# Computational Perception
## 15-485/785

### February 5, 2008

### Auditory Coding 1

# What are the problems of sensory coding?

- What should the sensor sense?

- How is energy transduced?

- How to deal with noise?

- How to compress dynamic range?

- How to prevent the sensor from being damaged?
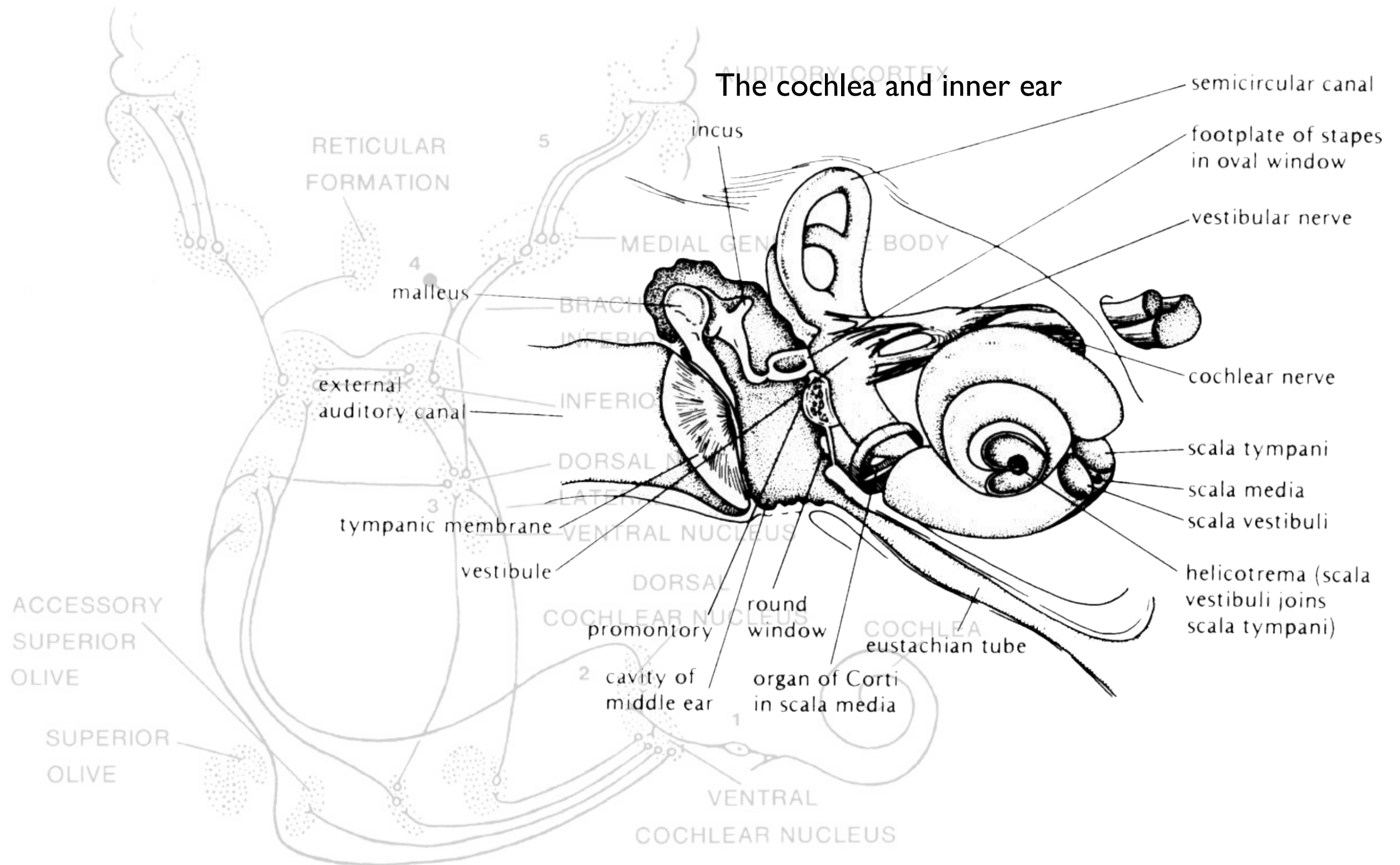
# Two approaches to the study of systems

1. Experimental/behavioral approach:

    - describe and characterize behavior

    - understand range and limitations

    - investigate system properties and organization

    - develop theories to better understand functional roles

2. Theoretical/computational approach:

    - define problem

    - develop models and algorithms

    - understand range and limitations

    - develop more algorithms: more general/specialized; faster/less resources

# The complexity of the auditory system

The cochlea and inner ear

*from Warren, 1999*

# What principles should guide the choice of representation?

Unsupervised approaches:

- find useful "features"

- adapted to the patterns of interest

- useful in a wide range of tasks

Supervised approaches:

- Maximize performance on given task

At low-levels, we have to use unsupervised approaches.

