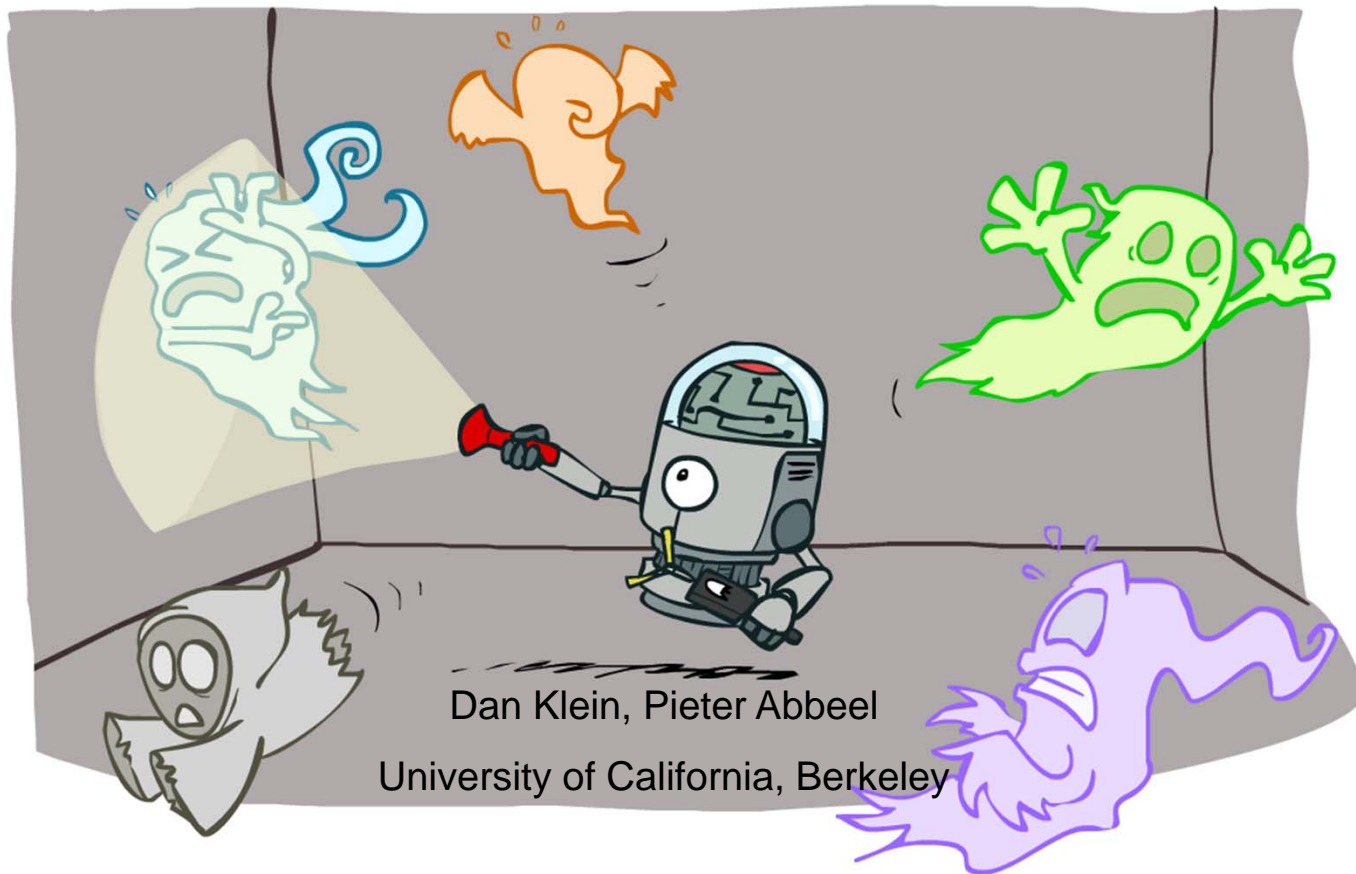


CS 188: Artificial Intelligence

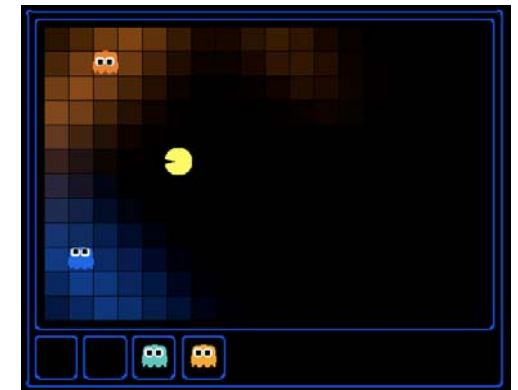
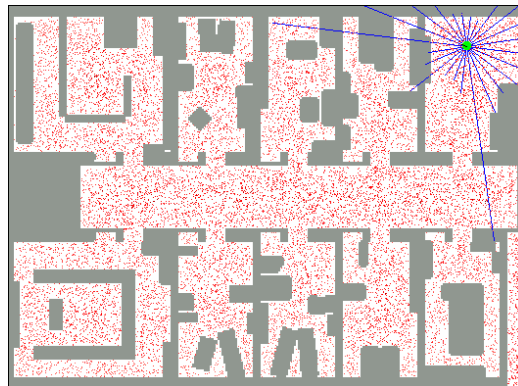
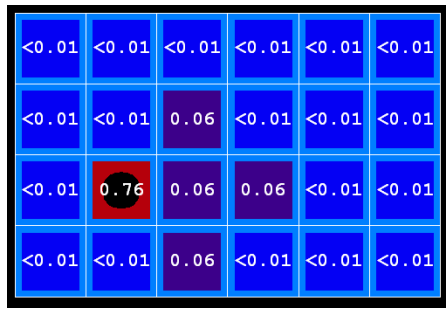
Advanced HMMs



Today

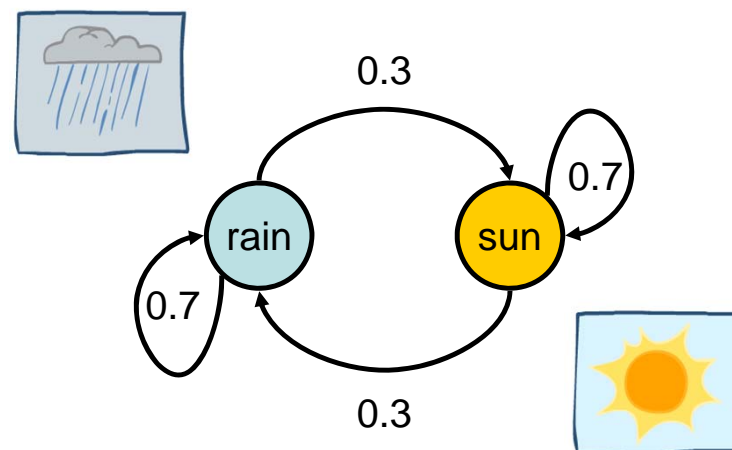
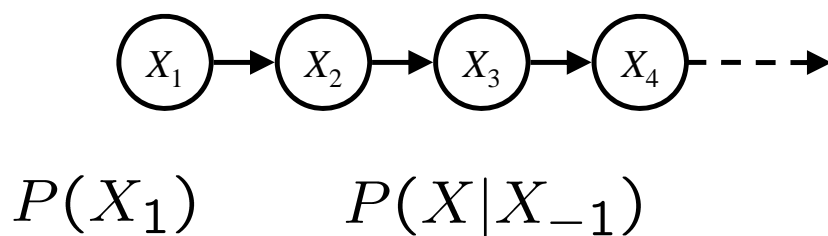
- HMMs
 - Demo bonanza!
 - Most-likely-explanation queries
- Speech recognition
 - A massive HMM!
 - Details of this section not required
- Start machine learning

Demo Bonanza!

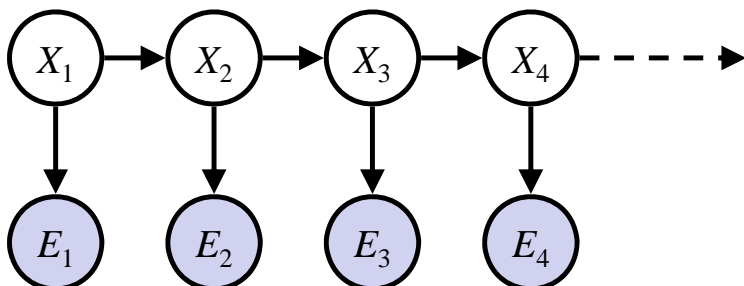


Recap: Reasoning Over Time

Markov models



Hidden Markov models



$$P(E|X)$$

X	E	P
rain	umbrella	0.9
rain	no umbrella	0.1
sun	umbrella	0.2
sun	no umbrella	0.8

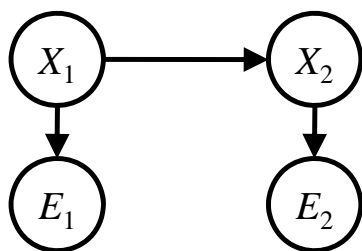
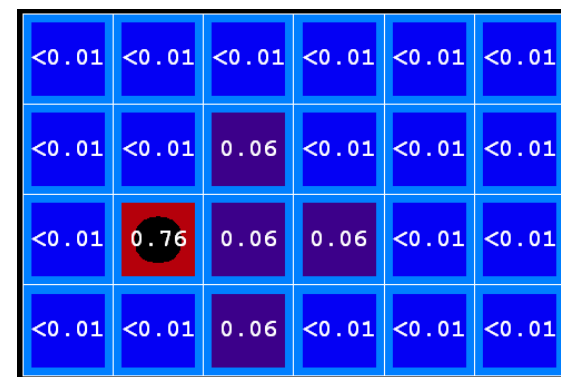
Recap: Filtering

Elapse time: compute $P(X_t | e_{1:t-1})$

$$P(x_t | e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1} | e_{1:t-1}) \cdot P(x_t | x_{t-1})$$

