

# CSE 573: Artificial Intelligence

## Autumn 2010

Lecture 11: Hidden Markov Models II  
11/4/2010

Luke Zettlemoyer

Many slides over the course adapted from either Dan Klein,  
Stuart Russell or Andrew Moore

# Outline

---

- Probabilistic sequence models (and inference)
  - (Review) Markov Chains
  - Hidden Markov Models
  - Particle Filters
  - Most Probable Explanations
  - Dynamic Bayesian networks

# Ghostbusters, Revisited

- Let's say we have two distributions:
  - Prior distribution** over ghost location:  $P(G)$ 
    - Let's say this is uniform
  - Sensor reading model:  $P(R | G)$ 
    - Given: we know what our sensors do
    - $R$  = reading color measured at  $(1,1)$
    - E.g.  $P(R = \text{yellow} | G=(1,1)) = 0.1$
- We can calculate the **posterior distribution**  $P(G|r)$  over ghost locations given a reading using Bayes' rule:

$$P(g|r) \propto P(r|g)P(g)$$

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

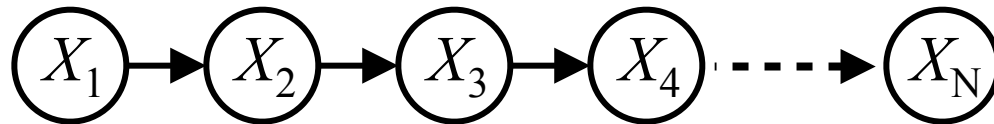
0.17	0.10	0.10
0.09	0.17	0.10
<0.01	0.09	0.17

# Recap: Markov Models

---

- A **Markov model** is:

- a MDP with no actions (and no rewards)
- a chain-structured Bayesian Network (BN)



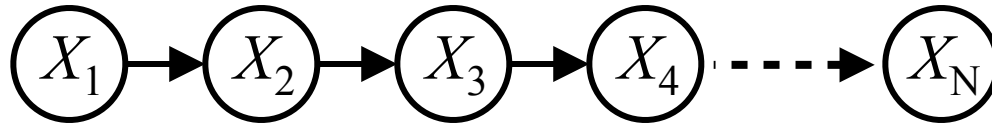
- A **Markov model** includes:

- Random variables  $X_t$  for all time steps  $t$  (the **state**)
- Parameters: called **transition probabilities** or dynamics, specify how the state evolves over time (also, initial probs)

$$P(X_1) \quad \text{and} \quad P(X_t | X_{t-1})$$

# Recap: Markov Models

---



- A **Markov model** defines
  - a joint probability distribution:

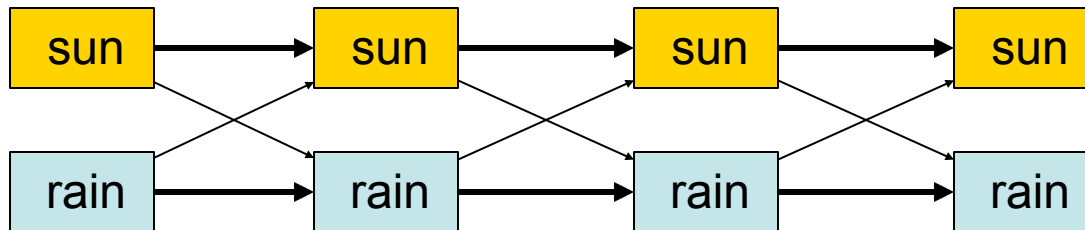
$$P(X_1, \dots, X_n) = P(X_1) \prod_{t=2}^N P(X_t | X_{t-1})$$

- One common inference problem:
  - Compute marginals  $P(X_t)$  for all time steps  $t$

# Recap: Mini-Forward Algorithm

---

- Question: What's  $P(X)$  on some day  $t$ ?
  - We don't need to enumerate every sequence!



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

$$P(x_1) = \text{known}$$

*Forward simulation*

# Recap: Stationary Distributions

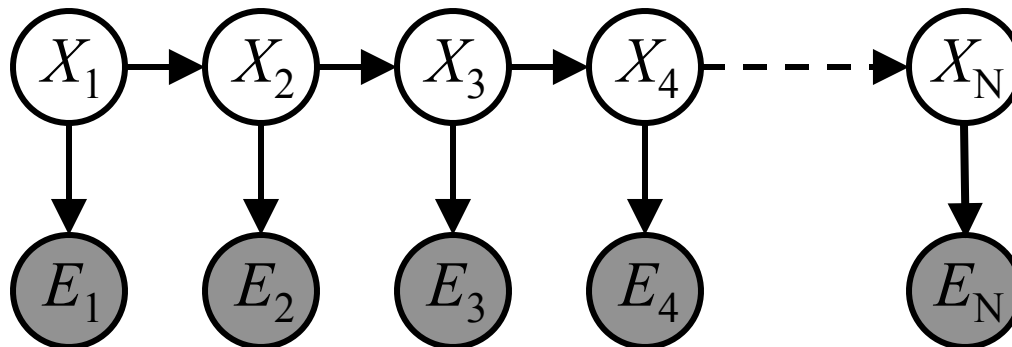
---

- If we simulate the chain long enough:
  - What happens?
  - Uncertainty accumulates
  - Eventually, we have no idea what the state is!
- Stationary distributions:
  - For most chains, the distribution we end up in is independent of the initial distribution
  - Called the **stationary distribution** of the chain
  - Usually, can only predict a short time out

# Hidden Markov Models

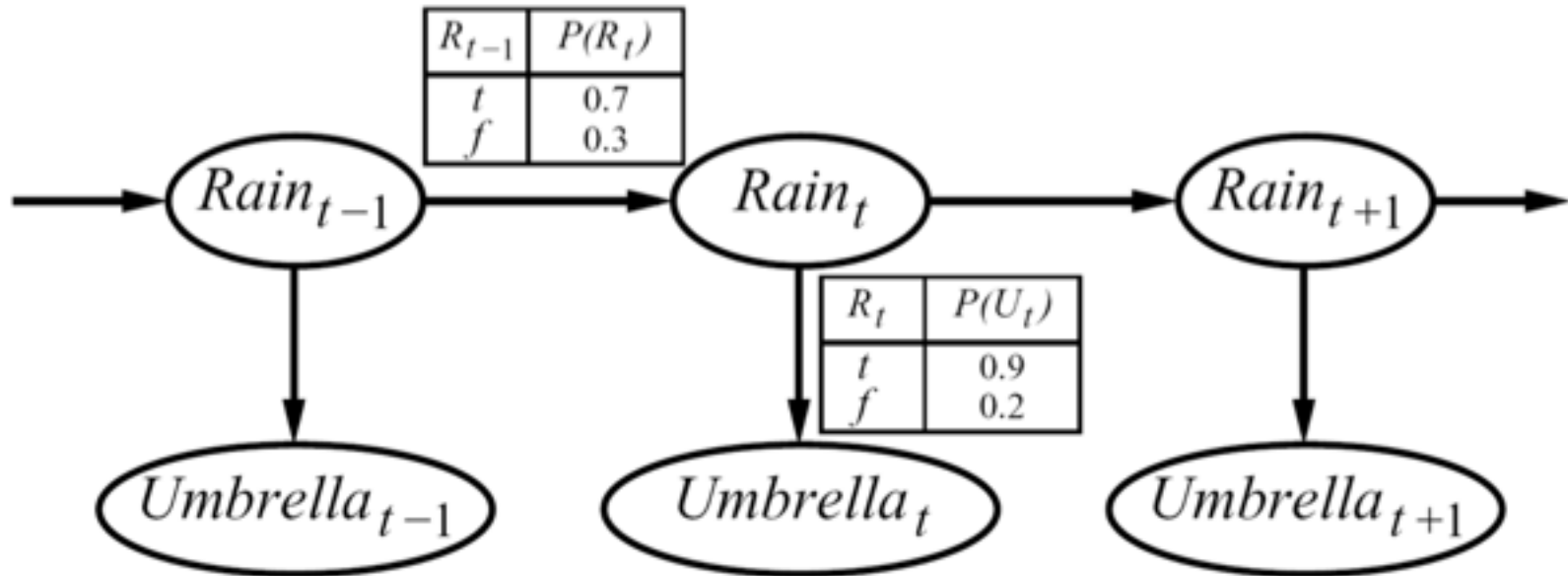
---

- Markov chains not so useful for most agents
  - Eventually you don't know anything anymore
  - Need observations to update your beliefs
- Hidden Markov models (HMMs)
  - Underlying Markov chain over states  $S$
  - You observe outputs (effects) at each time step
  - POMDPs without actions (or rewards).
  - As a Bayes' net:





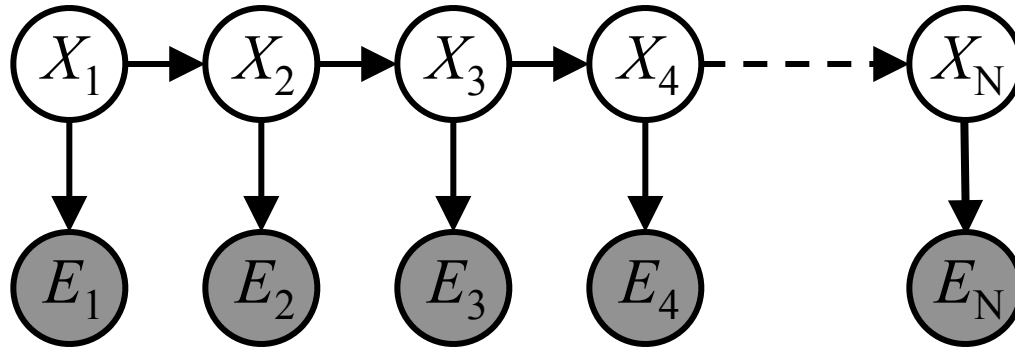
# Example



- An HMM is defined by:
  - Initial distribution:  $P(X_1)$
  - Transitions:  $P(X_t|X_{t-1})$
  - Emissions:  $P(E|X)$