# Linear superposition

Goal is to describe the data to desired precision.
Code signal by linear superposition of basis functions:

$$\mathbf{x} = \vec{a}_1 s_1 + \vec{a}_2 s_2 + \cdots + \vec{a}_L s_L + \vec{\epsilon}$$

$$= \mathbf{A}\mathbf{s} + \boldsymbol{\epsilon}$$

- $x(t)$ is represented by a vector $\mathbf{x}$
- $\vec{a}_i$ are the *basis vectors*
- $\mathbf{A}$ is the *basis* (could be Fourier, wavelet, etc.)
- $s_i$ are the *coefficients*

Can solve for $\hat{\mathbf{s}}$ in the no noise case

$$\hat{\mathbf{s}} = \mathbf{A}^{-1}\mathbf{x}$$

# An information theoretic approach

Want algorithm to choose optimal $\mathbf{A}$ (basis matrix).

Generative model for data is:
$$\mathbf{x} = \mathbf{A}\mathbf{s} + \boldsymbol{\epsilon}$$

Probability of pattern $\mathbf{x}$ given representation $\mathbf{s}$

$$P(\mathbf{x}|\mathbf{A}, \mathbf{s}) \quad \sim \quad f(\mathbf{x} - \mathbf{A}\mathbf{s}, \boldsymbol{\Sigma}, I)$$

# Learning objective

Objective: maximize coding efficiency
  $\Rightarrow$ maximize probability of data ensemble

Probability of pattern ensemble is:

$$P(\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N | \mathbf{A}) = \prod_k P(\mathbf{x}_k | \mathbf{A})$$

# Optimal coding of an acoustic waveform



- We do *not* assume a Fourier or spectral representation.
- Goal:
    *Predict optimal transformation of acoutsic waveform*
    *from statistics of the acoustic environment.*
- Use a simple model: bank of linear filters