Observe: compute $P(X_t | e_{1:t})$

$$P(x_t | e_{1:t}) \propto P(x_t | e_{1:t-1}) \cdot P(e_t | x_t)$$



Belief: $\langle P(\text{rain}), P(\text{sun}) \rangle$

$P(X_1)$ $\langle 0.5, 0.5 \rangle$ *Prior on X_1*

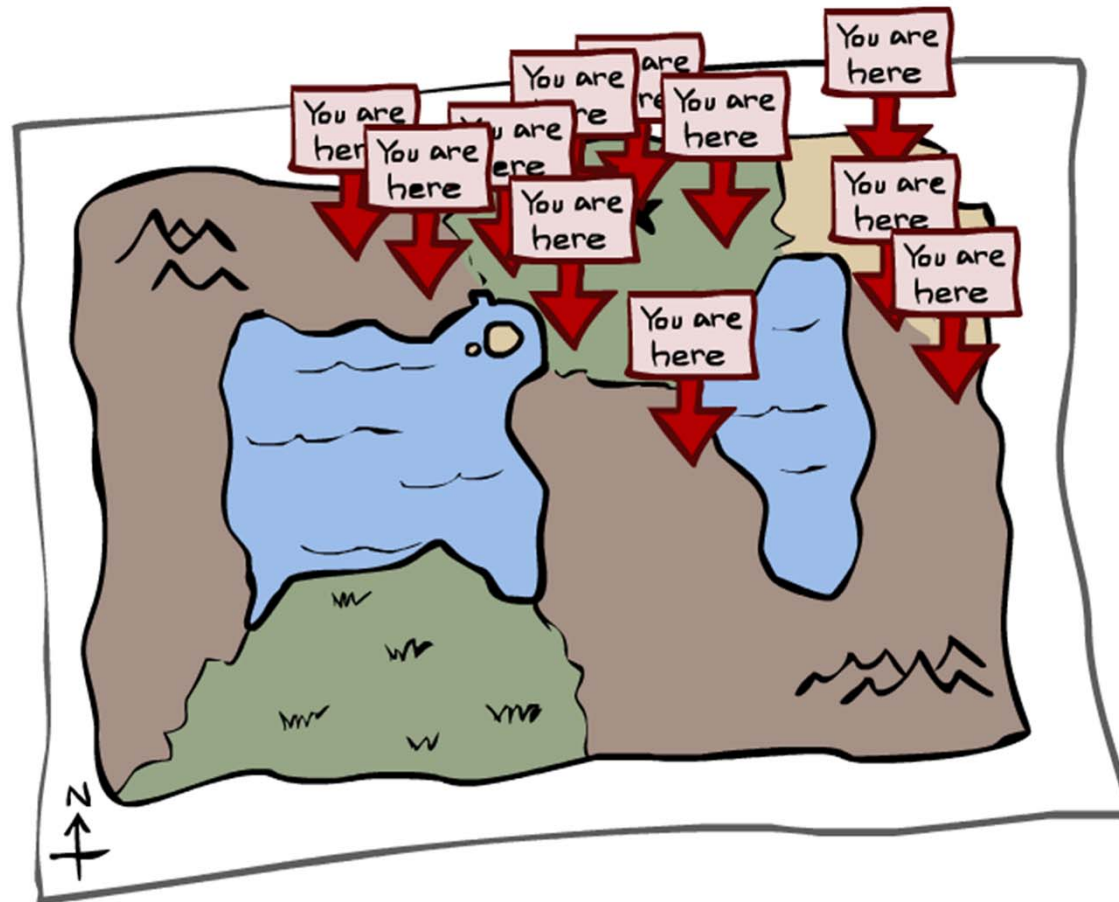
$P(X_1 | E_1 = \text{umbrella})$ $\langle 0.82, 0.18 \rangle$ *Observe*

$P(X_2 | E_1 = \text{umbrella})$ $\langle 0.63, 0.37 \rangle$ *Elapse time*

$P(X_2 | E_1 = \text{umb}, E_2 = \text{umb})$ $\langle 0.88, 0.12 \rangle$ *Observe*

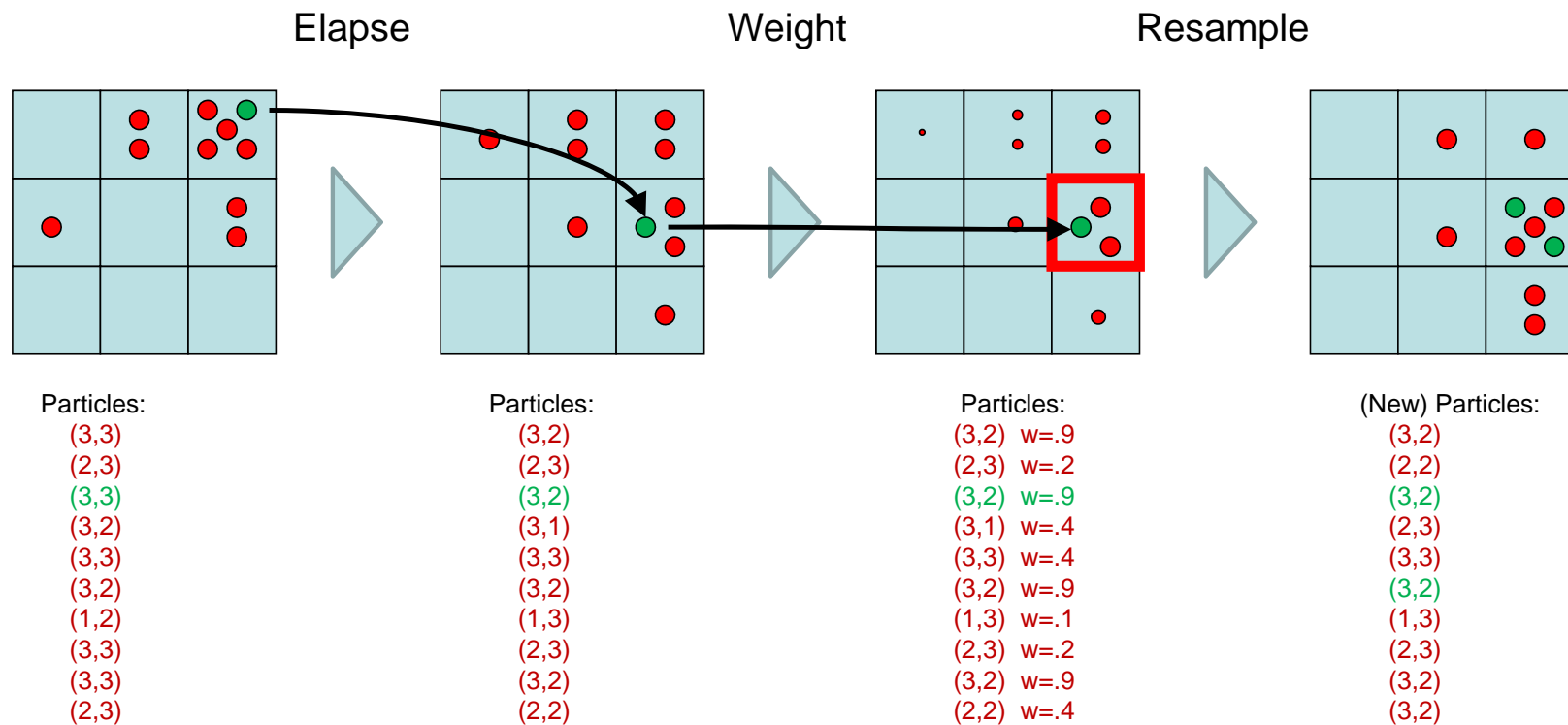
[demo: exact filtering]

Particle Filtering



Recap: Particle Filtering

- Particles: track samples of states rather than an explicit distribution

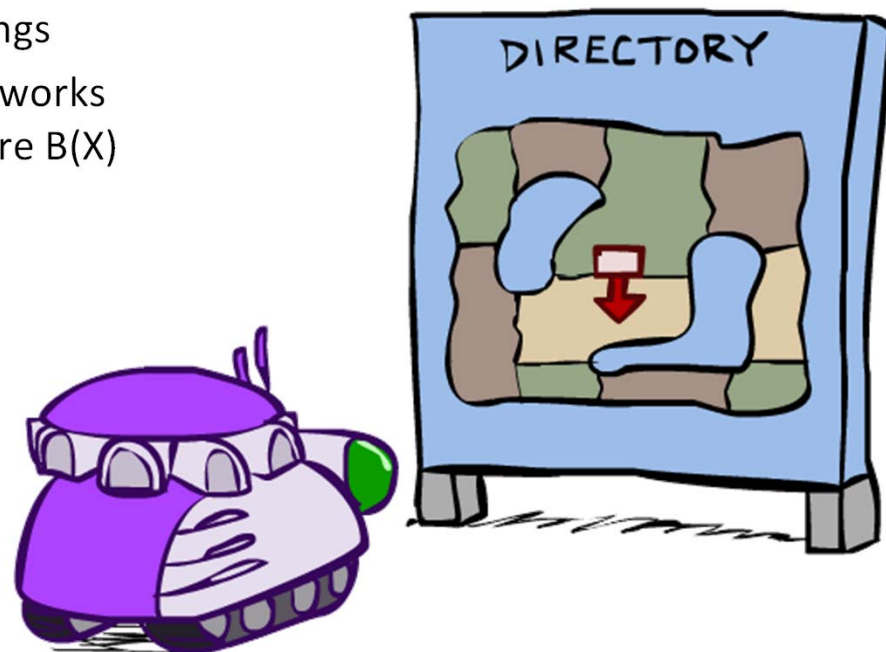
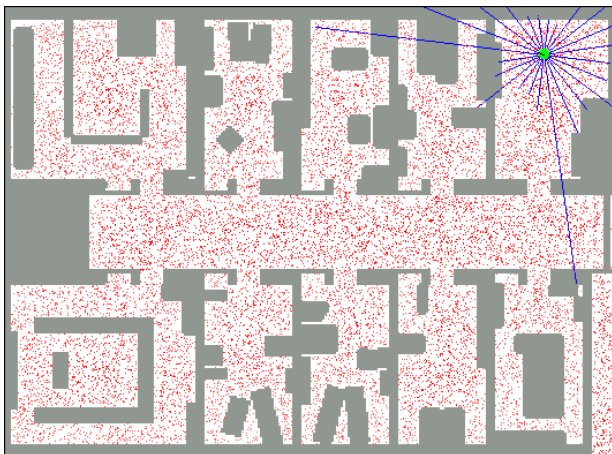


