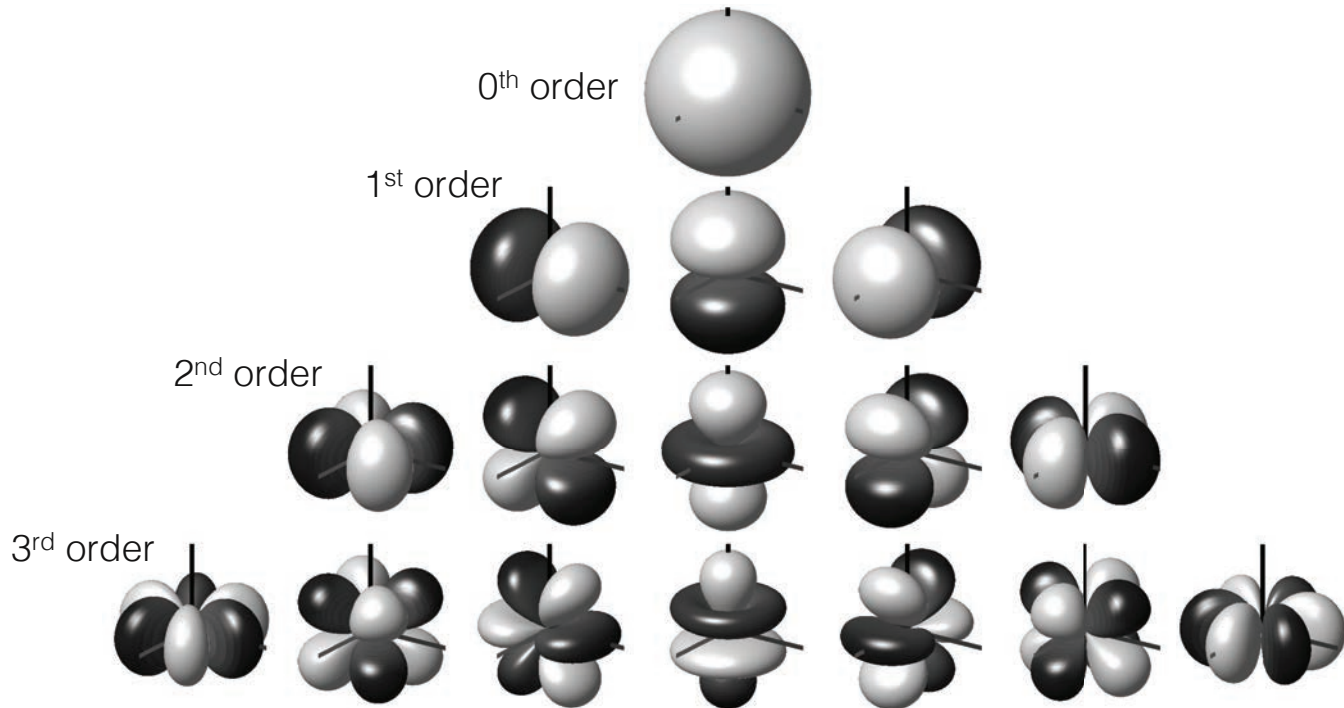
# Ambisonics

- idea: represent sound incident at a point (i.e. the listener) with some directional information
- using all angles  $\theta, \phi$  is impractical – need too many sound channels (one for each direction)
- some lower-frequency (in direction) components may be sufficient → directional basis representation to the rescue!