# Coding patterns with a statistical model
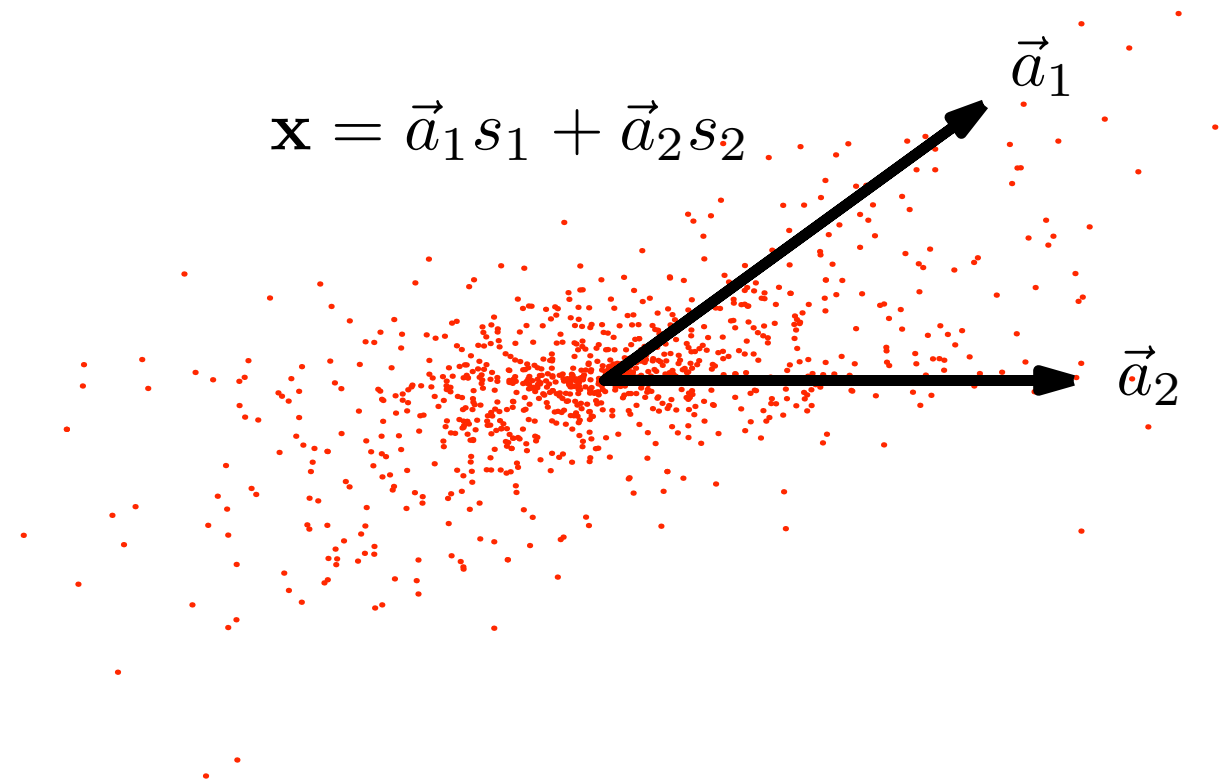
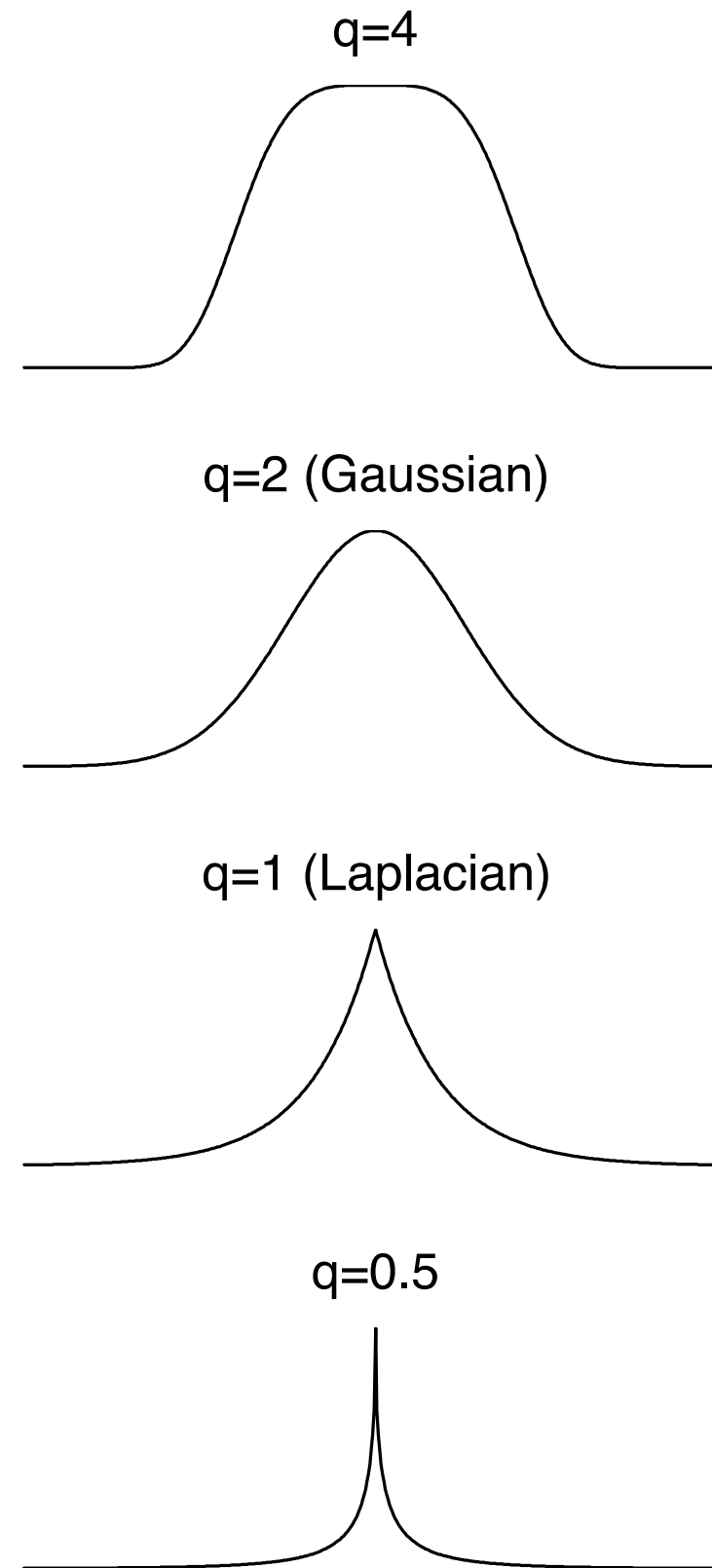*Goal*: Encode the patters to desired precision:

$$\begin{aligned} \mathbf{x} &= \vec{a}_1 s_1 + \cdots + \vec{a}_L s_L + \vec{\epsilon} \\ &= \mathbf{A}\mathbf{s} + \boldsymbol{\epsilon} \end{aligned}$$

$$\mathbf{x} = \vec{a}_1 s_1 + \vec{a}_2 s_2$$

$\vec{a}_1$

$\vec{a}_2$

*Posterior*:

$$P(\mathbf{s}|\mathbf{x}, \mathbf{A}) = \frac{P(\mathbf{s})P(\mathbf{x}|\mathbf{s}, \mathbf{A})}{P(\mathbf{x}|\mathbf{A})}$$

*Prior*: $s_i$'s are *independent* and *sparse*:

$$P(\mathbf{s}) = \prod_i P(s_i)$$

$$P(s_i) \propto \exp\left[-\left|\frac{s_i}{\lambda_i}\right|^{q_i}\right]$$

# Coding patterns with a statistical model

*Goal*: Encode the patters to desired precision:

$$\begin{aligned} \mathbf{x} &= \vec{a}_1 s_1 + \cdots + \vec{a}_L s_L + \vec{\epsilon} \\ &= \mathbf{A}\mathbf{s} + \boldsymbol{\epsilon} \end{aligned}$$

*Posterior*:

$$P(\mathbf{s}|\mathbf{x}, \mathbf{A}) = \frac{P(\mathbf{s})P(\mathbf{x}|\mathbf{s}, \mathbf{A})}{P(\mathbf{x}|\mathbf{A})}$$

*Prior*: $s_i$'s are *independent* and *sparse*:

$$P(\mathbf{s}) = \prod_i P(s_i)$$

$$P(s_i) \propto \exp\left[-\left|\frac{s_i}{\lambda_i}\right|^{q_i}\right]$$

q=4

q=2 (Gaussian)

q=1 (Laplacian)

q=0.5

# Coding patterns with a statistical model

*Goal*: Encode the patters to desired precision:

$$\mathbf{x} = \vec{a}_1 s_1 + \cdots + \vec{a}_L s_L + \vec{\epsilon}$$
$$= \mathbf{A}\mathbf{s} + \boldsymbol{\epsilon}$$

*Posterior*:

$$P(\mathbf{s}|\mathbf{x}, \mathbf{A}) = \frac{P(\mathbf{s})P(\mathbf{x}|\mathbf{s}, \mathbf{A})}{P(\mathbf{x}|\mathbf{A})}$$

*Prior*: $s_i$'s are *independent* and *sparse*:

$$P(\mathbf{s}) = \prod_i P(s_i)$$

$$P(s_i) \propto \exp\left[-\left|\frac{s_i}{\lambda_i}\right|^{q_i}\right]$$

*Likelihood*: Assume $\epsilon \sim$ Gaussian,

$$P(\mathbf{x}|\mathbf{s}, \Sigma) \propto \exp\left[-\frac{1}{2}\boldsymbol{\epsilon}^T \Sigma^{-1} \boldsymbol{\epsilon}\right]$$

*Inference*: use the MAP value:

$$\hat{\mathbf{s}} = \arg\max_{\mathbf{s}} P(\mathbf{s}|\mathbf{x}, \mathbf{A})$$

Simple special case: no noise (ICA)

$$\hat{\mathbf{s}} = \mathbf{A}^{-1}\mathbf{x}$$

Inference (or recognition or coding):

*finds most efficient representation of pattern $\mathbf{x}$ in a given basis $\mathbf{A}$*

# Learning: Optimizing the model parameters

Learning objective:

*maximize coding efficiency*
$\Rightarrow$ maximize $P(\mathbf{x}|\mathbf{A})$ over $\mathbf{A}$.