[demo: particle filtering]

Robot Localization

- In robot localization:

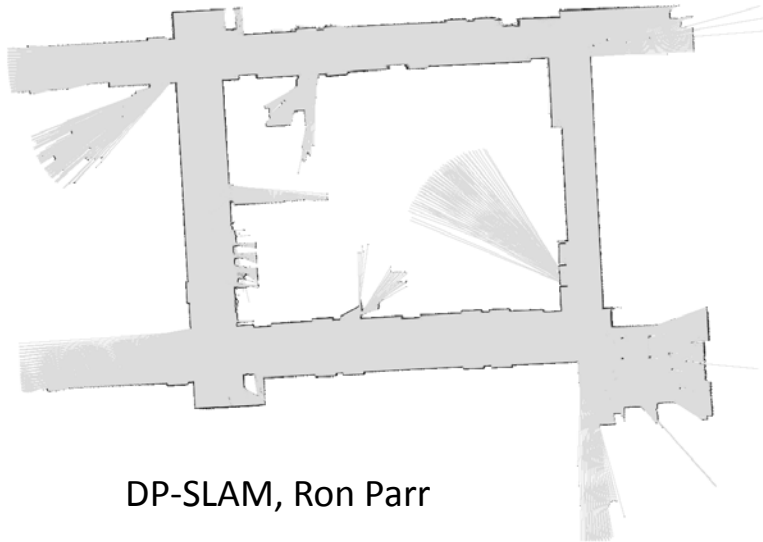
- We know the map, but not the robot's position
- Observations may be vectors of range finder readings
- State space and readings are typically continuous (works basically like a very fine grid) and so we cannot store $B(X)$
- Particle filtering is a main technique



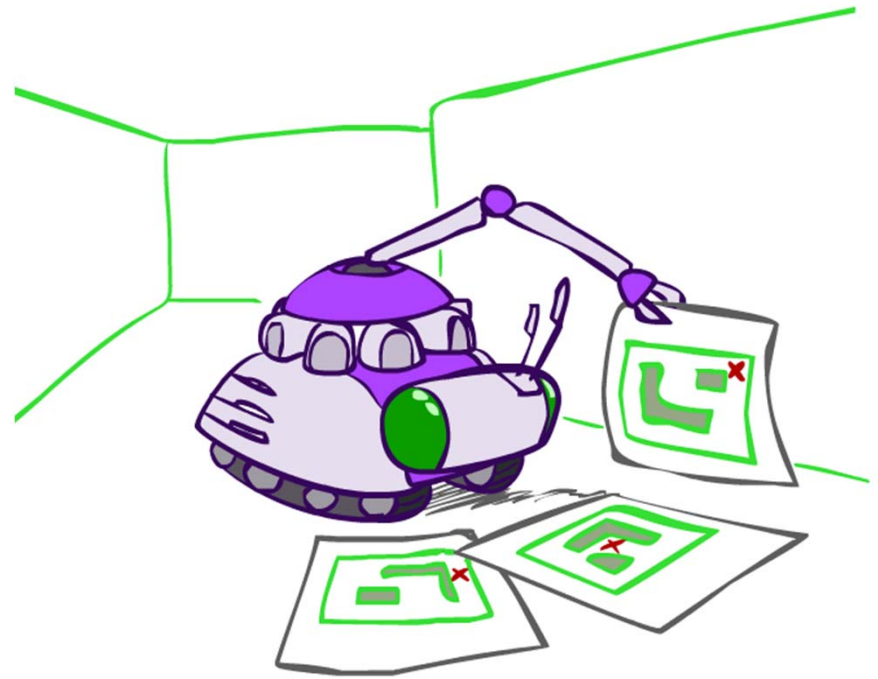
Robot Mapping

- SLAM: Simultaneous Localization And Mapping

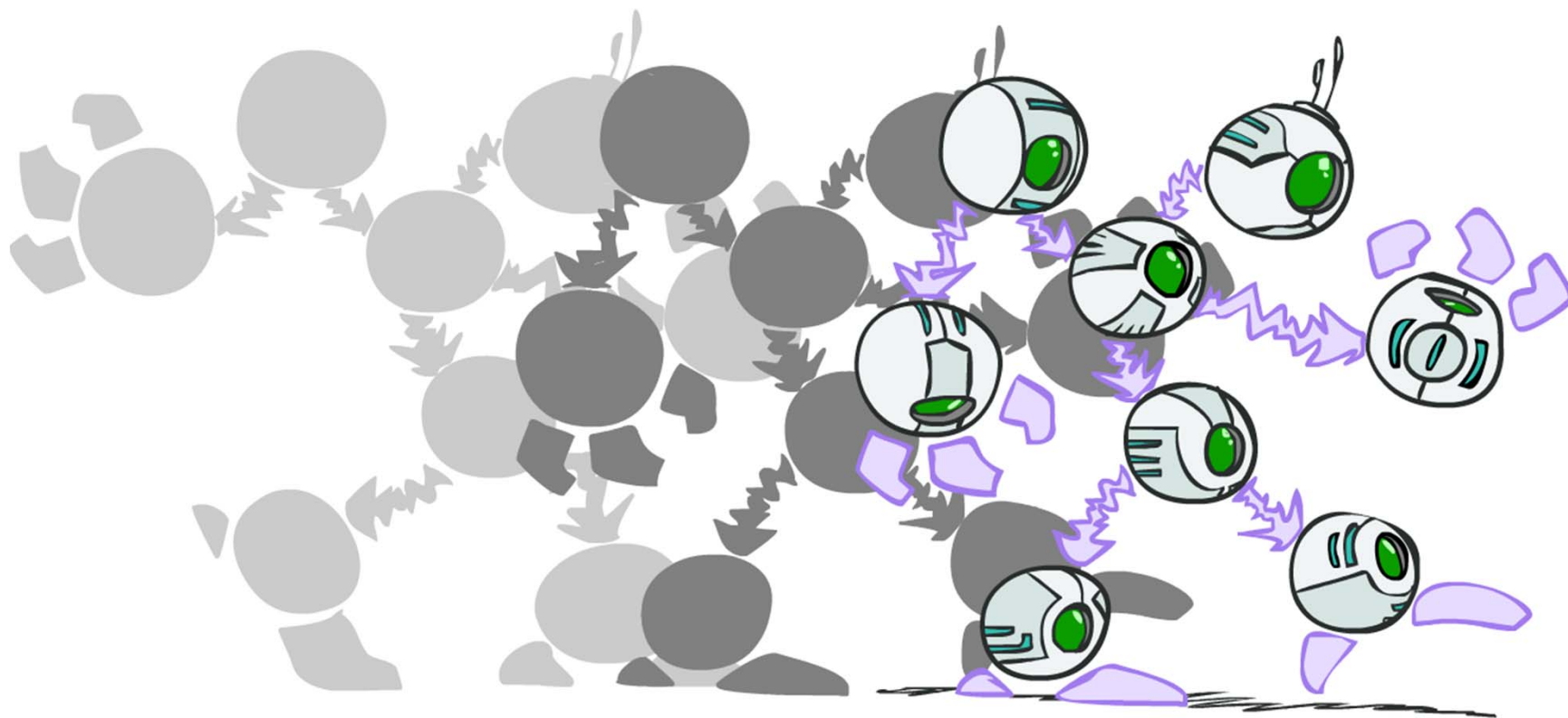
- We do not know the map or our location
- State consists of position AND map!
- Main techniques: Kalman filtering (Gaussian HMMs) and particle methods



DP-SLAM, Ron Parr

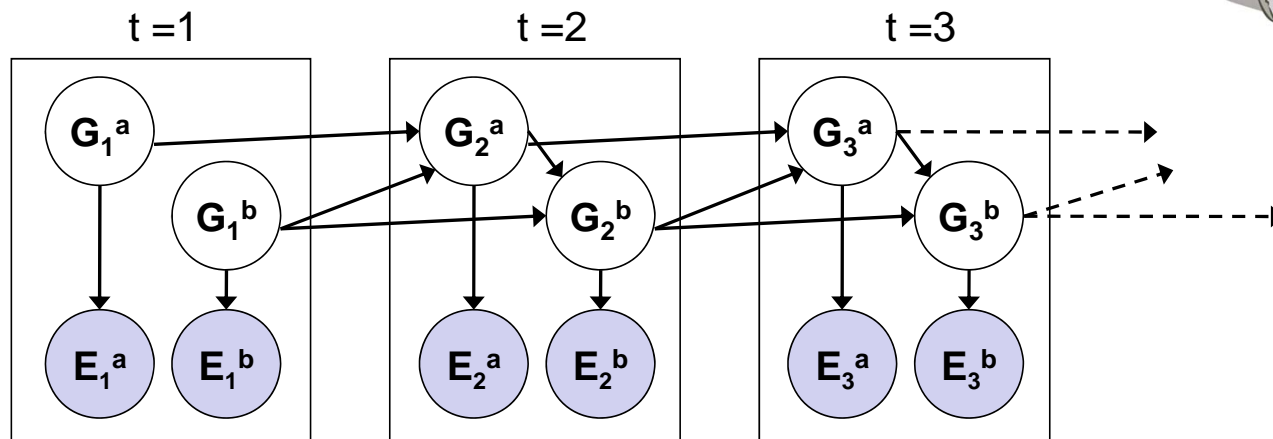
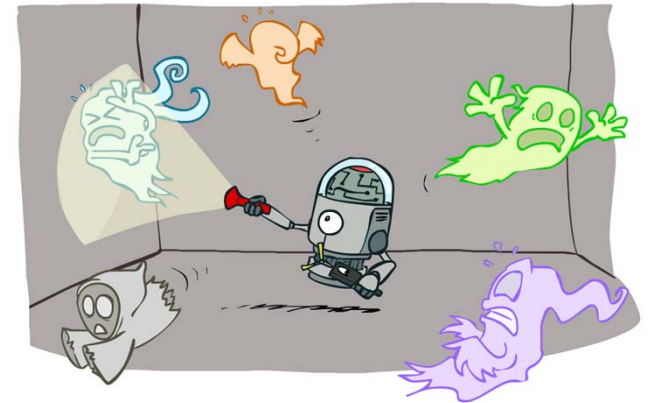