# Hidden Markov Models

---



- Defines a joint probability distribution:

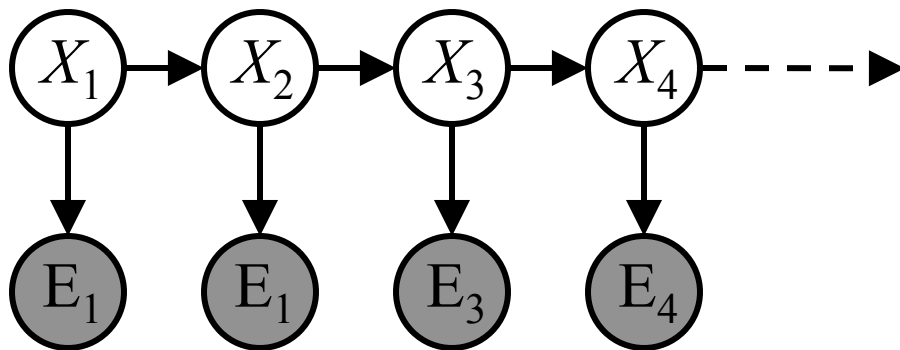
$$P(X_1, \dots, X_n, E_1, \dots, E_n) =$$

$$P(X_{1:n}, E_{1:n}) =$$

$$P(X_1)P(E_1|X_1) \prod_{t=2}^N P(X_t|X_{t-1})P(E_t|X_t)$$

# Ghostbusters HMM

- $P(X_1) = \text{uniform}$
- $P(X'|X) = \text{usually move clockwise, but sometimes move in a random direction or stay in place}$
- $P(E|X) = \text{same sensor model as before: red means close, green means far away.}$



1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

$P(X_1)$

1/6	1/6	1/2
0	1/6	0
0	0	0

$P(X'|X=\langle 1,2 \rangle)$

# HMM Computations

---

- Given
  - joint  $P(X_{1:n}, E_{1:n})$
  - evidence  $E_{1:n} = e_{1:n}$
- Inference problems include:
  - **Filtering**, find  $P(X_t | e_{1:t})$  for all  $t$
  - **Smoothing**, find  $P(X_t | e_{1:n})$  for all  $t$
  - **Most probable explanation**, find
$$x^*_{1:n} = \operatorname{argmax}_{x_{1:n}} P(x_{1:n} | e_{1:n})$$

# Real HMM Examples

---

- **Speech recognition HMMs:**
  - Observations are acoustic signals (continuous valued)
  - States are specific positions in specific words (so, tens of thousands)
- **Machine translation HMMs:**
  - Observations are words (tens of thousands)
  - States are translation options
- **Robot tracking:**
  - Observations are range readings (continuous)
  - States are positions on a map (continuous)

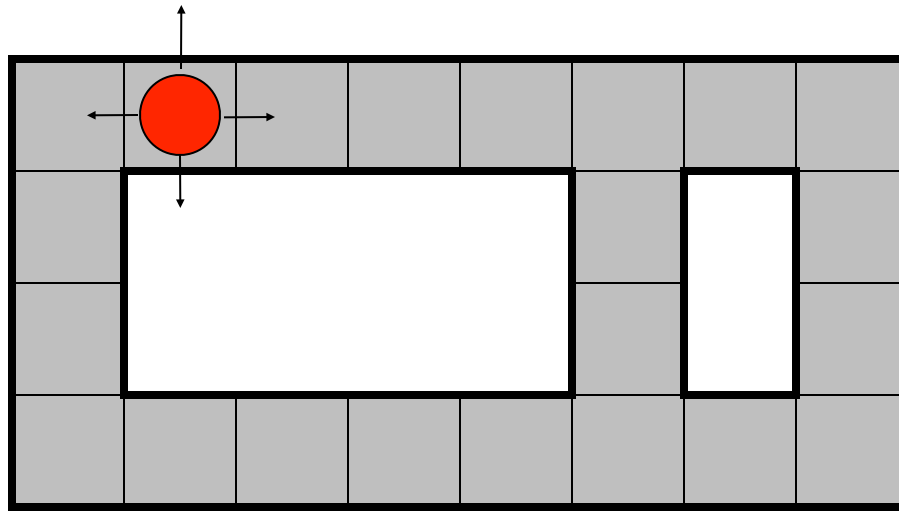
# Filtering / Monitoring

---

- Filtering, or monitoring, is the task of tracking the distribution  $B(X)$  (the belief state) over time
- We start with  $B(X)$  in an initial setting, usually uniform
- As time passes, or we get observations, we update  $B(X)$
- The Kalman filter was invented in the 60's and first implemented as a method of trajectory estimation for the Apollo program

# Example: Robot Localization

*Example from  
Michael Pfeiffer*



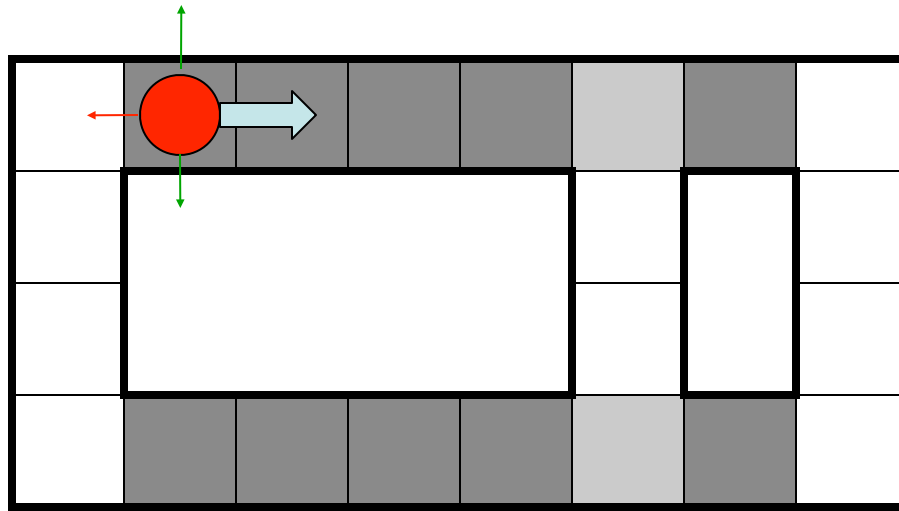
t=0

Sensor model: never more than 1 mistake

Motion model: may not execute action with small prob.

