

Speech & Audio Processing & Recognition

<http://www.ee.columbia.edu/~dpwe/courses/e6820-2001-01/>

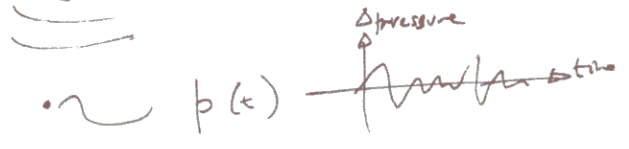
What's it about?

SOUND

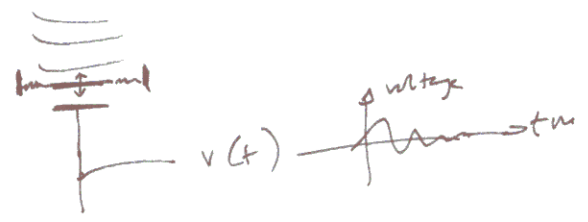
Mechanical vibration
 ↓
 air molecules
 ↓
 radiating wave (Δ pressure)



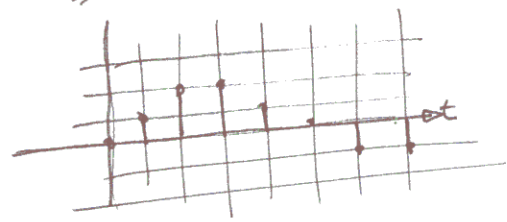
SOUND: pressure varies at a point in space



↓
 push against surface
 (diaphragm of mic)
 ↓
 voltage $v(t)$
 ↓
 processing



DIGITAL: representing sound on a computer → series of numbers
 ⇒ discrete time, discrete value



→ DSP

$$\hat{p}[n] = Q[p(nT)]$$

Annotations: 'quantize' points to the Q operator; 'sample interval' points to the nT term.

+ can convert back to sound:

Speaker passes current through magnetic field
 → force
 → moves "cone"
 → pushes air
 → radiates wave ...



Why sound?

because we hear

①

Why hearing? because it's useful (evolutionary advantage)

- info at a distance
- consistent/reliable
- correlates with physical force
- complements with

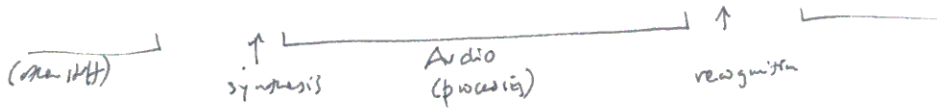
footsteps ex →

- common in higher animals
- "what" info? depends on what we are about → 4

Speech & Audio Processing & Recognition

is founded on Acoustic Communication:

eg. source-channel model:



What is Sound?

• air vibration → small changes in pressure
 hence, mechanical motion



pushes air molecules
 → radiating wave

pushes eg. diaphragm

→ electrical signal or pressure $p(t)$
 Computers • sample $p(t)$ at discrete times
 • quantize to finite levels → $p[n] = \alpha(p(nT))$
 Generate sound by mixing speaker or voltage

~~Audio is interesting because we hear~~

We hear because it's useful

ie. evolutionary advantage

- information at a distance
- relatively consistent/reliable
- correlated with physical force
- complementary to vision.
- universal?

⇒ what we hear reflects what is useful

ECOLOGICAL constraints.

What is the scope of audio processing?

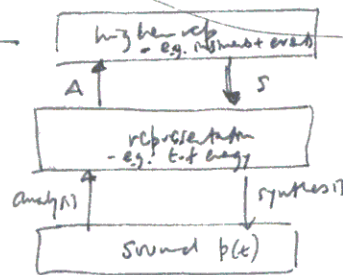


(inf) subset of audio
 AUDIO
 music
 - sigs
 - info
 - mix
 env.
 - filters
 - eqs
 - mod
 synth
 - change
 - SFX

Sp Synth	GSM bc	Vx mod	ASR
• Synthesis • phys. models	WBAC	• reverb • time-scale mod • SFX	• music ID • transcrip
• filter • models		•	• robust • Audio CBR
• sound design			• also detect.