Dynamic Bayes Nets



Dynamic Bayes Nets (DBNs)

- We want to track multiple variables over time, using multiple sources of evidence
- Idea: Repeat a fixed Bayes net structure at each time
- Variables from time t can condition on those from $t-1$



- Dynamic Bayes nets are a generalization of HMMs

DBN Particle Filters

- A particle is a complete sample for a time step
- **Initialize:** Generate prior samples for the $t=1$ Bayes net
 - Example particle: $\mathbf{G}_1^a = (3,3)$ $\mathbf{G}_1^b = (5,3)$
- **Elapse time:** Sample a successor for each particle
 - Example successor: $\mathbf{G}_2^a = (2,3)$ $\mathbf{G}_2^b = (6,3)$
- **Observe:** Weight each entire sample by the likelihood of the evidence conditioned on the sample
 - Likelihood: $P(\mathbf{E}_1^a | \mathbf{G}_1^a) * P(\mathbf{E}_1^b | \mathbf{G}_1^b)$
- **Resample:** Select prior samples (tuples of values) in proportion to their likelihood

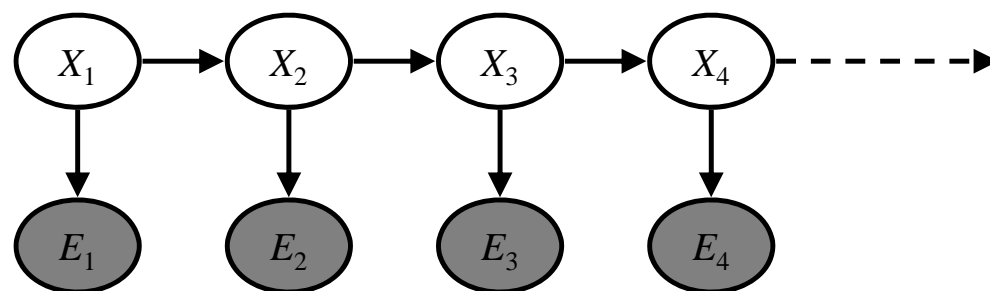
Most Likely Explanation



HMMs: MLE Queries

- HMMs defined by

- States X
- Observations E
- Initial distribution: $P(X_1)$
- Transitions: $P(X|X_{-1})$
- Emissions: $P(E|X)$



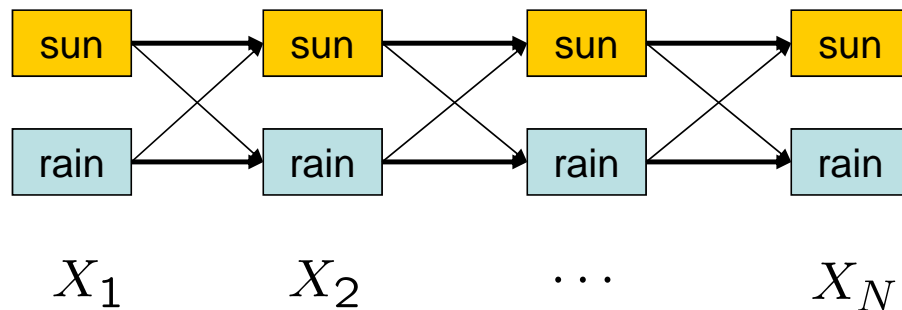
- New query: most likely explanation:

$$\arg \max_{x_{1:t}} P(x_{1:t}|e_{1:t})$$

- New method: the Viterbi algorithm

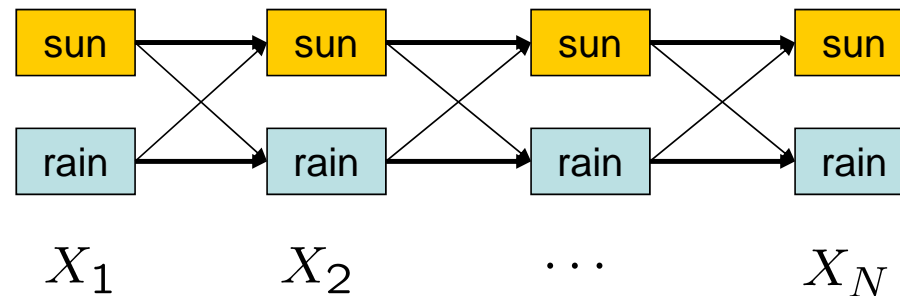
State Trellis

- State trellis: graph of states and transitions over time



- Each arc represents some transition $x_{t-1} \rightarrow x_t$
- Each arc has weight $P(x_t|x_{t-1})P(e_t|x_t)$
- Each path is a sequence of states (evidence)
- The product of weights on a path is that sequence's probability along with the evidence
- Forward algorithm computes sums of paths, Viterbi computes best paths

Forward / Viterbi Algorithms



Forward Algorithm (Sum)

$$f_t[x_t] = P(x_t, e_{1:t})$$

$$= P(e_t|x_t) \sum_{x_{t-1}} P(x_t|x_{t-1}) f_{t-1}[x_{t-1}]$$

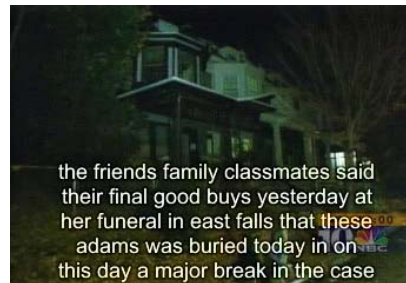
Viterbi Algorithm (Max)

$$m_t[x_t] = \max_{x_{1:t-1}} P(x_{1:t-1}, x_t, e_{1:t})$$

$$= P(e_t|x_t) \max_{x_{t-1}} P(x_t|x_{t-1}) m_{t-1}[x_{t-1}]$$

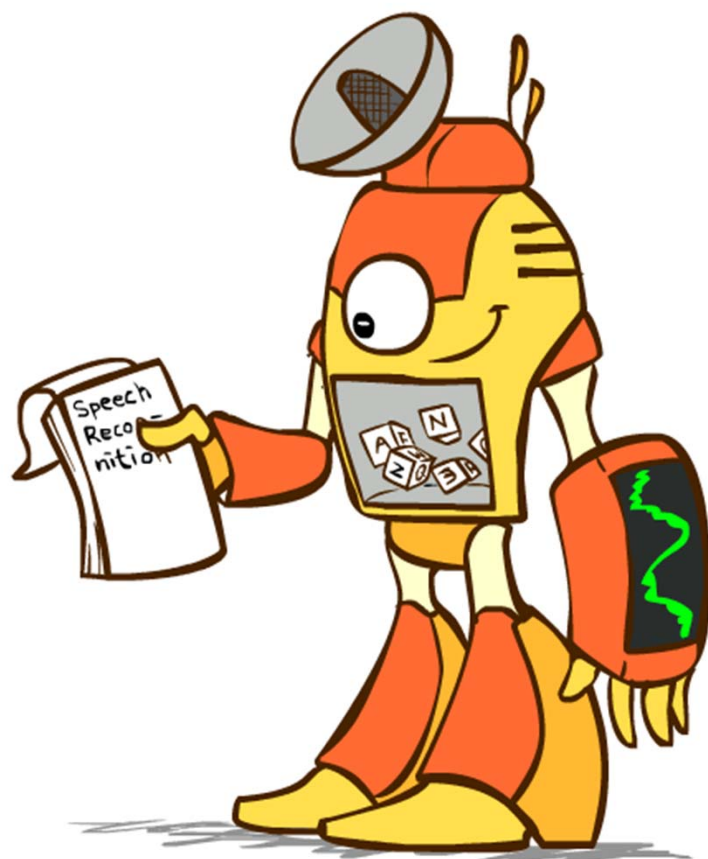
Natural Language

- Speech technologies (e.g. Siri)
 - Automatic speech recognition (ASR)
 - Text-to-speech synthesis (TTS)
 - Dialog systems
- Language processing technologies
 - Question answering
 - Machine translation

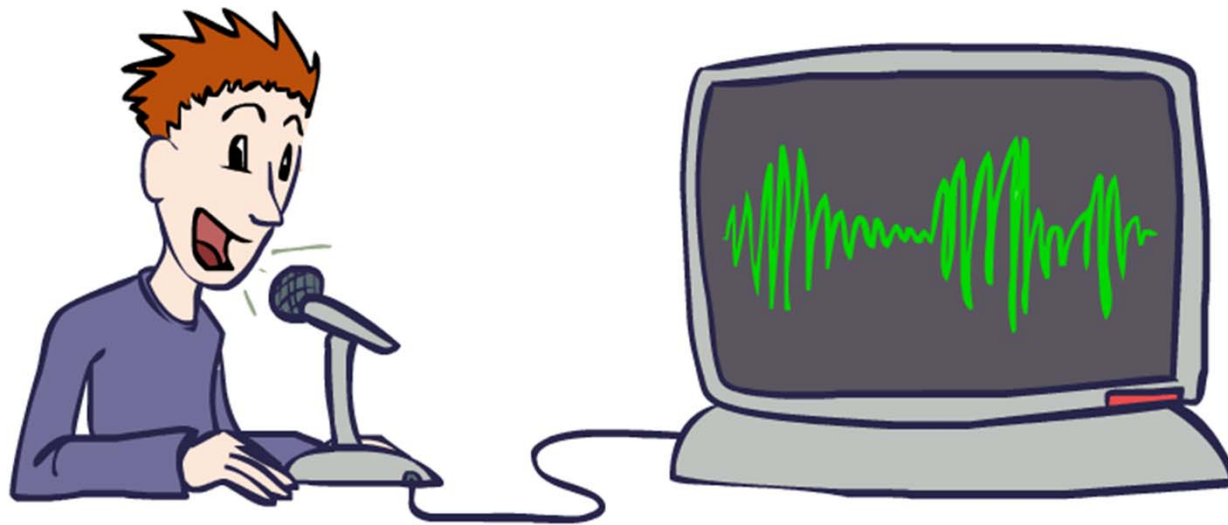


- Web search
- Text classification, spam filtering, etc...