# Example: Robot Localization

---



Prob

0

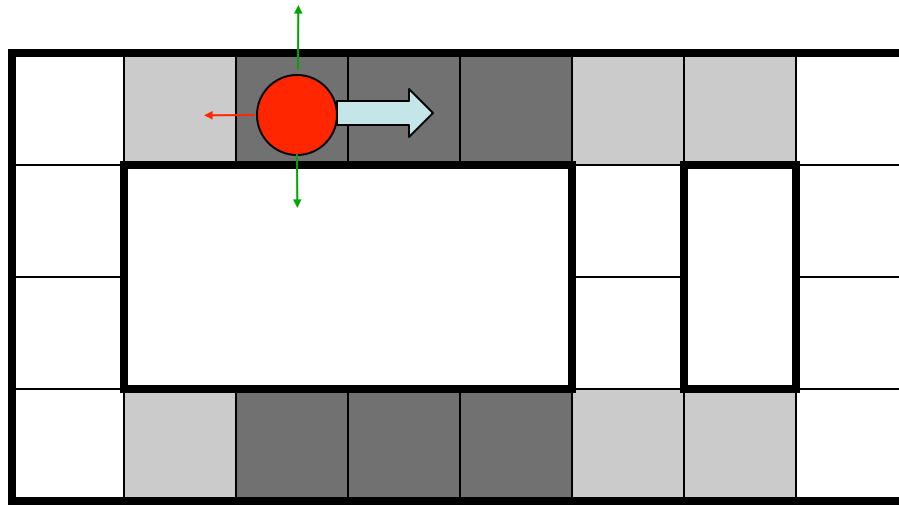
1

$t=1$



# Example: Robot Localization

---



Prob



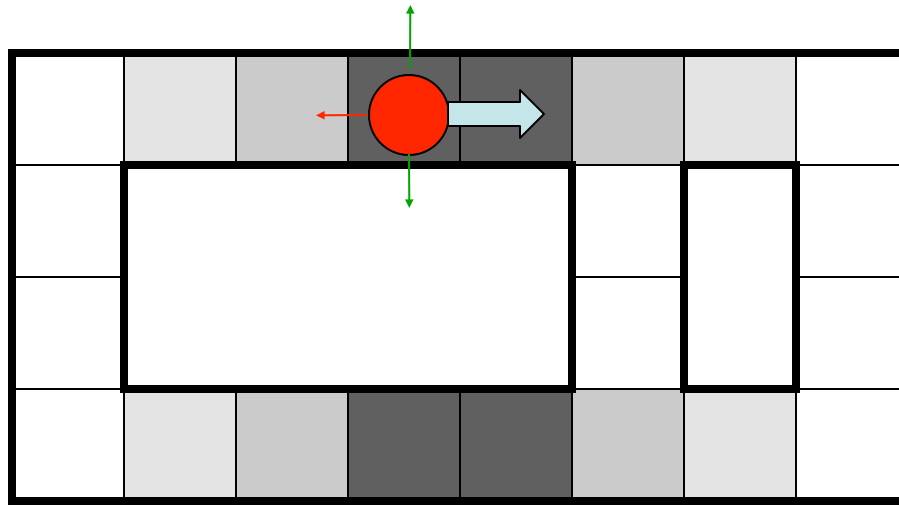
0

1

t=2

# Example: Robot Localization

---



Prob



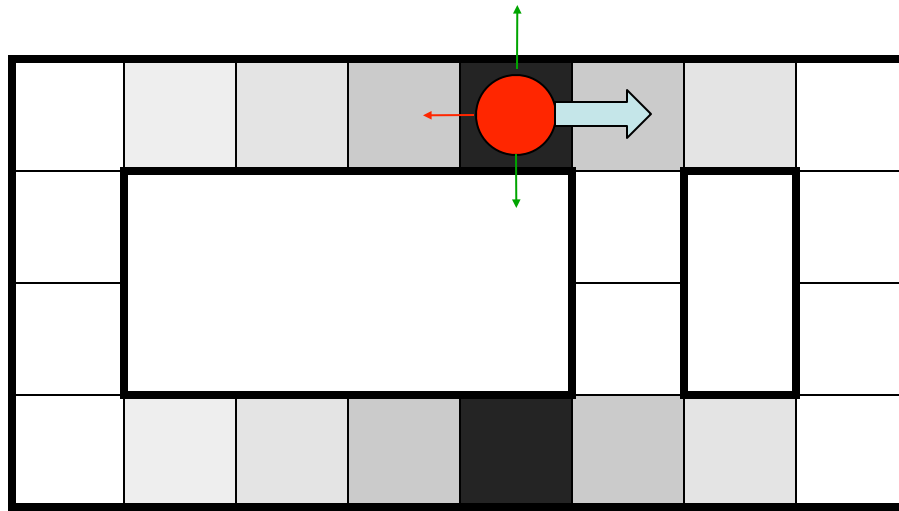
0

1

t=3

# Example: Robot Localization

---



Prob

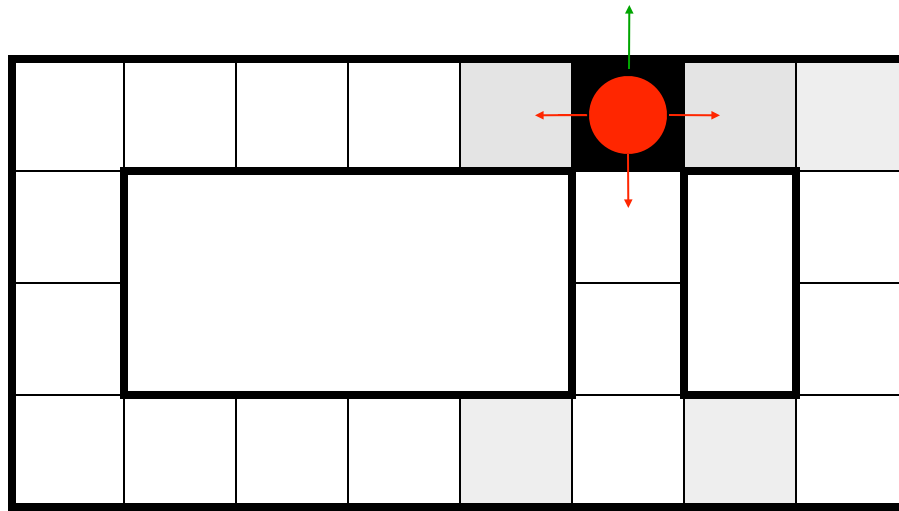
0

1

$t=4$

# Example: Robot Localization

---



Prob

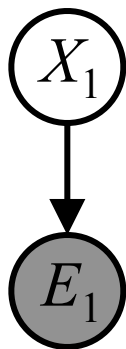
0

1

t=5

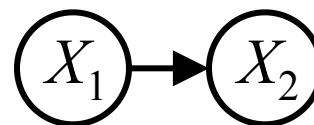
# Inference Recap: Simple Cases

---



$$P(X_1|e_1)$$

$$\begin{aligned} P(x_1|e_1) &= P(x_1, e_1)/P(e_1) \\ &\propto_{X_1} P(x_1, e_1) \\ &= P(x_1)P(e_1|x_1) \end{aligned}$$



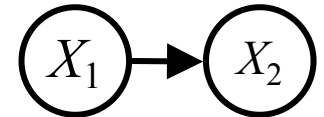
$$P(X_2)$$

$$\begin{aligned} P(x_2) &= \sum_{x_1} P(x_1, x_2) \\ &= \sum_{x_1} P(x_1)P(x_2|x_1) \end{aligned}$$

# Online Belief Updates

---

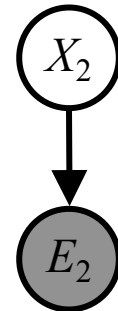
- Every time step, we start with current  $P(X \mid \text{evidence})$
- We update for time:



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

- We update for evidence:

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



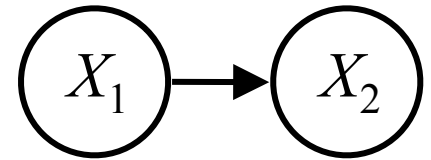
# Passage of Time

---

- Assume we have current belief  $P(X \mid \text{evidence to date})$

$$B(X_t) = P(X_t | e_{1:t})$$

- Then, after one time step passes:



$$P(X_{t+1} | e_{1:t}) = \sum_{x_t} P(X_{t+1} | x_t) P(x_t | e_{1:t})$$

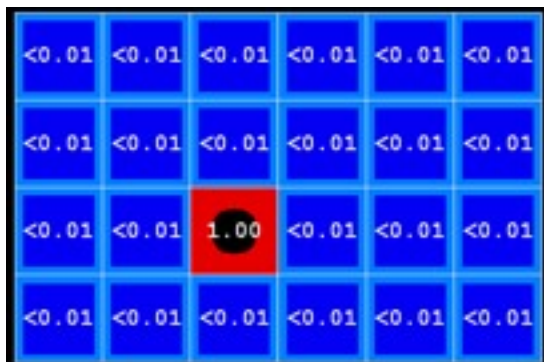
- Or, compactly:

$$B'(X') = \sum_x P(X' | x) B(x)$$

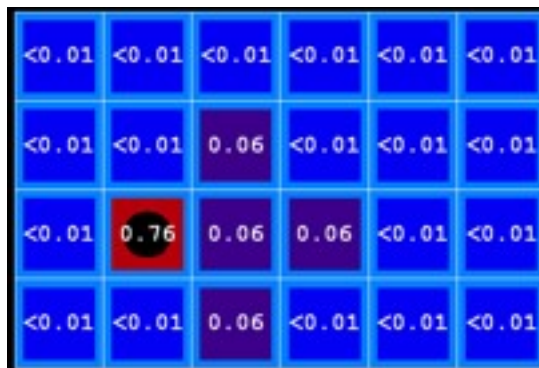
- Basic idea: beliefs get “pushed” through the transitions
  - With the “B” notation, we have to be careful about what time step  $t$  the belief is about, and what evidence it includes

# Example: Passage of Time

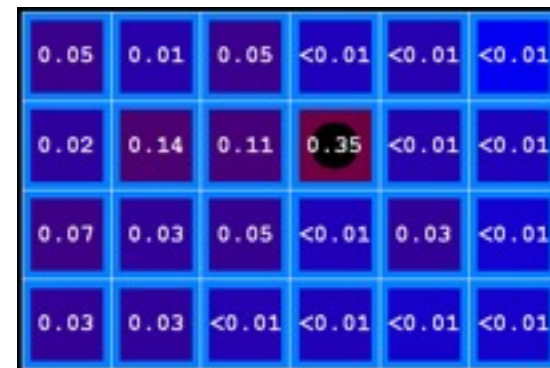
- As time passes, uncertainty “accumulates”



T = 1



T = 2



T = 5

$$B'(X') = \sum_x P(X'|x)B(x)$$

Transition model: ghosts usually go clockwise



# Observation

---

- Assume we have current belief  $P(X \mid \text{previous evidence})$ :

$$B'(X_{t+1}) = P(X_{t+1} | e_{1:t})$$

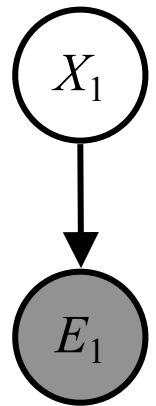
- Then:

$$P(X_{t+1} | e_{1:t+1}) \propto P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t})$$

- Or:

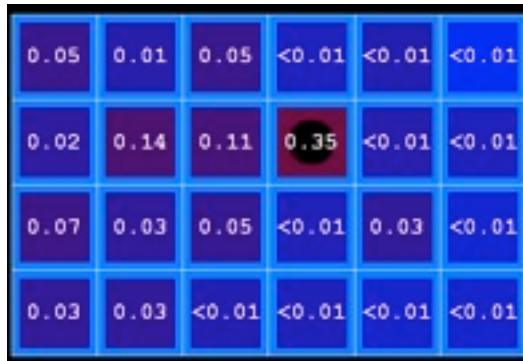
$$B(X_{t+1}) \propto P(e | X) B'(X_{t+1})$$

- Basic idea: beliefs reweighted by likelihood of evidence
- Unlike passage of time, we have to renormalize

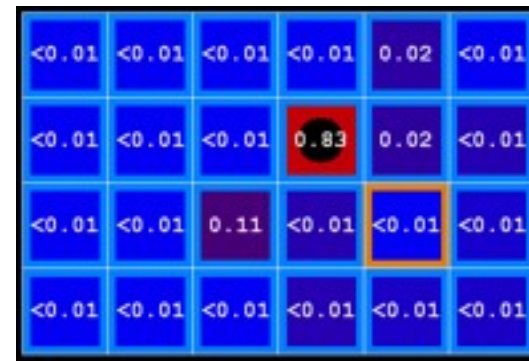


# Example: Observation

- As we get observations, beliefs get reweighted, uncertainty “decreases”