Probability of pattern ensemble is:

$$P(\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N|\mathbf{A}) = \prod_k P(\mathbf{x}_k|\mathbf{A})$$

$P(\mathbf{x}|\mathbf{A})$ is obtained by marginalization:

$$P(\mathbf{x}|\mathbf{A}) = \int d\mathbf{s}\, P(\mathbf{x}|\mathbf{A}, \mathbf{s})P(\mathbf{s})$$

$$= \frac{P(\mathbf{s})}{|\det \mathbf{A}|}$$

Use *independent component analysis* (ICA) to learn $\mathbf{A}$:

$$\Delta\mathbf{A} \propto \mathbf{A}\mathbf{A}^T \frac{\partial}{\partial \mathbf{A}} \log P(\mathbf{x}|\mathbf{A})$$

$$= -\mathbf{A}(\mathbf{z}\mathbf{s}^T - \mathbf{I}),$$

where $\mathbf{z} = (\log P(\mathbf{s}))'$. Assume generalized Gaussians:

$$P(s_i) \sim \mathcal{N}^{q_i}(s_i|\mu, \sigma).$$

This learning rule:

- *learns the feature set that captures the most structure*
- *optimizes basis to maximize the efficiency of the code*

# Learning the optimal codes

Goal:

*Predict optimal transformation of sound waveform from statistics of the acoustic environment*

Learning procedure:

- random sound segments (8 msec)

- optimize features using ICA

## *What sounds to use?*

What tasks are auditory systems adapted to do?

- localization ⇒ environmental sounds

- communication ⇒ vocalizations

- general sound recognition

Use a variety of sound ensembles:

- non-harmonic *environmental sounds* (e.g. footsteps, stream sounds, etc.)

- *animal vocalizations* (rainforest mammals, e.g. chirps, screeches, cries, etc.)

- *speech* (samples from 100 male & female speakers from the TIMIT corpus)

# Natural sounds

| | environmental sounds | |
|---|---|---|
| vocalizations | transient | ambient |
| fox | walking on leaves | rustling leaves |
| squirrel | cracking branches | stream by waterfall |

# Natural sounds

|  | vocalizations | environmental sounds | |
|---|---|---|---|
|  |  | transient | ambient |
|  | **fox** | walking on leaves | rustling leaves |
|  | squirrel | cracking branches | stream by waterfall |

# Natural sounds

|  | environmental sounds | |
| vocalizations | transient | ambient |
| --- | --- | --- |
| fox | walking on leaves | rustling leaves |
| **squirrel** | cracking branches | stream by waterfall |

# Natural sounds

|  | environmental sounds | |
| vocalizations | transient | ambient |
| --- | --- | --- |
| fox | **walking on leaves** | rustling leaves |
| squirrel | cracking branches | stream by waterfall |

# Natural sounds

| vocalizations | environmental sounds | |
| | transient | ambient |
| --- | --- | --- |
| fox | walking on leaves | rustling leaves |
| squirrel | **cracking branches** | stream by waterfall |

# Natural sounds

| | environmental sounds | |
| vocalizations | transient | ambient |
| --- | --- | --- |
| fox | walking on leaves | **rustling leaves** |
| squirrel | cracking branches | stream by waterfall |

# Natural sounds

|  | environmental sounds | |
| vocalizations | transient | ambient |
| --- | --- | --- |
| fox | walking on leaves | rustling leaves |
| squirrel | cracking branches | **stream by waterfall** |

Coding Cost

$10^0$ $10^1$ $10^2$ $10^3$ $10^4$

Iteration #

# Optimal linear filters for natural sounds



environmental sounds



vocalizations



speech

The optimal code depends on the class of sounds being encoded:

- a wavelet-like transform is best for environmental sounds
- a Fourier-like transform is best for vocalizations
- an intermediate transform is best for speech *or general natural sounds*