Speech Recognition



Digitizing Speech



Speech in an Hour

- Speech input is an acoustic waveform

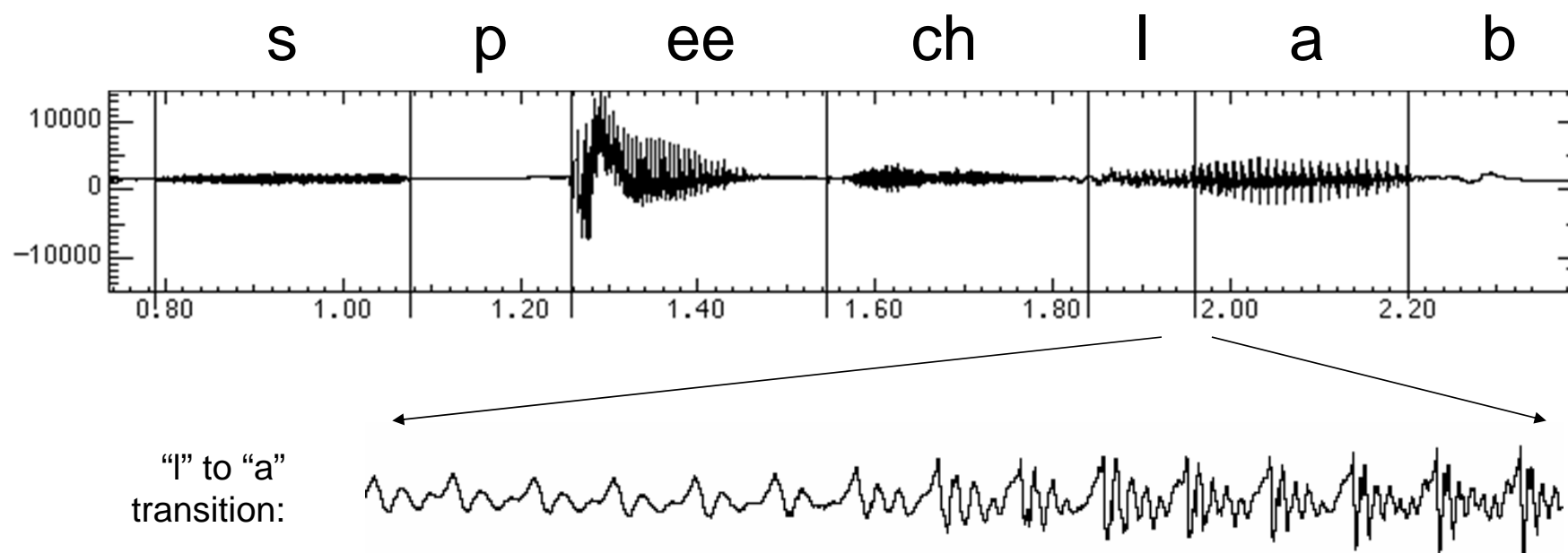
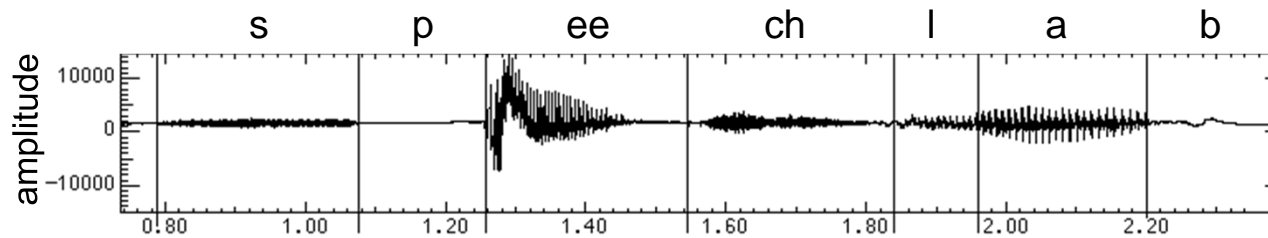


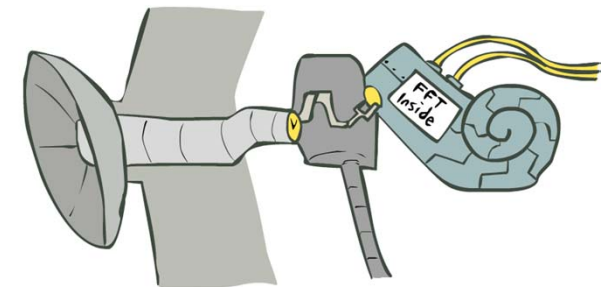
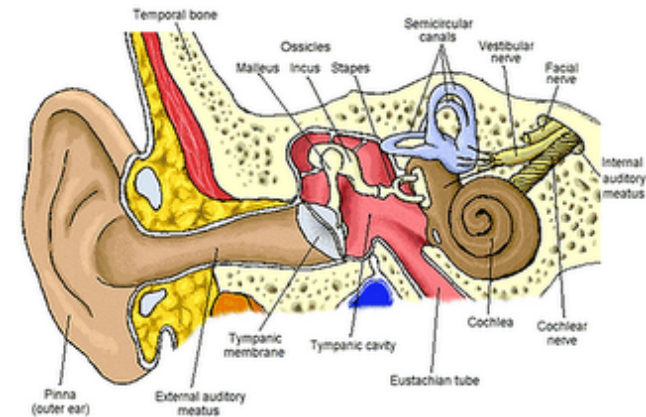
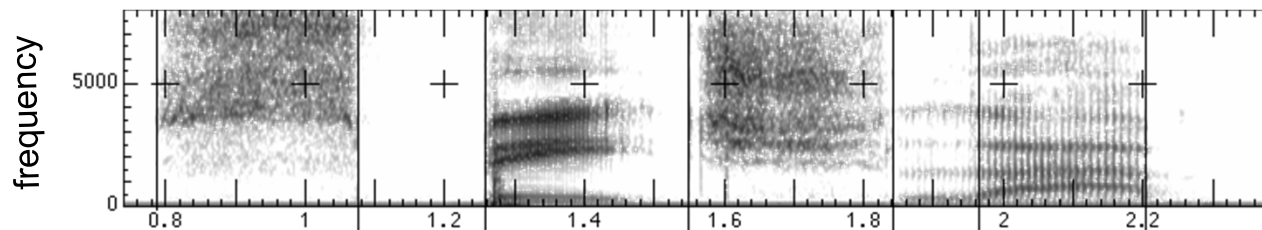
Figure: Simon Arnfield, <http://www.psyc.leeds.ac.uk/research/cogn/speech/tutorial/>

Spectral Analysis

- Frequency gives pitch; amplitude gives volume
 - Sampling at ~8 kHz (phone), ~16 kHz (mic) (kHz=1000 cycles/sec)

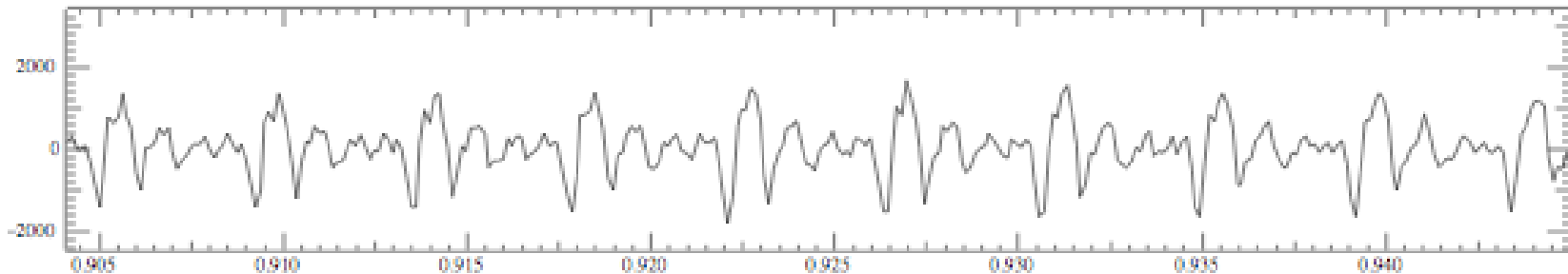


- Fourier transform of wave displayed as a spectrogram
 - Darkness indicates energy at each frequency



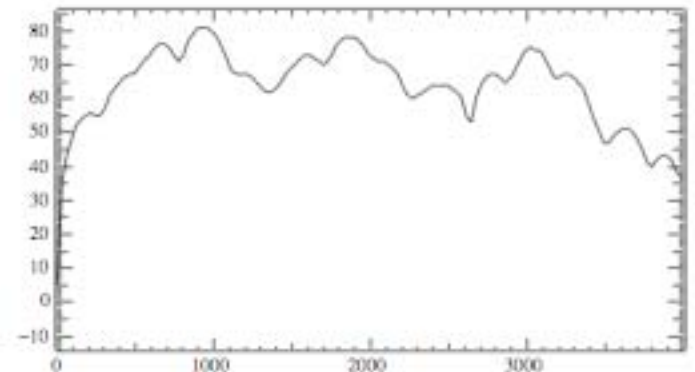
Human ear figure: depion.blogspot.com

Part of [ae] from “lab”