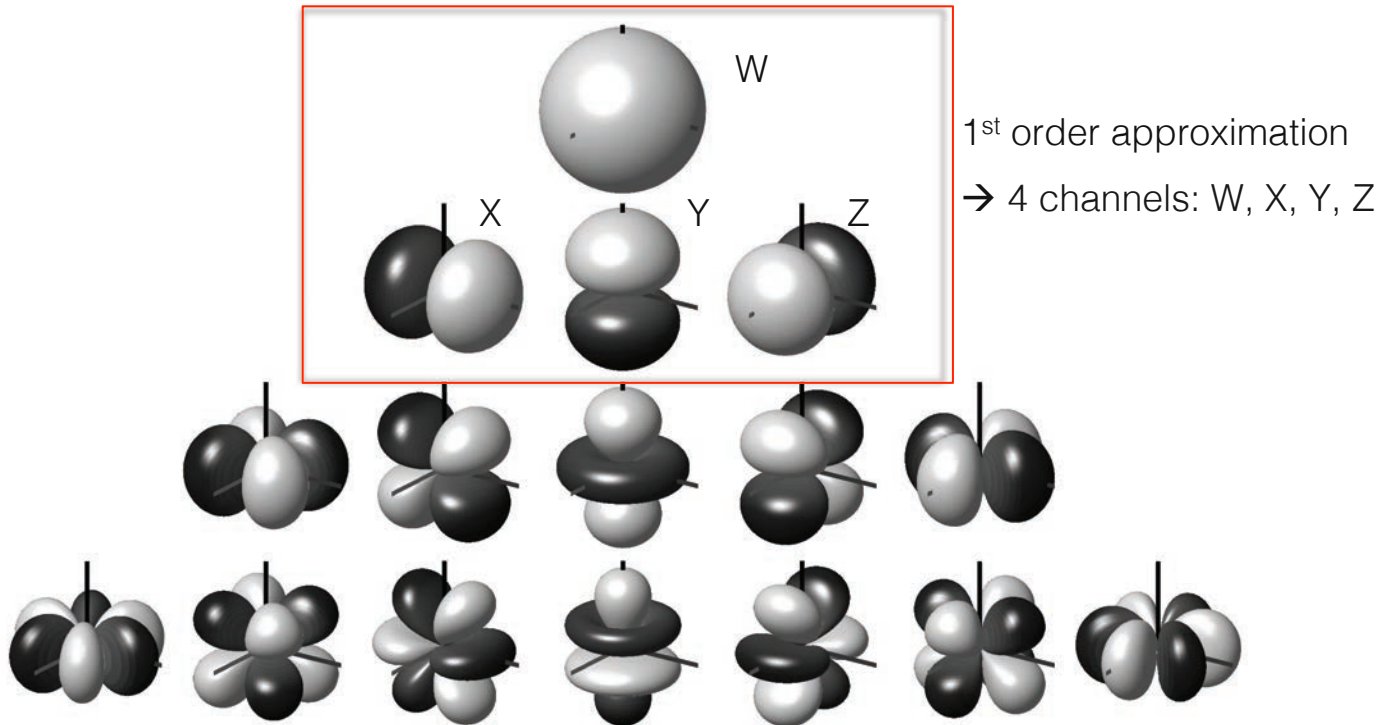
# Ambisonics – Spherical Harmonics

- use spherical harmonics!
- orthogonal basis functions on a sphere, i.e. full-sphere surround sound
- think Fourier transform acting on the directions of a sphere

# Ambisonics – Spherical Harmonics



# Ambisonics – Spherical Harmonics



# Ambisonics – Spherical Harmonics

- can easily convert a point sound source to the 4-channel ambisonics representation
- given azimuth and elevation  $\theta, \phi$ , compute W,X,Y,Z as

$$W = S \cdot \frac{1}{\sqrt{2}}$$



omnidirectional component (angle-independent)

$$X = S \cdot \cos \theta \cos \phi$$



“stereo in x”

$$Y = S \cdot \sin \theta \cos \phi$$



“stereo in y”

$$Z = S \cdot \sin \phi$$



“stereo in z”

# Ambisonics – Spherical Harmonics

- can also record 4-channel ambisonics via special microphone
- same format supported by YouTube VR and other platforms