Michael S. Lewicki ◇ Carnegie Mellon

# Characterizing the filter population

time-frequency distributions
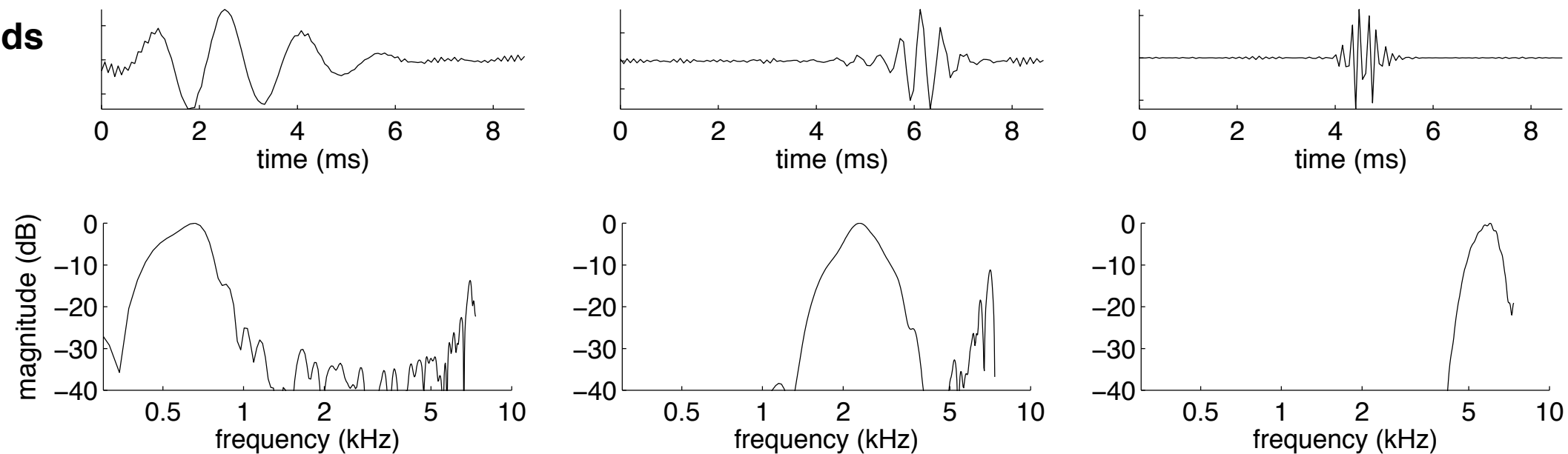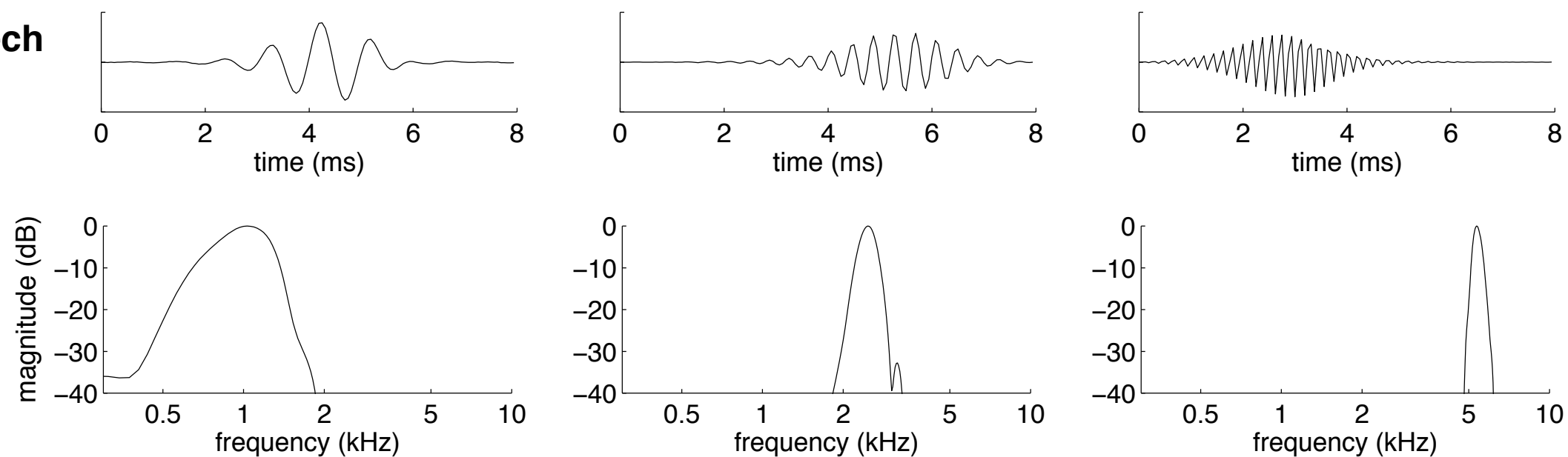
# Schematic time-frequency distributions



Fourier

typical wavelet

frequency

time

time

# Comparison to cat auditory nerve data



Data

Theory

Filter sharpness:

$$Q_{10\text{dB}} = f_c / w_{10\text{dB}}$$

+ vocalizations
o speech/combined
x environmental sounds

# Next time:
*non-linear coding*