- Complex wave repeating nine times

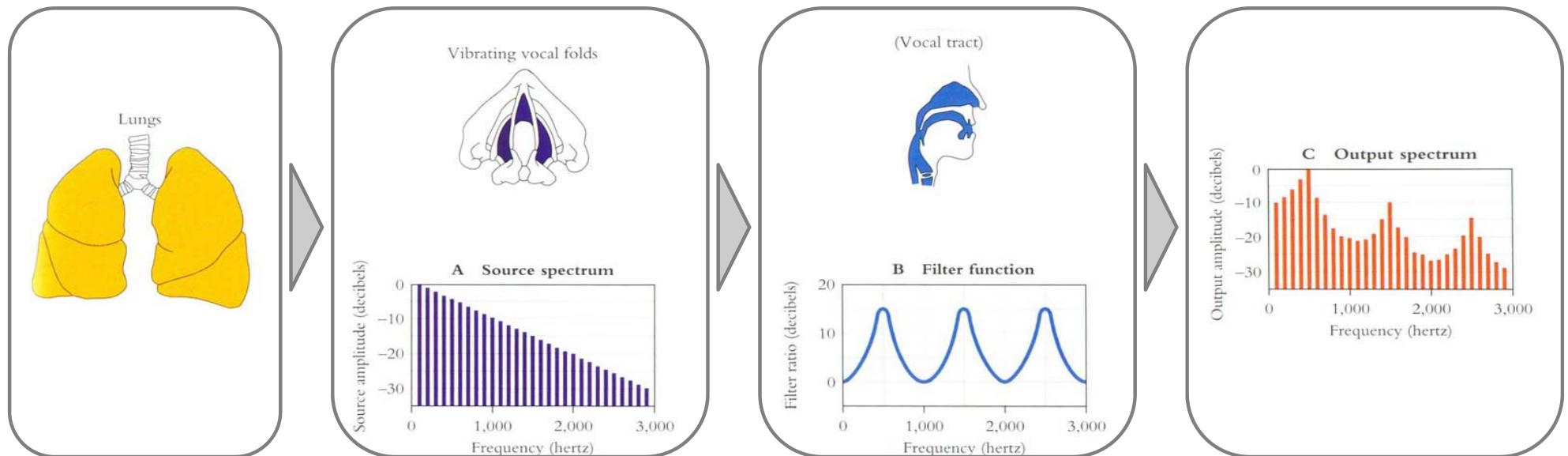
- Plus smaller wave that repeats 4x for every large cycle
- Large wave: freq of 250 Hz (9 times in .036 seconds)
- Small wave roughly 4 times this, or roughly 1000 Hz



Why These Peaks?

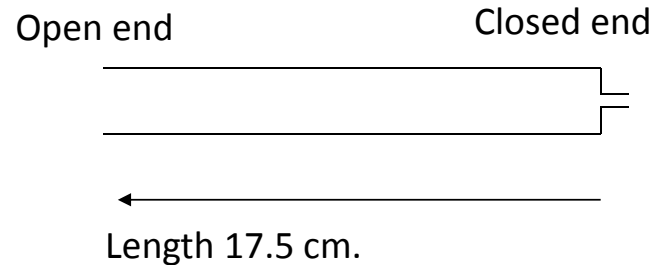
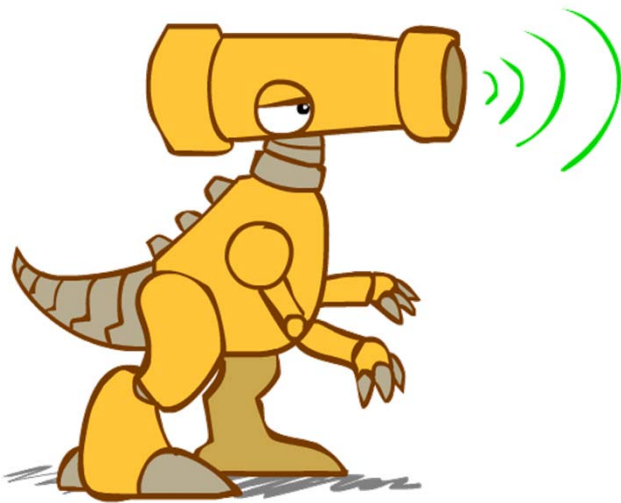
■ Articulator process:

- Vocal cord vibrations create harmonics
- The mouth is an amplifier
- Depending on shape of mouth, some harmonics are amplified more than others



Resonances of the Vocal Tract

- The human vocal tract as an open tube



- Air in a tube of a given length will tend to vibrate at resonance frequency of tube
- Constraint: Pressure differential should be maximal at (closed) glottal end and minimal at (open) lip end

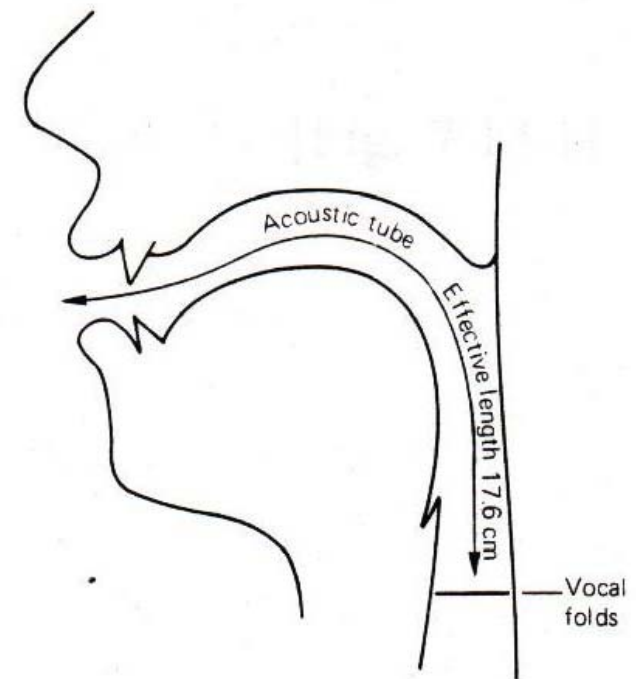
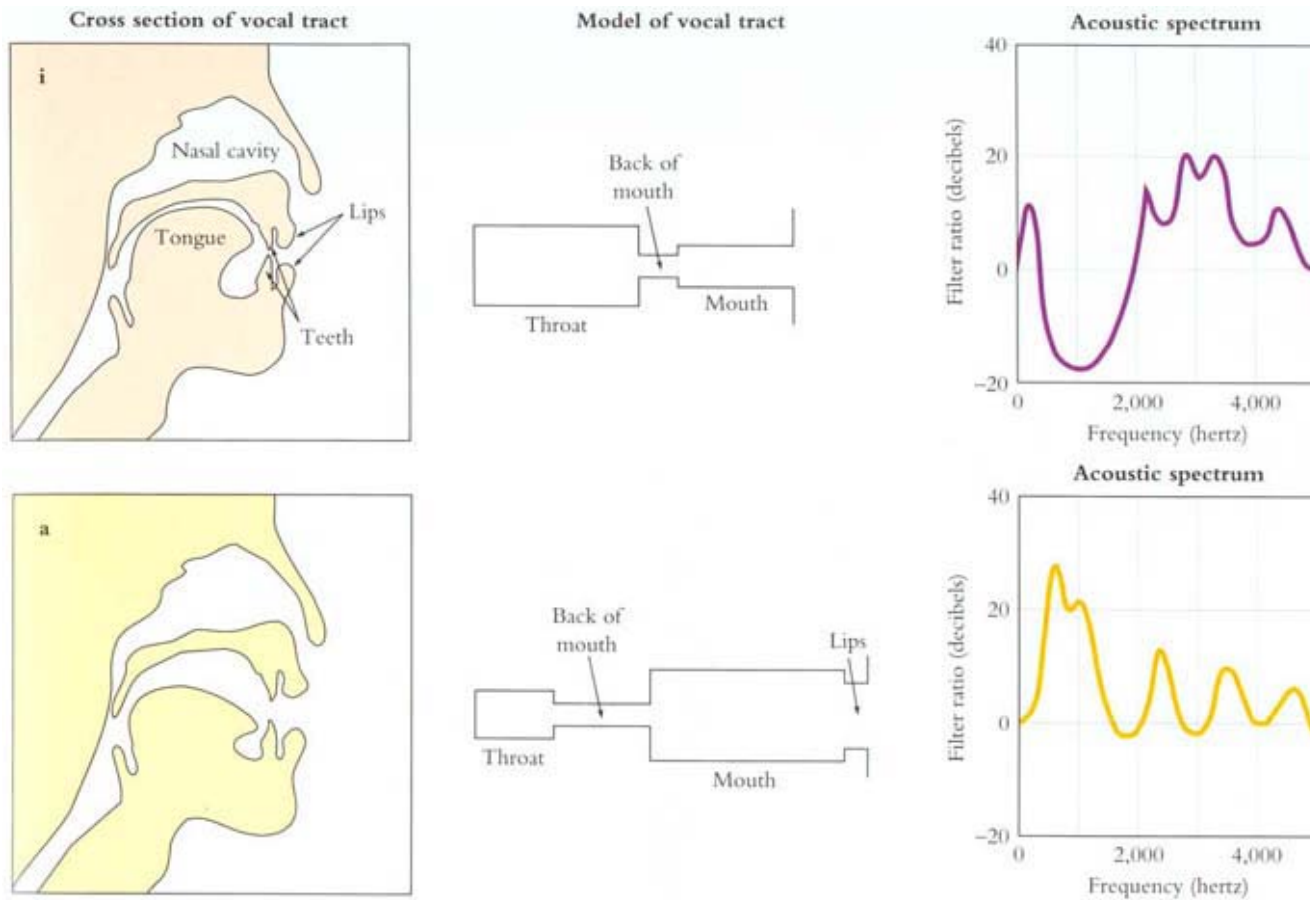