# Ambisonics – Spherical Harmonics

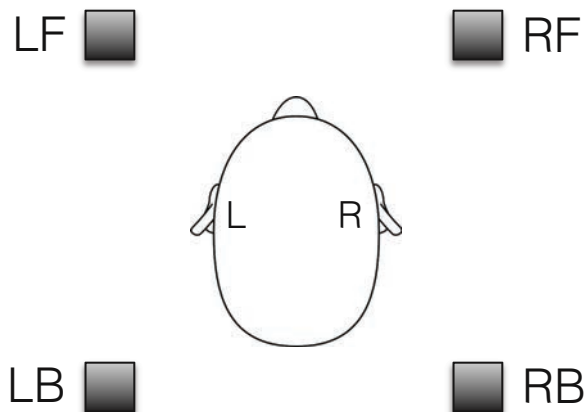
- easiest way to render ambisonics: convert W,X,Y,Z channels into 4 virtual speaker positions
- for a regularly-spaced square setup, this results in

$$LF = (2W + X + Y)\sqrt{8}$$

$$LB = (2W - X + Y)\sqrt{8}$$

$$RF = (2W + X - Y)\sqrt{8}$$

$$RB = (2W - X - Y)\sqrt{8}$$

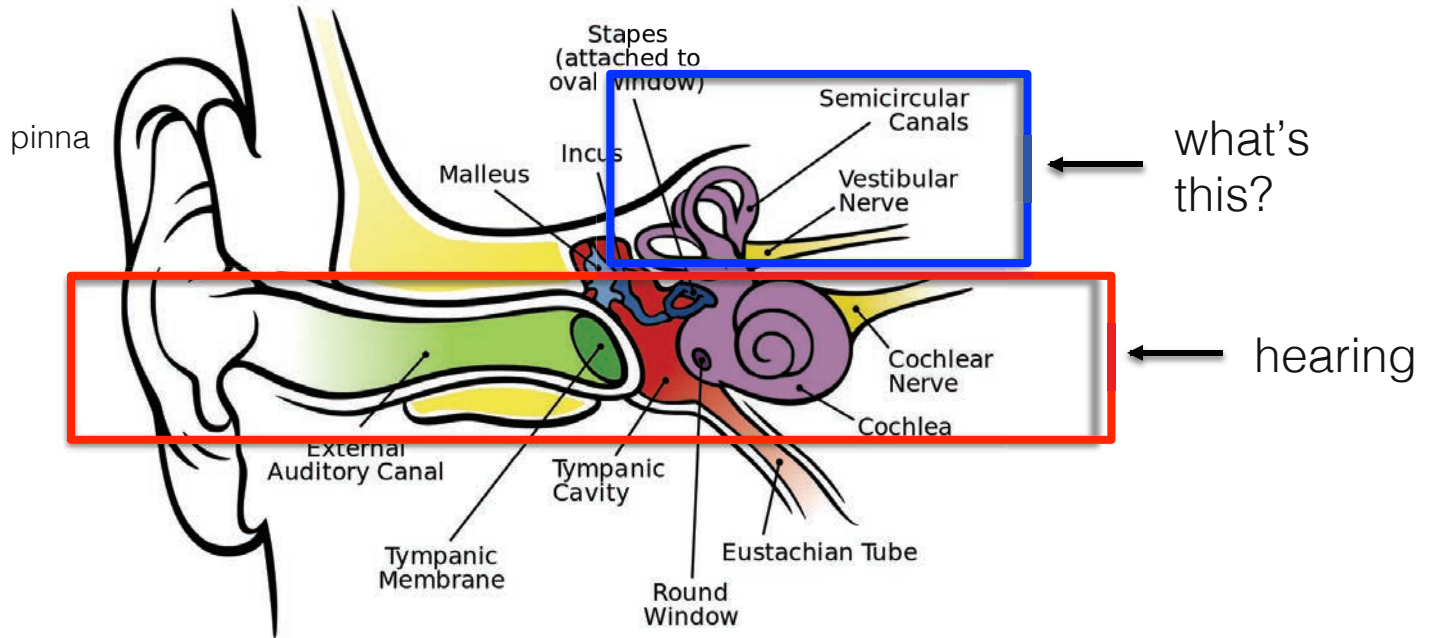


Audio perception happens mostly in the inner ear

What else is happening there?