Before observation



After observation

$$B(X) \propto P(e|X)B'(X)$$

# The Forward Algorithm

---

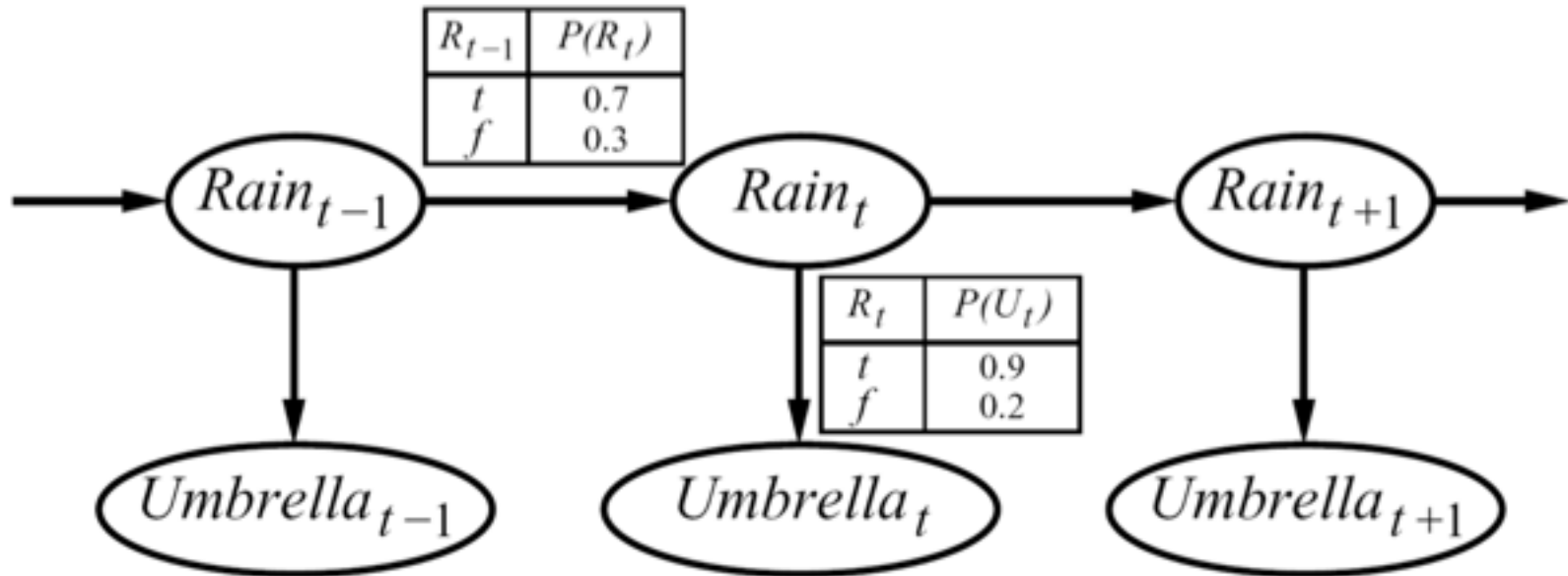
- We to know:  $B_t(X) = P(X_t|e_{1:t})$
- We can derive the following updates

$$\begin{aligned} P(x_t|e_{1:t}) &\propto_X P(x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, e_{1:t-1}) P(x_t|x_{t-1}) P(e_t|x_t) \\ &= P(e_t|x_t) \sum_{x_{t-1}} P(x_t|x_{t-1}) P(x_{t-1}, e_{1:t-1}) \end{aligned}$$

- To get  $B_t(X)$ , compute each entry and normalize

# Example: Run the Filter

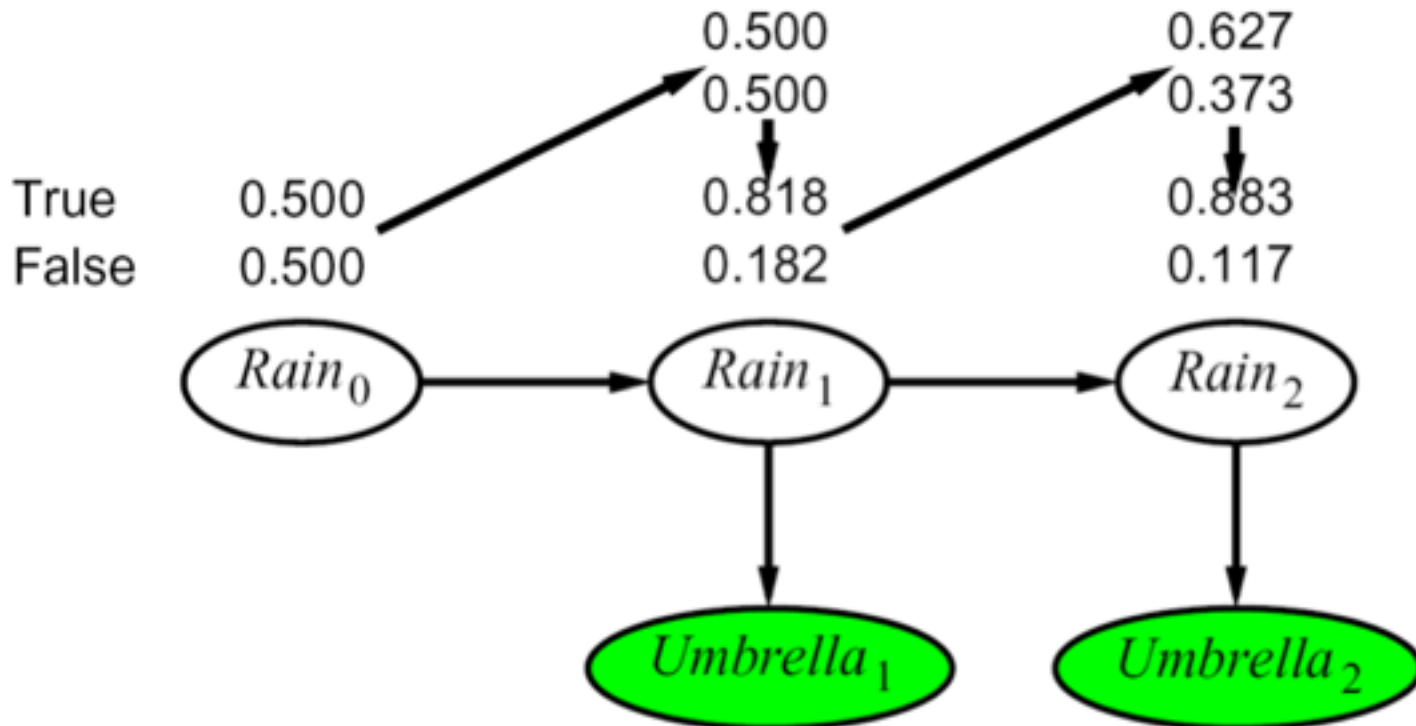
---



- An HMM is defined by:
  - Initial distribution:  $P(X_1)$
  - Transitions:  $P(X_t|X_{t-1})$
  - Emissions:  $P(E|X)$

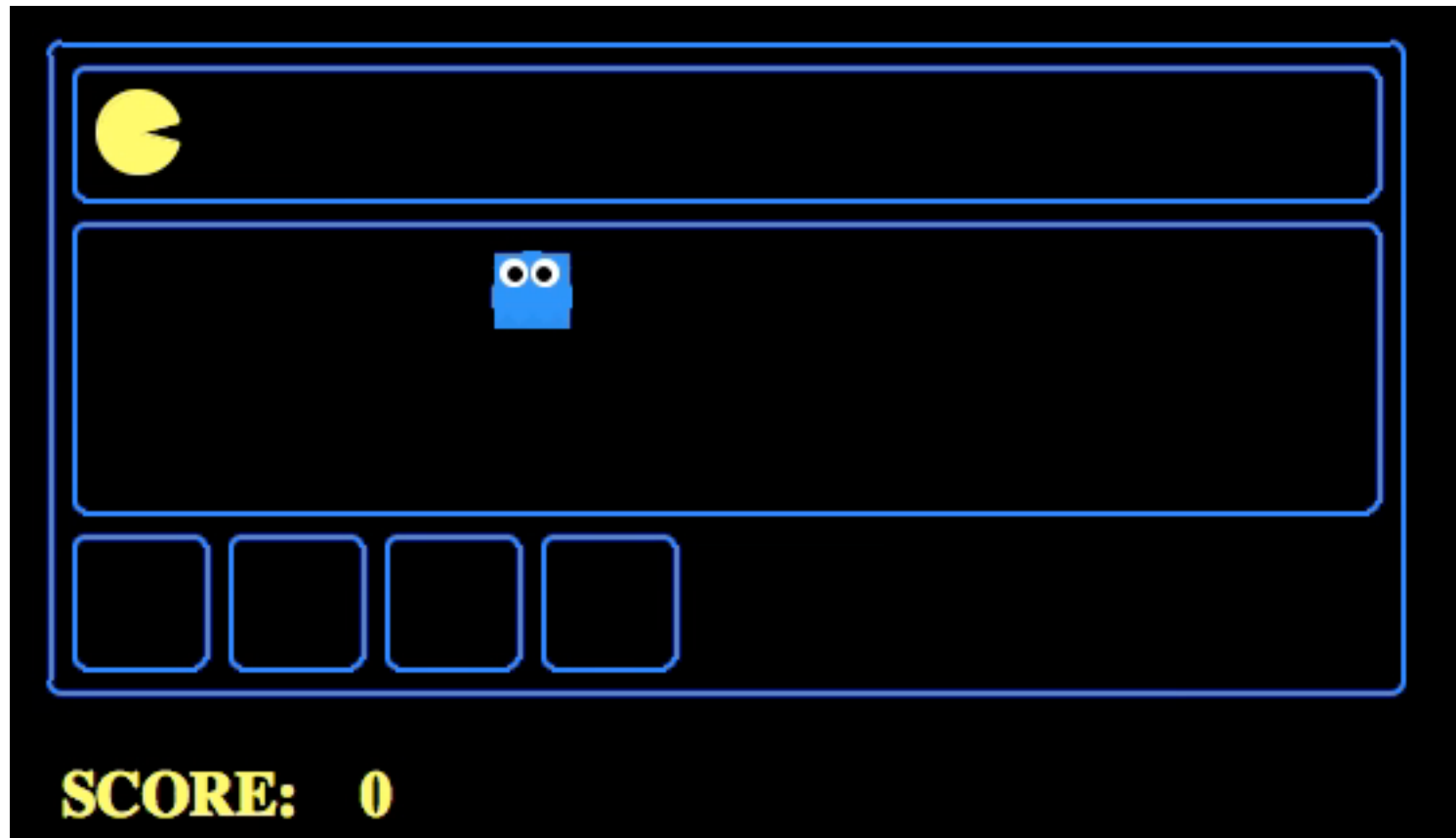
# Example HMM

---



# Example Pac-man

---



# Summary: Filtering

---

- Filtering is the inference process of finding a distribution over  $X_T$  given  $e_1$  through  $e_T$  :  $P( X_T | e_{1:t} )$
- We first compute  $P( X_1 | e_1 )$ :  $P(x_1|e_1) \propto P(x_1) \cdot P(e_1|x_1)$
- For each  $t$  from 2 to  $T$ , we have  $P( X_{t-1} | e_{1:t-1} )$
- **Elapse time:** compute  $P( X_t | e_{1:t-1} )$