Figure: W. Barry Speech Science slides

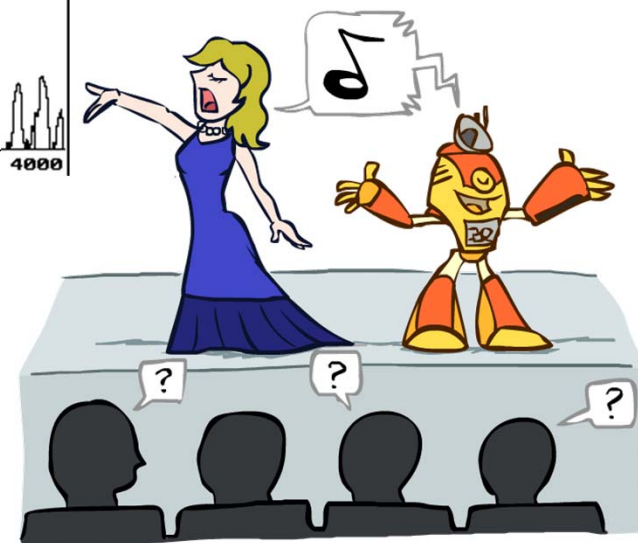
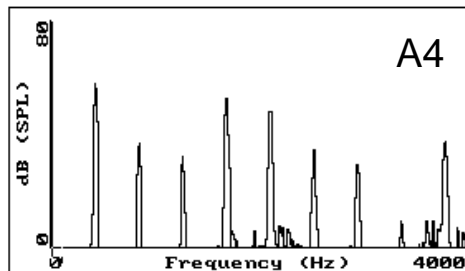
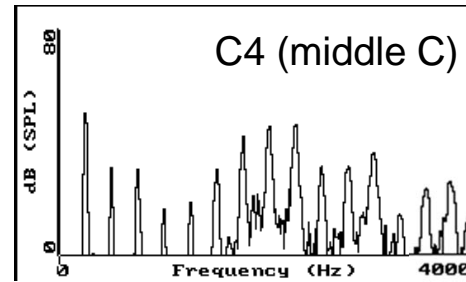
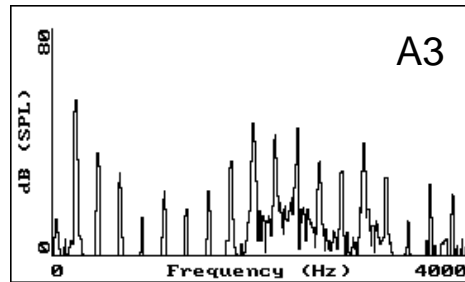
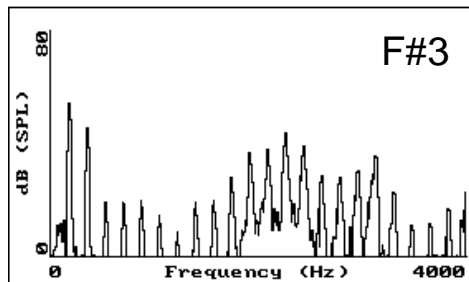
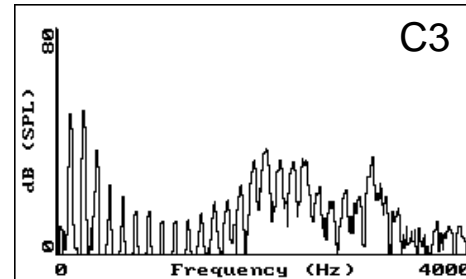
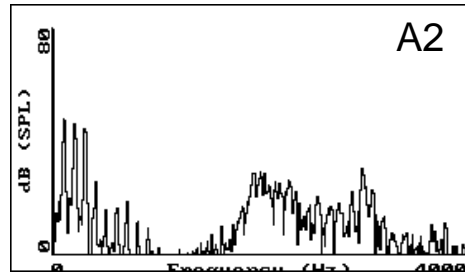
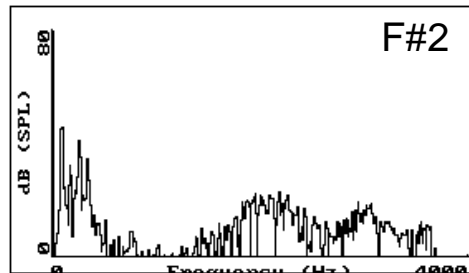
Spectrum Shapes



[demo]

Figure: Mark Liberman

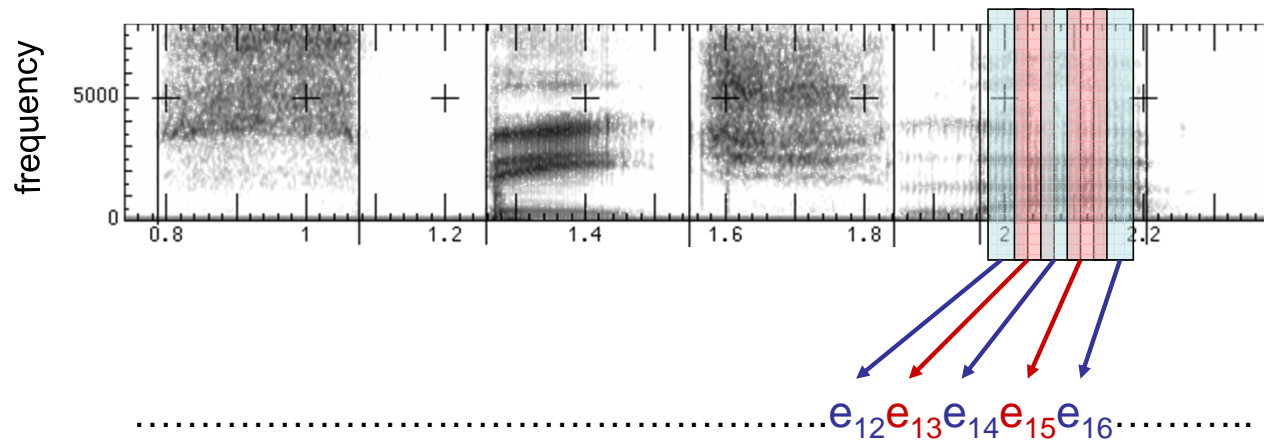
Vowel [i] sung at successively higher pitches



Graphs: Ratree Wayland

Acoustic Feature Sequence

- Time slices are translated into acoustic feature vectors (~39 real numbers per slice)



- These are the observations E , now we need the hidden states X

Speech State Space

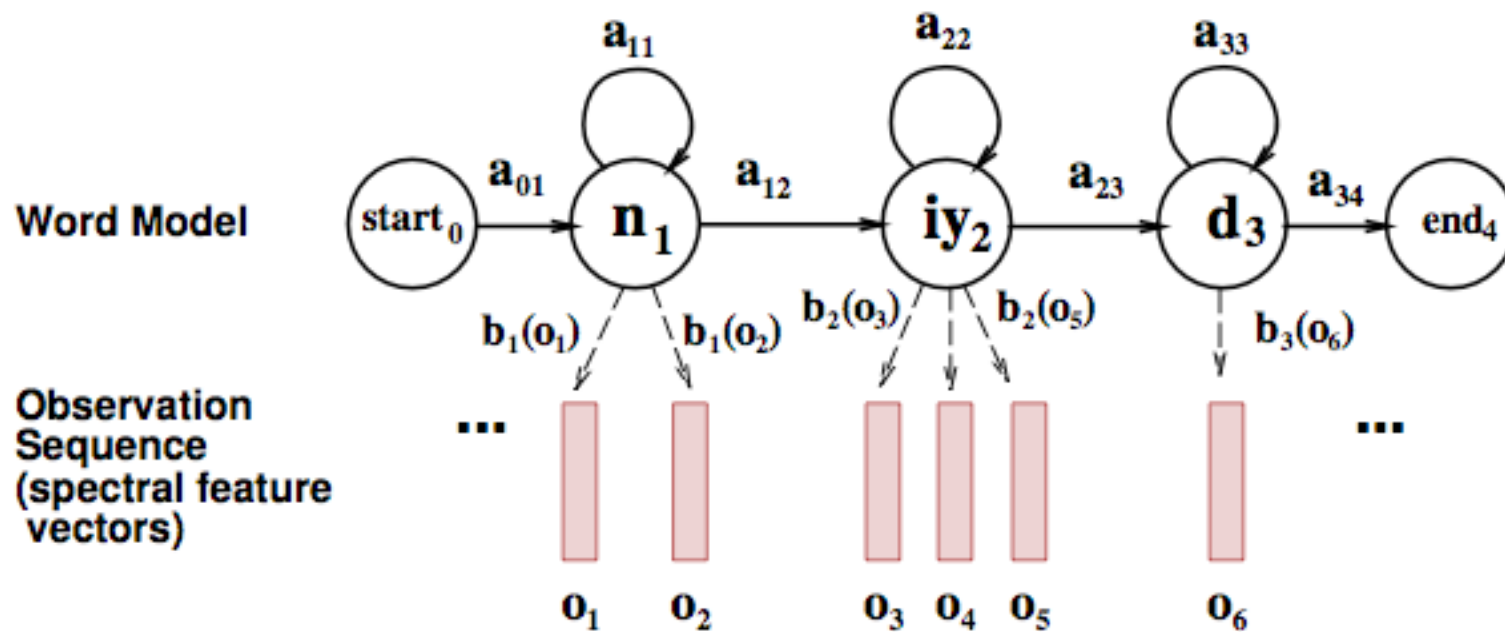
- HMM Specification

- $P(E|X)$ encodes which acoustic vectors are appropriate for each phoneme (each kind of sound)
- $P(X|X')$ encodes how sounds can be strung together

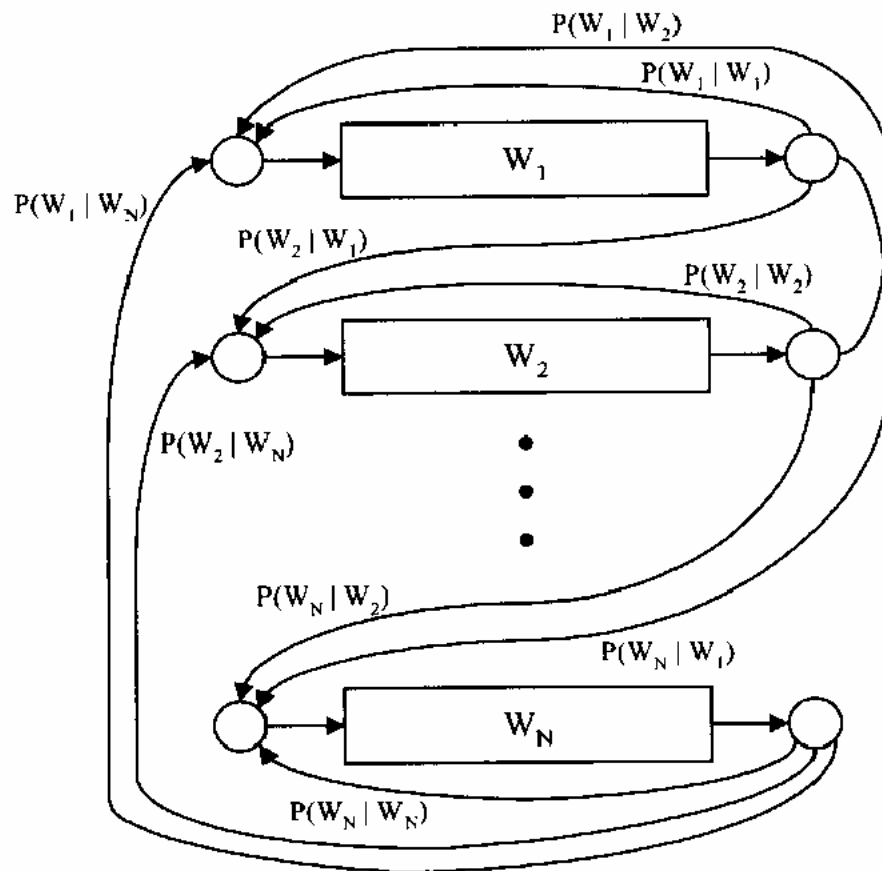
- State Space

- We will have one state for each sound in each word
- Mostly, states advance sound by sound
- Build a little state graph for each word and chain them together to form the state space X