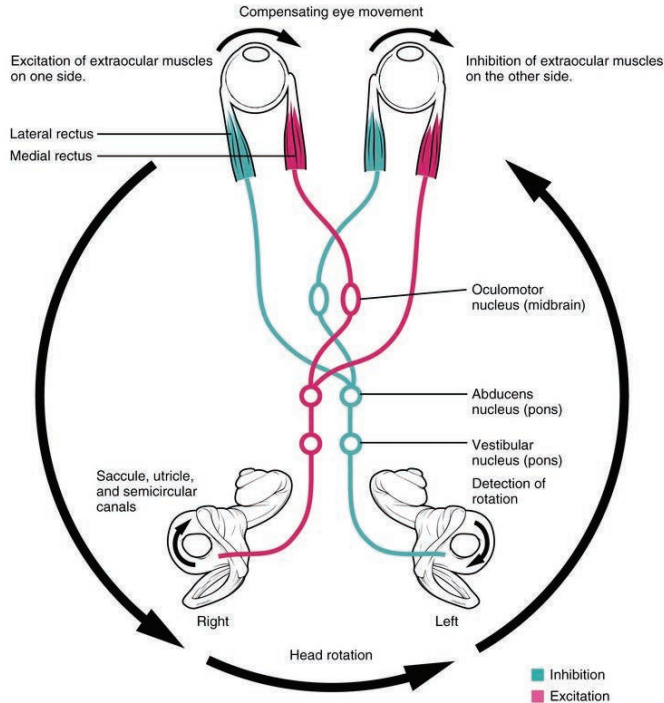
# The Inner Ear



# Brief Overview of the Vestibular System

- provides sense of balance & gravity
- like IMUs – one in each ear!
- in each ear, sense linear (3 dof from otolithic organs) and angular (3 dof from 3 semicircular canals) acceleration via hair cells

# Vestibulo-Ocular Reflex (VOR)



- vestibular system and ocular system are directly coupled in a feedback system
- enables low-latency “optical image stabilization” of the visual system with head motion

# Motion Sickness

3 types of motion sickness (all related to visual-vestibular conflict theory):

1. Motion sickness caused by motion that is felt but not seen
2. Motion sickness caused by motion that is seen but not felt
3. Motion sickness caused when both systems detect motion but they do not correspond.

# Motion Sickness

3 types of motion sickness (all related to visual-vestibular conflict theory):

1. Motion sickness caused by motion that is felt but not seen
2. Motion sickness caused by motion that is seen but not felt
3. Motion sickness caused when both systems detect motion but they do not correspond.

Example: car and sea sickness

# Motion Sickness

3 types of motion sickness (all related to visual-vestibular conflict theory):

1. Motion sickness caused by motion that is felt but not seen
2. Motion sickness caused by motion that is seen but not felt
3. Motion sickness caused when both systems detect motion but they do not correspond.

Example: VR sickness or visually-induced motion sickness (VIMS)

# Motion Sickness

3 types of motion sickness (all related to visual-vestibular conflict theory):

1. Motion sickness caused by motion that is felt but not seen
2. Motion sickness caused by motion that is seen but not felt
3. Motion sickness caused when both systems detect motion but they do not correspond.

Example: motion in low gravity

# References and Further Reading

- Google's take on spatial audio: <https://developers.google.com/vr/concepts/spatial-audio>

HRTF:

- Algazi, Duda, Thompson, Avendado "The CIPIC HRTF Database", Proc. 2001 IEEE Workshop on Applications of Signal Processing to Audio and Electroacoustics
- download CIPIC HRTF database here: <http://interface.cipic.ucdavis.edu/sound/hrtf.html>

Resources by Google:

- <https://github.com/GoogleChrome/omnitone>
- <https://developers.google.com/vr/concepts/spatial-audio>
- <https://opensource.googleblog.com/2016/07/omnitone-spatial-audio-on-web.html>
- <http://googlechrome.github.io/omnitone/#home>
- <https://github.com/google/spatial-media/>