

# TEMPORAL PROBABILITY MODELS

## CHAPTER 15, SECTIONS 1–5

## Outline

- ◇ Time and uncertainty
- ◇ Inference: filtering, prediction, smoothing
- ◇ Hidden Markov models
- ◇ Kalman filters (a brief mention)
- ◇ Dynamic Bayesian networks
- ◇ Particle filtering

# Time and uncertainty

The world changes; we need to track and predict it

Diabetes management vs vehicle diagnosis

Basic idea: copy state and evidence variables for each time step

$X_t$  = set of unobservable state variables at time  $t$   
e.g., *BloodSugar<sub>t</sub>*, *StomachContents<sub>t</sub>*, etc.

$E_t$  = set of observable evidence variables at time  $t$   
e.g., *MeasuredBloodSugar<sub>t</sub>*, *PulseRate<sub>t</sub>*, *FoodEaten<sub>t</sub>*

This assumes **discrete time**; step size depends on problem

Notation:  $X^{a:b} = X_a, X_{a+1}, \dots, X_{b-1}, X_b$

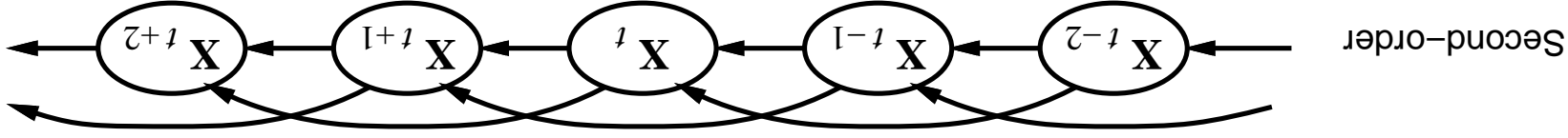
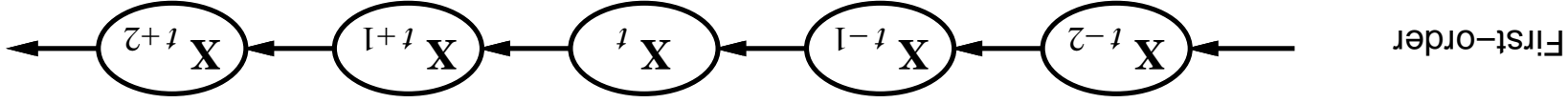
# Markov processes (Markov chains)

Construct a Bayes net from these variables: parents?

Markov assumption:  $X_t$  depends on **bounded** subset of  $X_{0:t-1}$

First-order Markov process:  $P(X_t | X_{0:t-1}) = P(X_t | X_{t-1})$

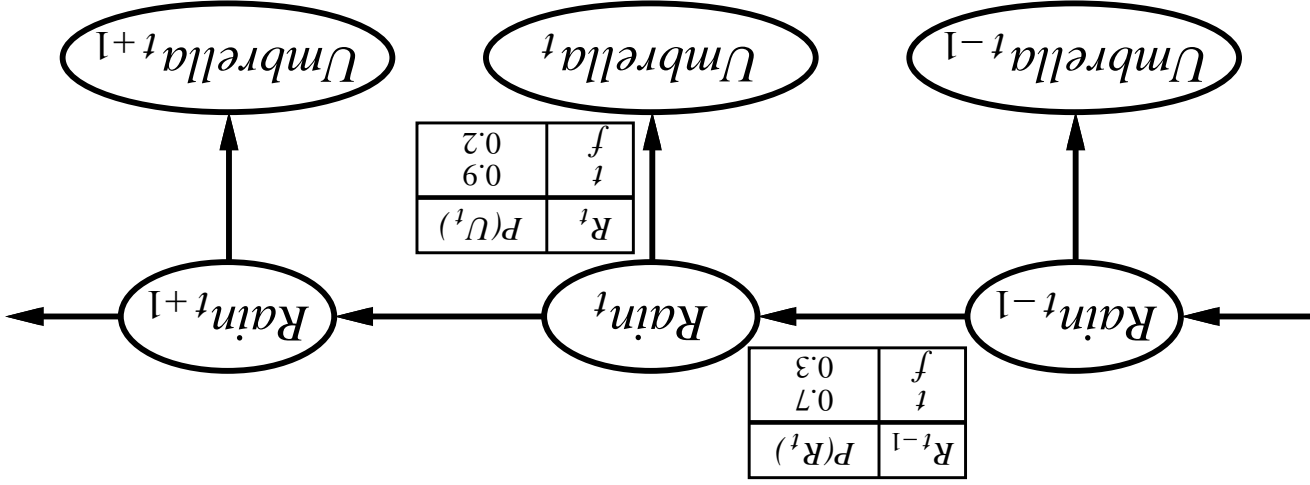
Second-order Markov process:  $P(X_t | X_{0:t-1}) = P(X_t | X_{t-2}, X_{t-1})$