States in a Word



Transitions with a Bigram Model



Training Counts

198015222 the first
 194623024 the same
 168504105 the following
 158562063 the world
 ...
 14112454 the door

 23135851162 the *

$$\hat{P}(\text{door}|\text{the}) = \frac{14112454}{23135851162}$$

$$= 0.0006$$

Figure: Huang et al, p. 618

Decoding

- Finding the words given the acoustics is an HMM inference problem
- Which state sequence $x_{1:T}$ is most likely given the evidence $e_{1:T}$?

$$x_{1:T}^* = \arg \max_{x_{1:T}} P(x_{1:T} | e_{1:T}) = \arg \max_{x_{1:T}} P(x_{1:T}, e_{1:T})$$

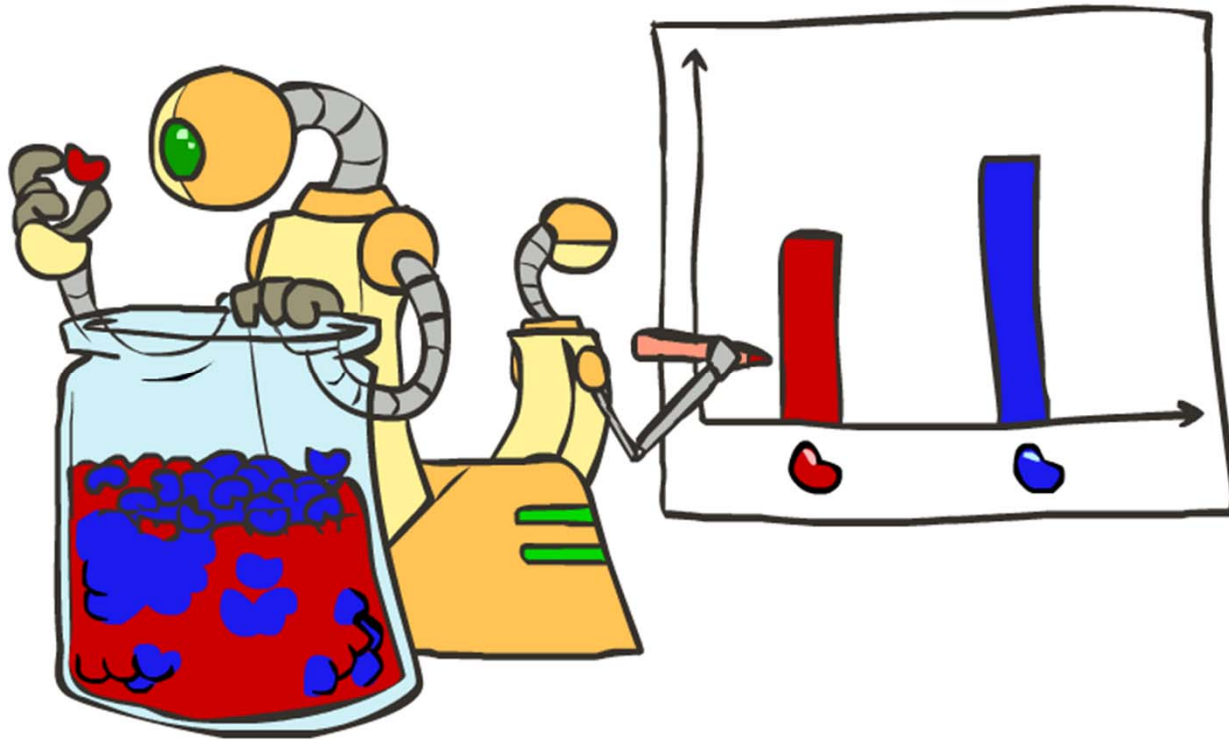
- From the sequence x , we can simply read off the words



End of Part II!

- Now we're done with our unit on probabilistic reasoning
- Last part of class: machine learning

Machine Learning



Machine Learning

- Up until now: how use a model to make optimal decisions
- Machine learning: how to acquire a model from data / experience
 - Learning parameters (e.g. probabilities)
 - Learning structure (e.g. BN graphs)
 - Learning hidden concepts (e.g. clustering)

Parameter Estimation

- Estimating the distribution of a random variable
- *Elicitation*: ask a human (why is this hard?)
- *Empirically*: use training data (learning!)
 - E.g.: for each outcome x , look at the *empirical rate* of that value:

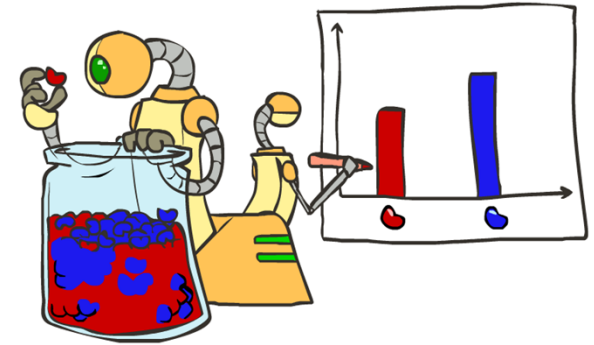
$$P_{\text{ML}}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

(r) (r) (b)

$$P_{\text{ML}}(\textcolor{red}{r}) = 2/3$$

- This is the estimate that maximizes the *likelihood of the data*

$$L(x, \theta) = \prod_i P_{\theta}(x_i)$$



Estimation: Smoothing

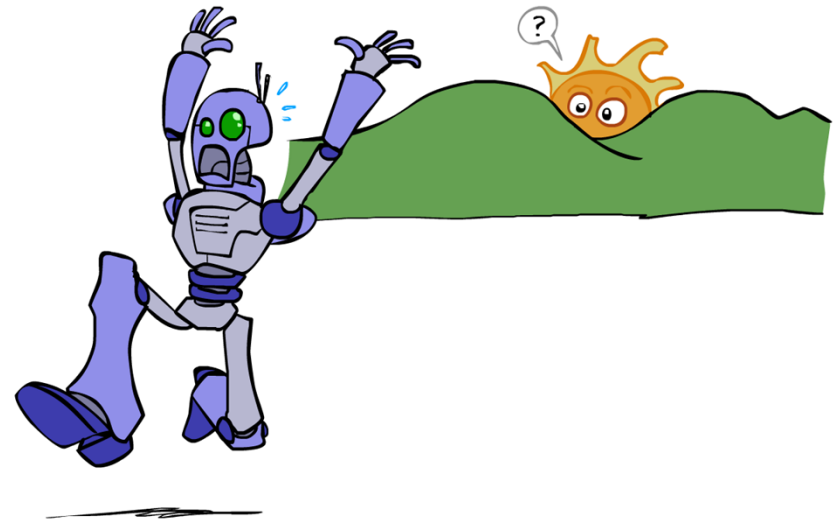
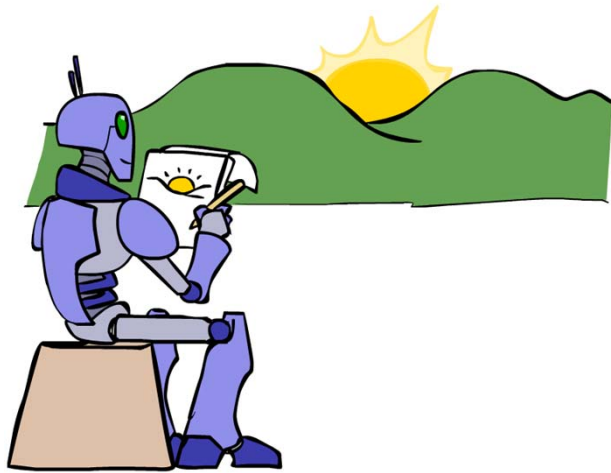
- Relative frequencies are the maximum likelihood estimates

$$\begin{aligned}\theta_{ML} &= \arg \max_{\theta} P(\mathbf{X}|\theta) \\ &= \arg \max_{\theta} \prod_i P_{\theta}(X_i)\end{aligned} \quad \Rightarrow \quad P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

- Another option is to consider the most likely parameter value given the data

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} P(\theta|\mathbf{X}) \\ &= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)/P(\mathbf{X}) \quad \Rightarrow \quad \text{????} \\ &= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)\end{aligned}$$

Smoothing



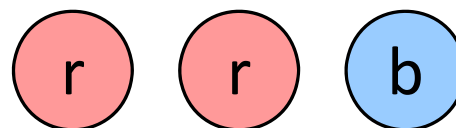
Estimation: Laplace Smoothing

- Laplace's estimate:

- Pretend you saw every outcome once more than you actually did

$$\begin{aligned} P_{LAP}(x) &= \frac{c(x) + 1}{\sum_x [c(x) + 1]} \\ &= \frac{c(x) + 1}{N + |X|} \end{aligned}$$

- Can derive this estimate with *Dirichlet priors* (see cs281a)



$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

Estimation: Laplace Smoothing

- Laplace's estimate (extended):

- Pretend you saw every outcome k extra times

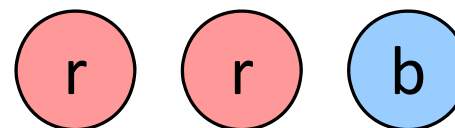
$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

- What's Laplace with $k = 0$?
- k is the **strength** of the prior

- Laplace for conditionals:

- Smooth each condition independently:

$$P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$$



$$P_{LAP,0}(X) =$$

$$P_{LAP,1}(X) =$$

$$P_{LAP,100}(X) =$$