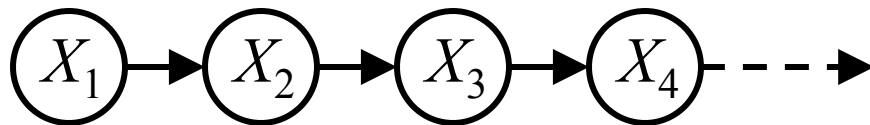
$$P(x_t|e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1}|e_{1:t-1}) \cdot P(x_t|x_{t-1})$$

- **Observe:** compute  $P(X_t | e_{1:t-1}, e_t) = P( X_t | e_{1:t} )$

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

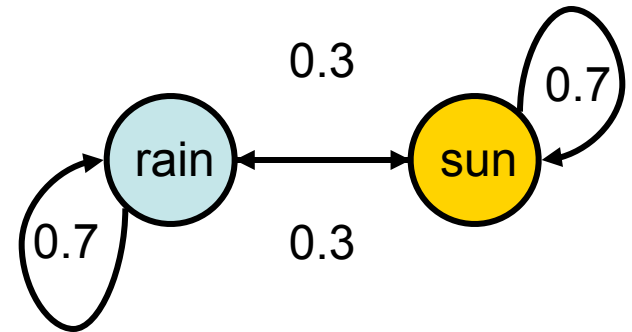
# Recap: Reasoning Over Time

- Stationary Markov models



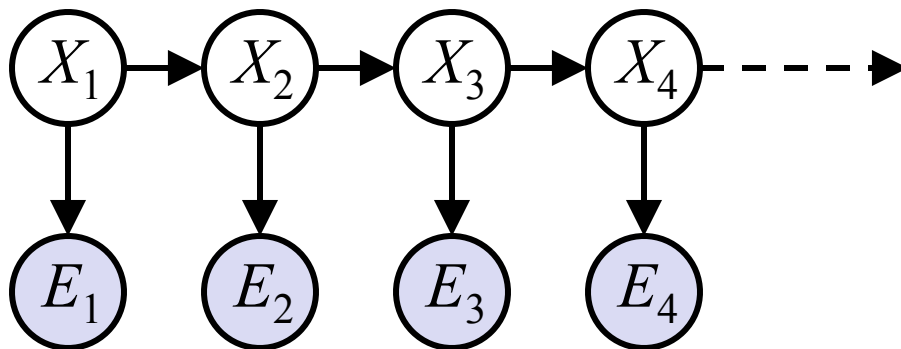
$$P(X_1)$$

$$P(X|X_{-1})$$



$$P(E|X)$$

- Hidden Markov models



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



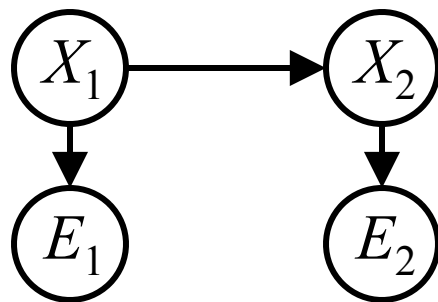
# Recap: Filtering

**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*

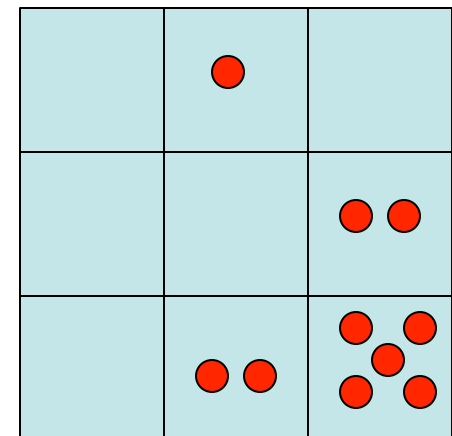
$P(X_2 | E_1 = \text{umbrella})$        $\langle 0.63, 0.37 \rangle$       *Elapse time*

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

# Particle Filtering

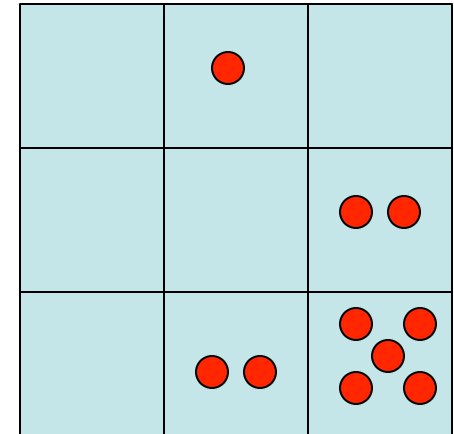
- Sometimes  $|X|$  is too big to use exact inference
  - $|X|$  may be too big to even store  $B(X)$
  - E.g.  $X$  is continuous
  - $|X|^2$  may be too big to do updates
- Solution: approximate inference
  - Track samples of  $X$ , not all values
  - Samples are called particles
  - Time per step is linear in the number of samples
  - But: number needed may be large
  - In memory: list of particles, not states
- This is how robot localization works in practice

0.0	0.1	0.0
0.0	0.0	0.2
0.0	0.2	0.5



# Representation: Particles

- Our representation of  $P(X)$  is now a list of  $N$  particles (samples)
  - Generally,  $N \ll |X|$
  - Storing map from  $X$  to counts would defeat the point
- $P(x)$  approximated by number of particles with value  $x$ 
  - So, many  $x$  will have  $P(x) = 0!$
  - More particles, more accuracy
- For now, all particles have a weight of 1



Particles:

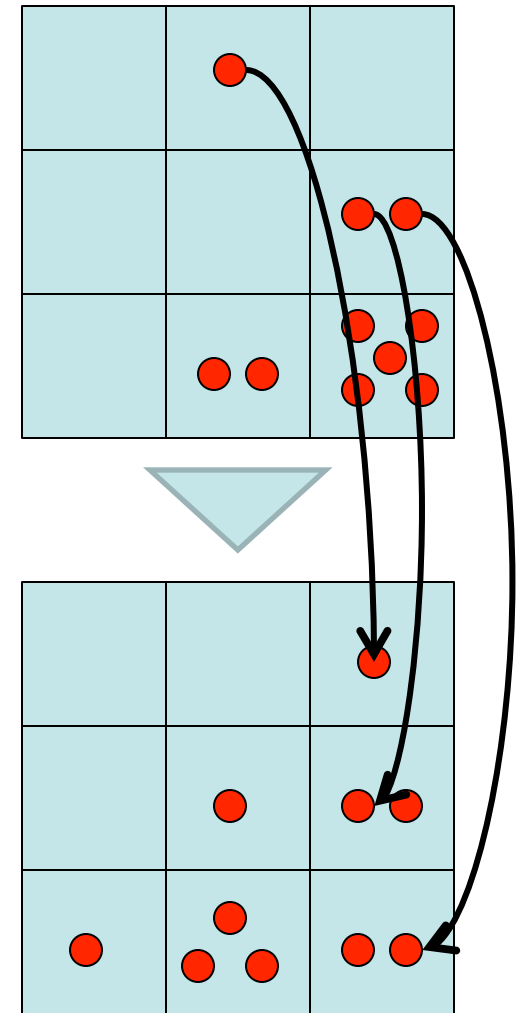
(3,3)  
(2,3)  
(3,3)  
(3,2)  
(3,3)  
(3,2)  
(2,1)  
(3,3)  
(3,3)  
(2,1)

# Particle Filtering: Elapse Time

- Each particle is moved by sampling its next position from the transition model

$$x' = \text{sample}(P(X'|x))$$

- This is like prior sampling – samples' frequencies reflect the transition probs
  - Here, most samples move clockwise, but some move in another direction or stay in place
- This captures the passage of time
    - If we have enough samples, close to the exact values before and after (consistent)



# Particle Filtering: Observe

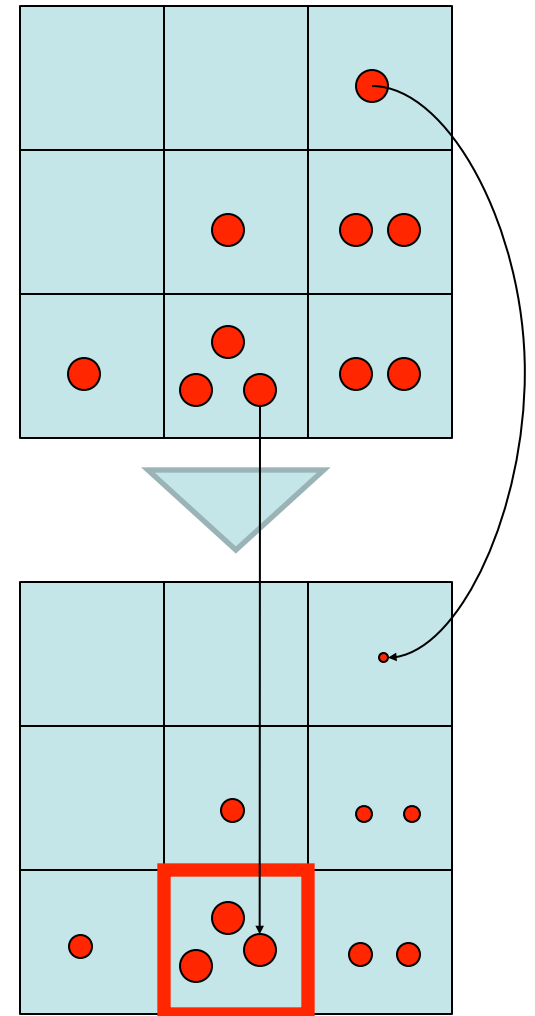
- Slightly trickier:

- We don't sample the observation, we fix it
- We weight our samples based on the evidence

$$w(x) = P(e|x)$$

$$B(X) \propto P(e|X)B'(X)$$

- Note that, as before, the weights/probabilities don't sum to one, since most have been downweighted (in fact they sum to an approximation of  $P(e)$ )

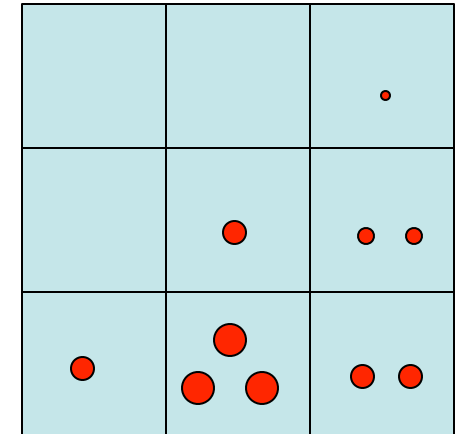


# Particle Filtering: Resample

- Rather than tracking weighted samples, we resample
- N times, we choose from our weighted sample distribution (i.e. draw with replacement)
- This is equivalent to renormalizing the distribution
- Now the update is complete for this time step, continue with the next one

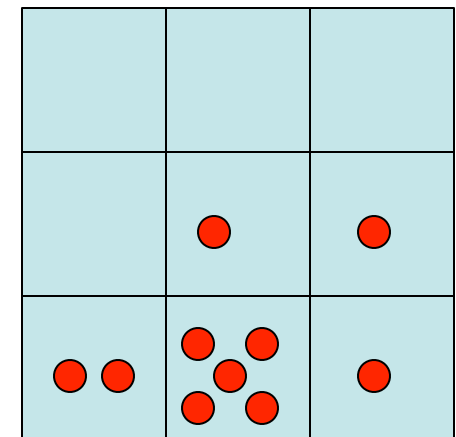
Old Particles:

(3,3)  $w=0.1$   
(2,1)  $w=0.9$   
(2,1)  $w=0.9$   
(3,1)  $w=0.4$   
(3,2)  $w=0.3$   
(2,2)  $w=0.4$   
(1,1)  $w=0.4$   
(3,1)  $w=0.4$   
(2,1)  $w=0.9$   
(3,2)  $w=0.3$



New Particles:

(2,1)  $w=1$   
(2,1)  $w=1$   
(2,1)  $w=1$   
(3,2)  $w=1$   
(2,2)  $w=1$   
(2,1)  $w=1$   
(1,1)  $w=1$   
(3,1)  $w=1$   
(2,1)  $w=1$   
(1,1)  $w=1$



# Summary: Particle Filtering

---

At each time step  $t$ , we have a set of  $N$  particles / samples

- Three step procedure, to move to time  $t+1$ :

1. **Sample transitions:** for each each particle  $x$ , sample next state

$$x' = \text{sample}(P(X'|x))$$

2. **Reweight:** for each particle, compute its weight

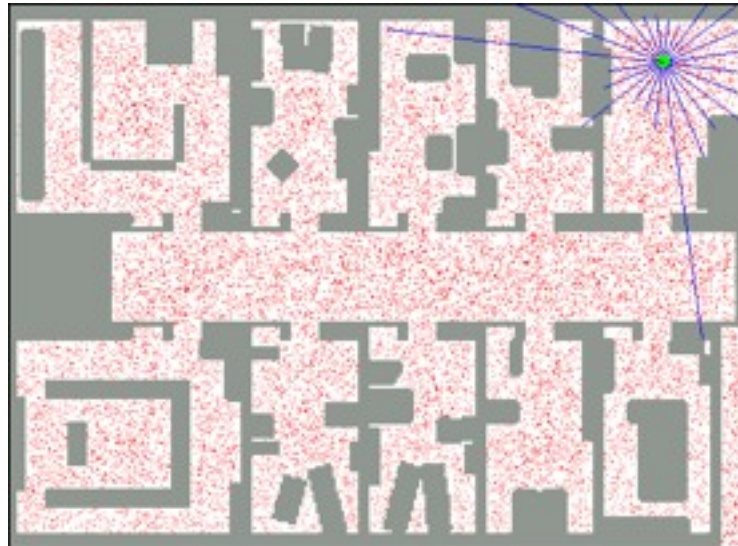
$$w(x) = P(e|x)$$

3. **Resample:** normalize the weights, and sample  $N$  new particles from the resulting distribution over states

# Robot Localization

---

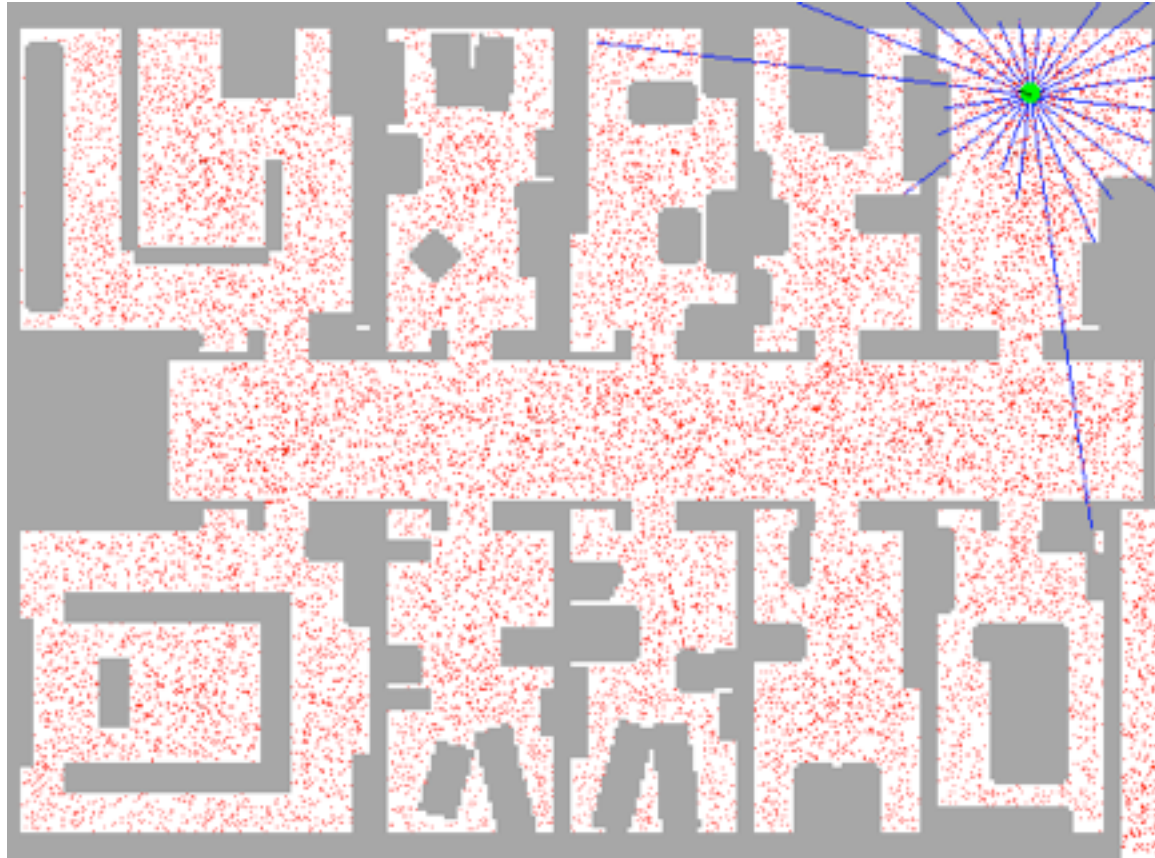
- In robot localization:
  - We know the map, but not the robot's position
  - Observations may be vectors of range finder readings
  - State space and readings are typically continuous (works basically like a very fine grid) and so we cannot store  $B(X)$
  - Particle filtering is a main technique