Sensor Markov assumption:  $P(E_t | X_{0:t}, E_{0:t-1}) = P(E_t | X_t)$

Stationary process: transition model  $P(X_t | X_{t-1})$  and sensor model  $P(E_t | X_t)$  fixed for all  $t$

# Example



First-order Markov assumption not exactly true in real world!

Possible fixes:

1. **Increase order** of Markov process
2. **Augment state**, e.g., add *Temp<sub>t</sub>*, *Pressure<sub>t</sub>*

Example: robot motion.

Augment position and velocity with *Battery<sub>t</sub>*

## Inference tasks

Filtering:  $P(\mathbf{X}_t | \mathbf{e}_{1:t})$

belief state—input to the decision process of a rational agent

Prediction:  $P(\mathbf{X}_{t+k} | \mathbf{e}_{1:t})$  for  $k > 0$

evaluation of possible action sequences;

like filtering without the evidence

Smoothing:  $P(\mathbf{X}_k | \mathbf{e}_{1:t})$  for  $0 \leq k < t$

better estimate of past states, essential for learning

Most likely explanation:  $\arg \max_{\mathbf{x}_{1:t}} P(\mathbf{x}_{1:t} | \mathbf{e}_{1:t})$

speech recognition, decoding with a noisy channel

# Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(X_{t+1}|e_{1:t+1}, P(X_t|e_{1:t}))$$