Synthesis compression modification recognition
 PROCESSING
 subset. classif
understand

Representation

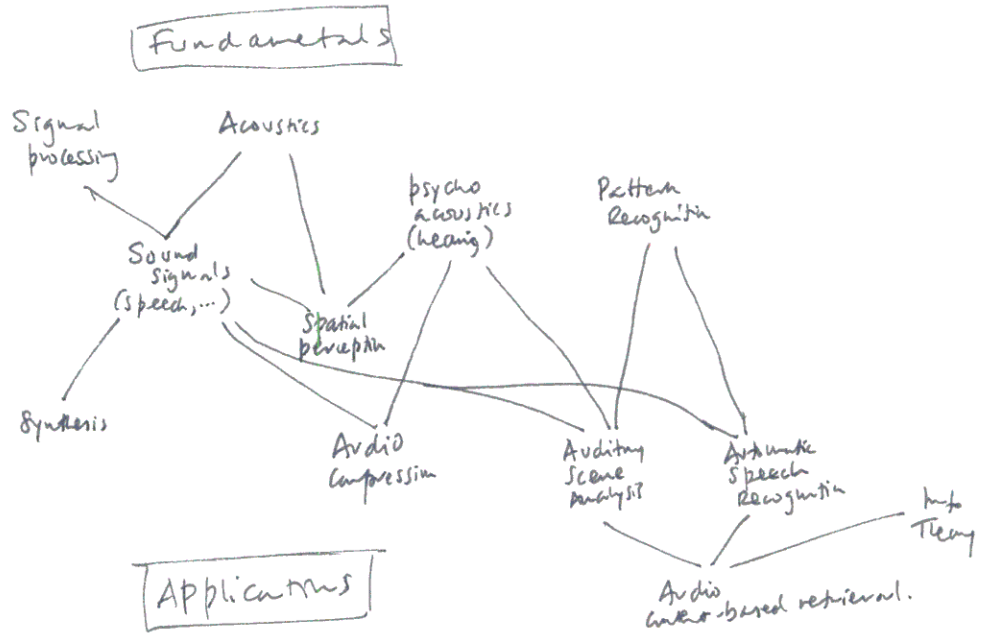


General info proc model

examples: speech vocog
 music synth

- Representations:
- permit modifications (pitch + form)
 - make something useful explicit (spectrogram)
 - efficient (compress)

Course overview



Goals

- ① Audio Proc - tools & techniques
- ② Advanced DSP - example applications
 - specific techniques } apply beyond audio
- ③ Practical experience - computer examples
 - reading papers
 - ~~projects~~ projects & HW

Structure

- ① lecture Class - (lecture)
 - demonstrations
 - discussions (papers, ideas) } student presentation.
 - project ITH
- ② Homework - code
 - problems, puzzles - Matlab
 - computer paper summaries
 - build up a web site
 + exams - fundamentals
 - qualitative
- ③ Project - any specific aspect (some ideas...)
 • lit vsch
 • PRACTICAL - Matlab (eg.)
 • several passes
 • evaluation
 • presentation (midterm?)
 - start early - ideas now!

- HW1 - G+M - into
- Matlab:
 - variable
 - plot
 - fine of loop
 - initial prog. length
 - Sum File

Course overview - summarize some of the key points
 → see where we're going
 → start thinking

- Slidepack
- index

① Speech signal

Spectrogram

- waveform - timescale - sounds
- spectrogram
- t/f tradeoff
- computer
- t-f energy of ear

t vs. f tradeoff.

RT Sgram
 SNACK?

- variation in time
- phonemes - stops, vowels, fricatives, glides
- linguistic variation.

- Sgram pair
- segmentation

+ waveform

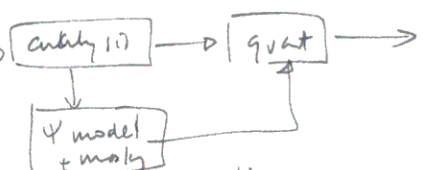
- limit examples → reg
- t/f tradeoff

→ ψ

- SWS
- restoration
- vowel seg.
- pulsation

② WBAC

Structure:



- coding: samples = rate + bits
- Q noise
- subband → channel mixing
- ψ masking - ~~channel~~ noise shaping

dsps

- blk drag
- masking image
- + index - BR
- NMR
- basic mpeg example
- vary bit rates
- look at sgram.

- examples at various levels + sgrams.

③ ASR

- basic blocks
- context models
- basic eq = SEARCH

words → grammar → phons → acous phons

$$W^* = \max_M (p(x|M))$$

- ASR into slides
- + SPRAC demo

- clear example
- --plus added noise

④ Audio CBR

- mixtures
- search (by example by attributes)

- muscle fish

- muscle fish slide
- + demo
- muscle fish demo.