# Robot Localization

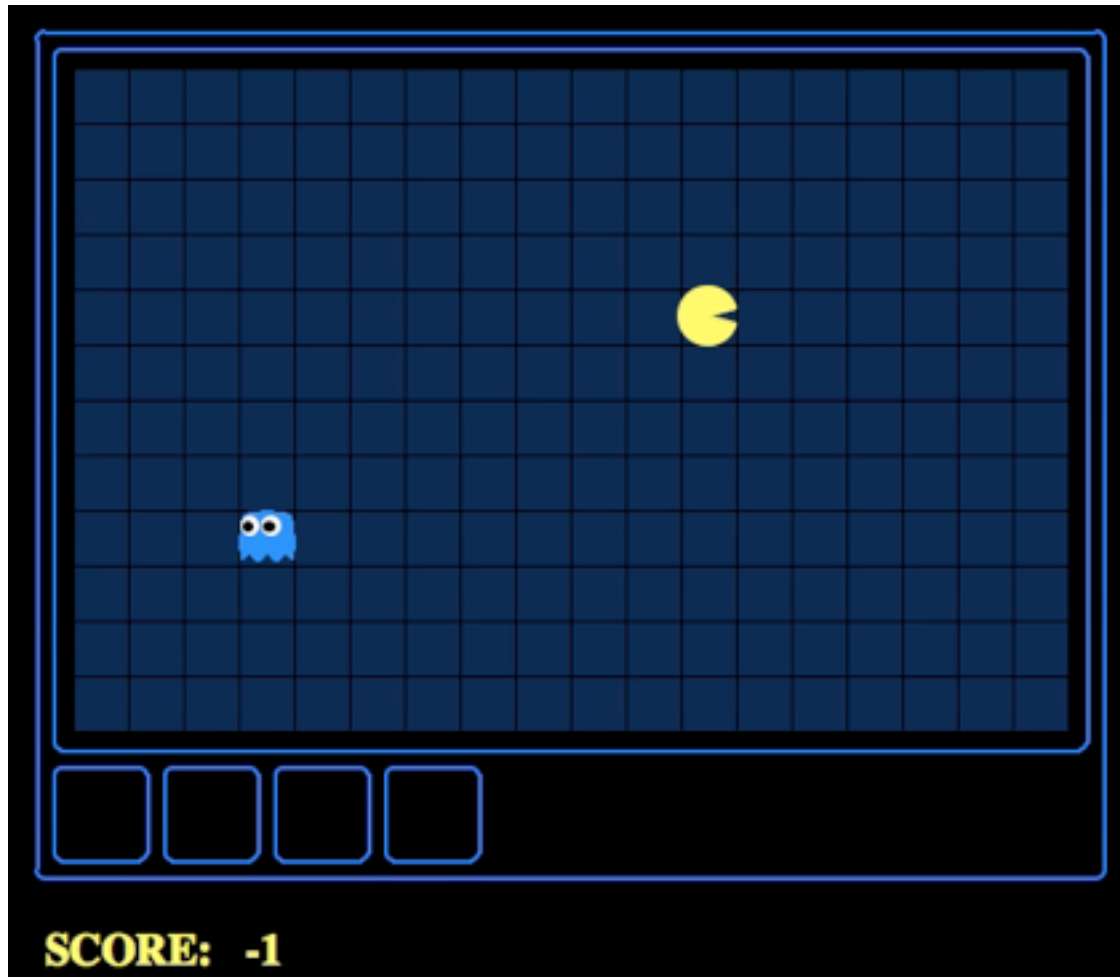
---



# Which Algorithm?

---

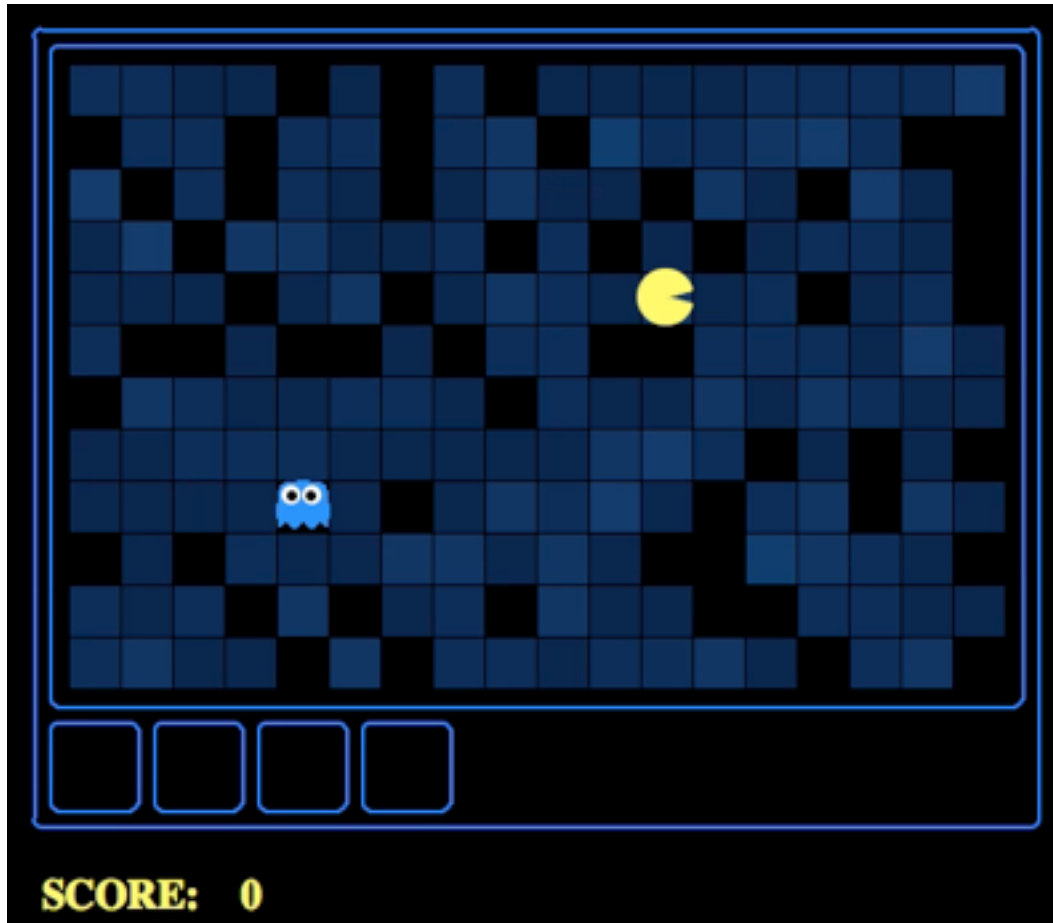
Exact filter, uniform initial beliefs



# Which Algorithm?

---

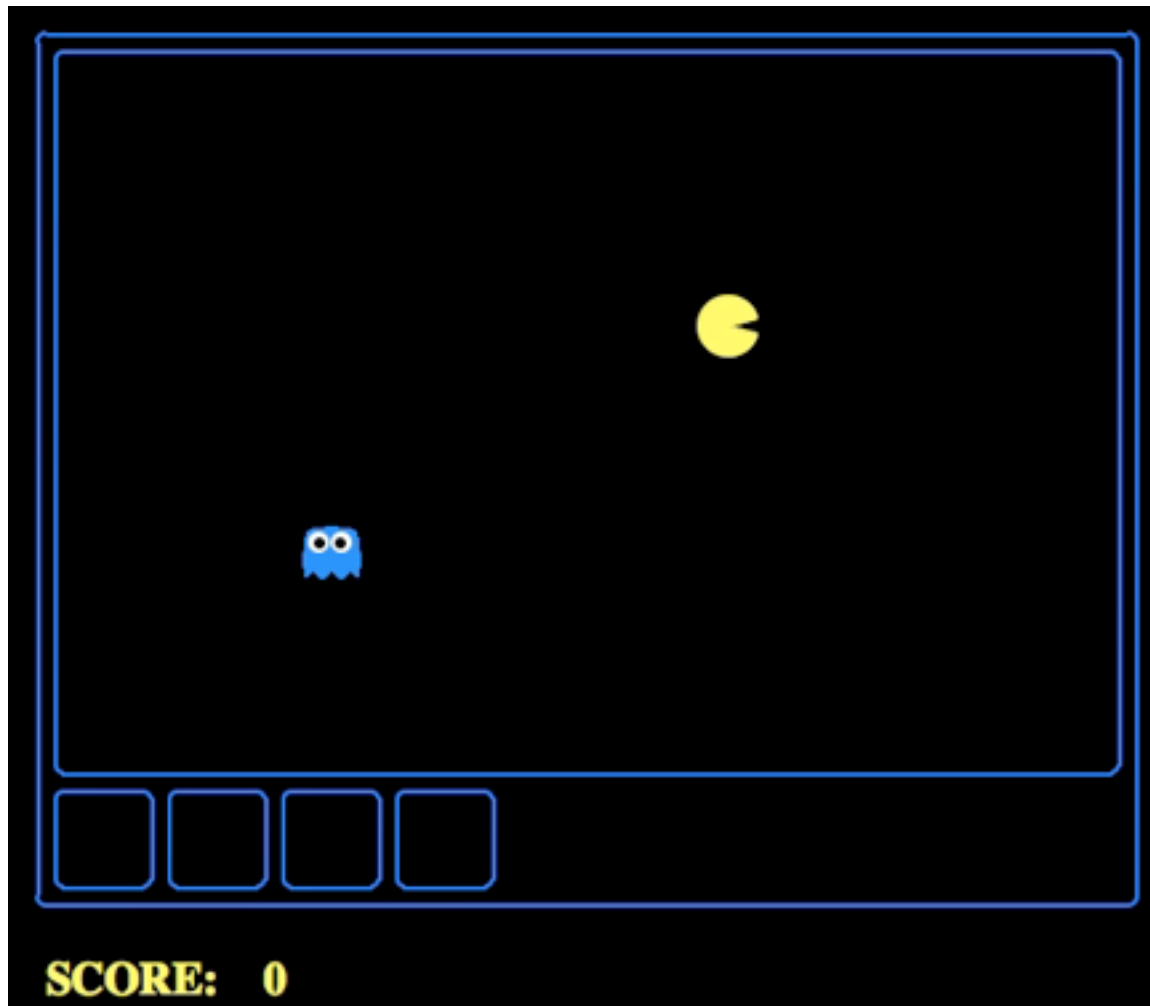
Particle filter, uniform initial beliefs, 300 particles



# Which Algorithm?

---

Particle filter, uniform initial beliefs, 25 particles



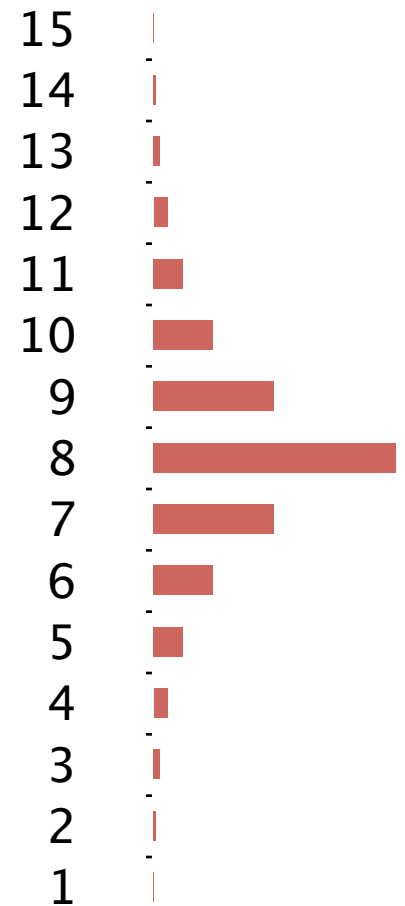
# P4: Ghostbusters

---

- **Plot:** Pacman's grandfather, Grandpac, learned to hunt ghosts for sport.
- He was blinded by his power, but could hear the ghosts' banging and clanging.
- **Transition Model:** All ghosts move randomly, but are sometimes biased
- **Emission Model:** Pacman knows a “noisy” distance to each ghost