$$P(X_{t+1}|e_{1:t+1}) = P(X_{t+1}|e_{1:t}, e_{t+1})$$

$$= \alpha P(e_{t+1}|X_{t+1}, e_{1:t}) P(X_{t+1}|e_{1:t})$$

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

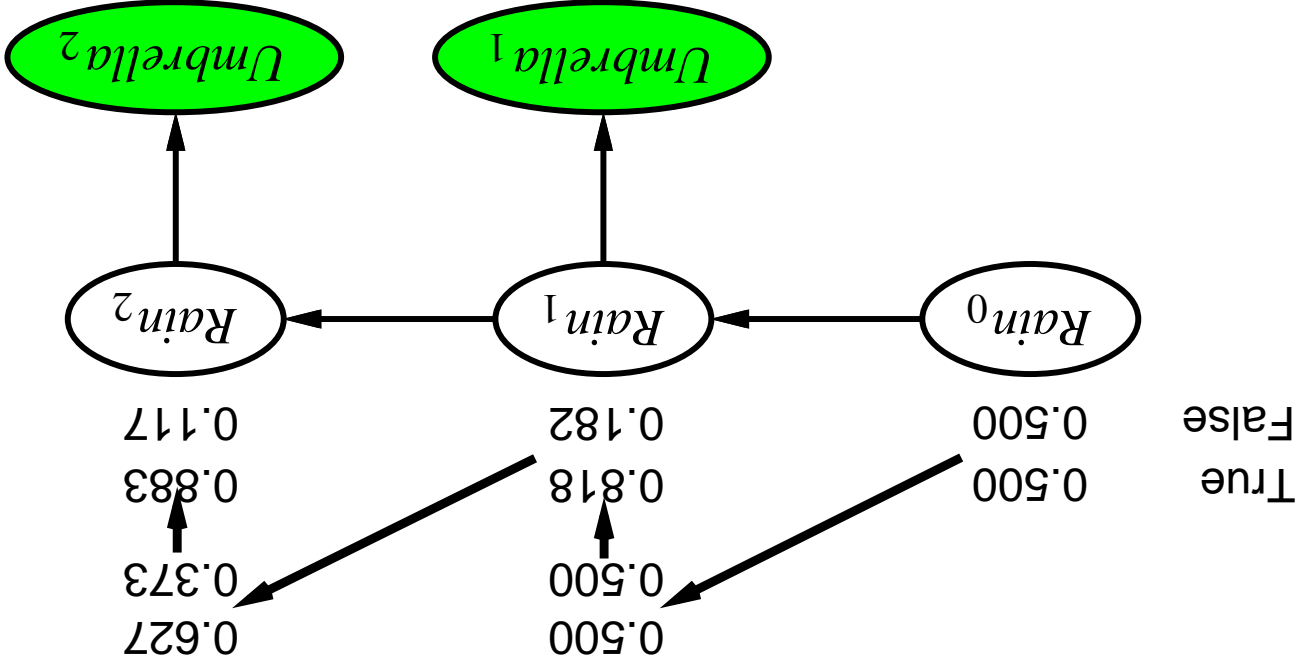
i.e., prediction + estimation. Prediction by summing out  $X_t$ :

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

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

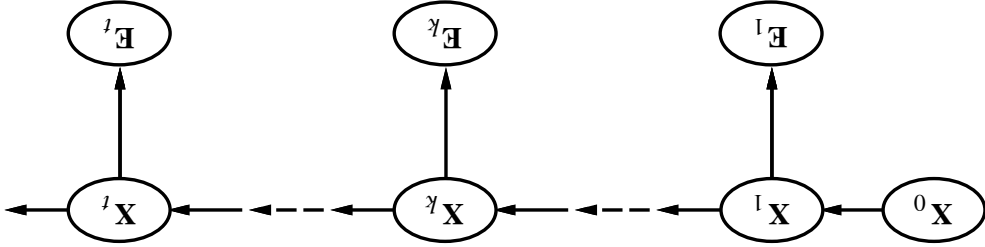
$f_{1:t+1} = \text{FORWARD}(f_{1:t}, e_{t+1})$  where  $f_{1:t} = P(X_t|e_{1:t})$   
 Time and space **constant** (independent of  $t$ )

Filtering example





# Smoothing



Divide evidence  $e_{1:t}$  into  $e_{1:k}, e_{k+1:t}$ :