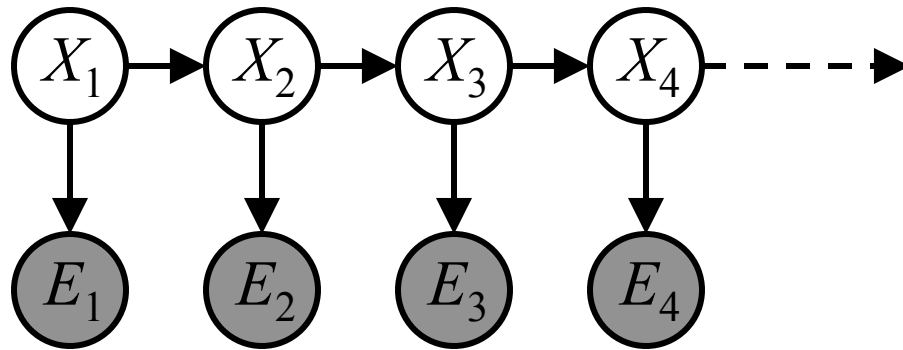
## Noisy distance prob

True distance = 8



# Best Explanation Queries

---

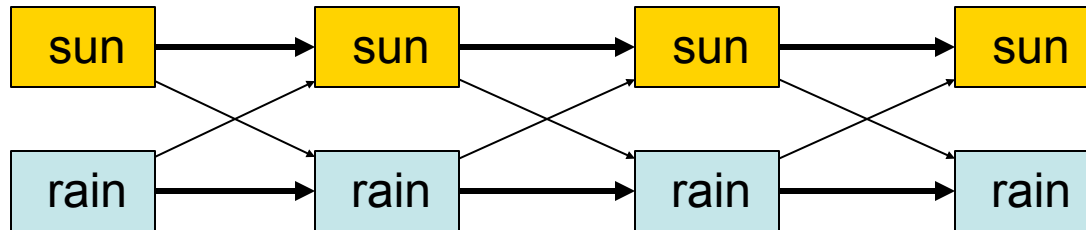


- Query: most likely seq:

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

# Viterbi Algorithm

---



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

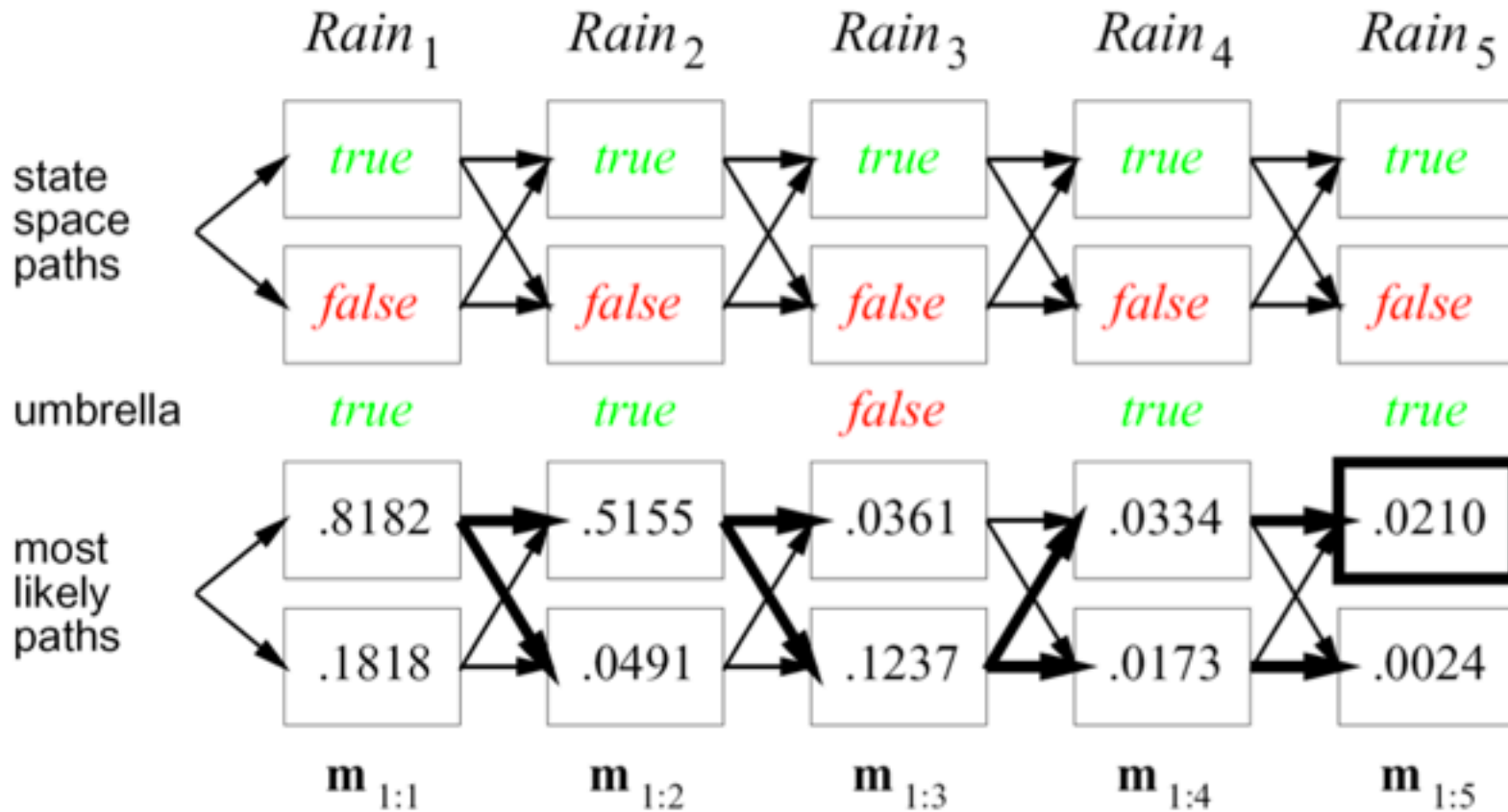
$$m_t[x_t] = \max_{x_{1:t-1}} P(x_{1:t-1}, x_t, e_{1:t})$$

$$= \max_{x_{1:t-1}} P(x_{1:t-1}, e_{1:t-1}) P(x_t|x_{t-1}) P(e_t|x_t)$$

$$= P(e_t|x_t) \max_{x_{t-1}} P(x_t|x_{t-1}) \max_{x_{1:t-2}} P(x_{1:t-1}, e_{1:t-1})$$

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

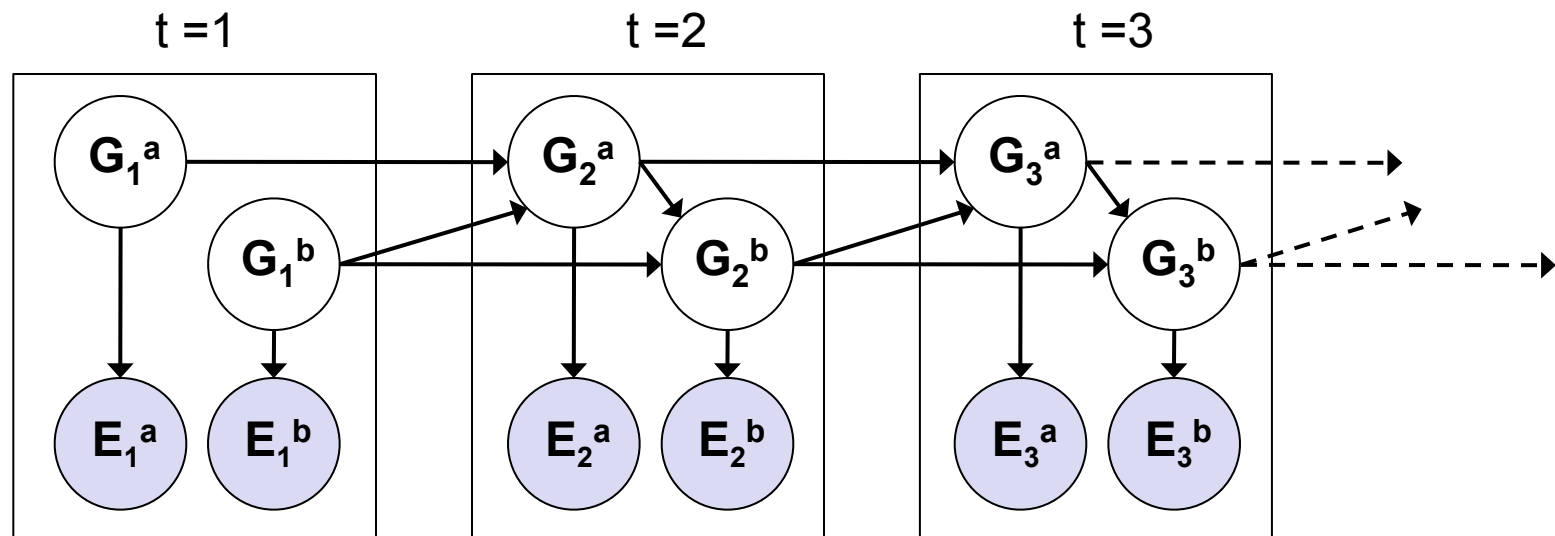
# Example





# 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  $t-1$



- Discrete valued dynamic Bayes nets are also 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