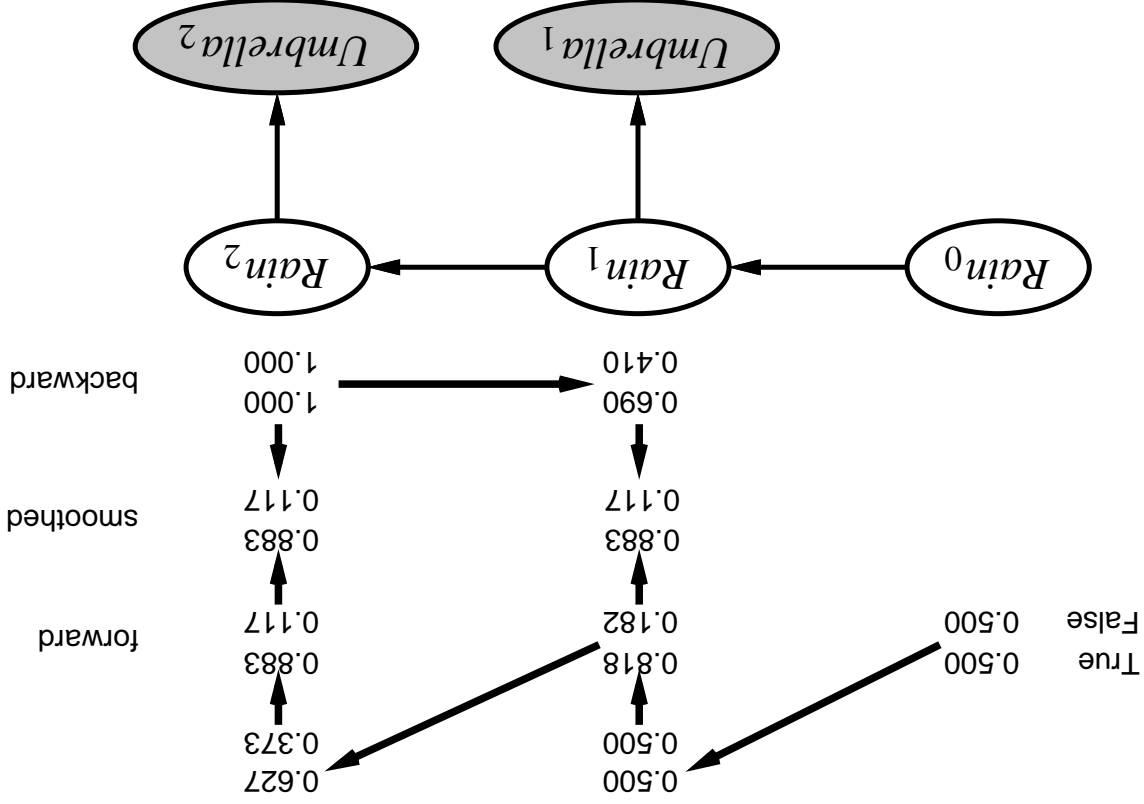
$$\begin{aligned}
 P(\mathbf{X}_k | e_{1:t}) &= P(\mathbf{X}_k | e_{1:k}, e_{k+1:t}) \\
 &= \alpha P(\mathbf{X}_k | e_{1:k}) P(e_{k+1:t} | \mathbf{X}_k) \\
 &= \alpha P(\mathbf{X}_k | e_{1:k}) P(e_{k+1:t} | \mathbf{X}_k) \\
 &= \alpha f_{1:k} b_{k+1:t}
 \end{aligned}$$

Backward message computed by a backwards recursion:

$$\begin{aligned}
 P(e_{k+1:t} | \mathbf{X}_k) &= \sum_{\mathbf{x}_{k+1}} P(e_{k+1:t} | \mathbf{X}_k, \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} | \mathbf{X}_k) \\
 &= \sum_{\mathbf{x}_{k+1}} P(e_{k+1:t} | \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} | \mathbf{X}_k) \\
 &= \sum_{\mathbf{x}_{k+1}} P(e_{k+1:t} | \mathbf{x}_{k+1}) P(e_{k+2:t} | \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} | \mathbf{X}_k)
 \end{aligned}$$

Forward-backward algorithm: cache forward messages along the way  
 Time linear in  $t$  (polytree inference), space  $O(t|F|)$

# Smoothing example



## Most likely explanation

Most likely sequence  $\neq$  sequence of most likely states!!!

Most likely path to each  $\mathbf{x}_{t+1}$

= most likely path to **some**  $\mathbf{x}_t$  plus one more step

$$\max_{\mathbf{x}_1 \dots \mathbf{x}_t} P(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \max_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t) \max_{\mathbf{x}_1 \dots \mathbf{x}_{t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t | \mathbf{e}_{1:t})$$

Identical to filtering, except  $\mathbf{f}_{1:t}$  replaced by

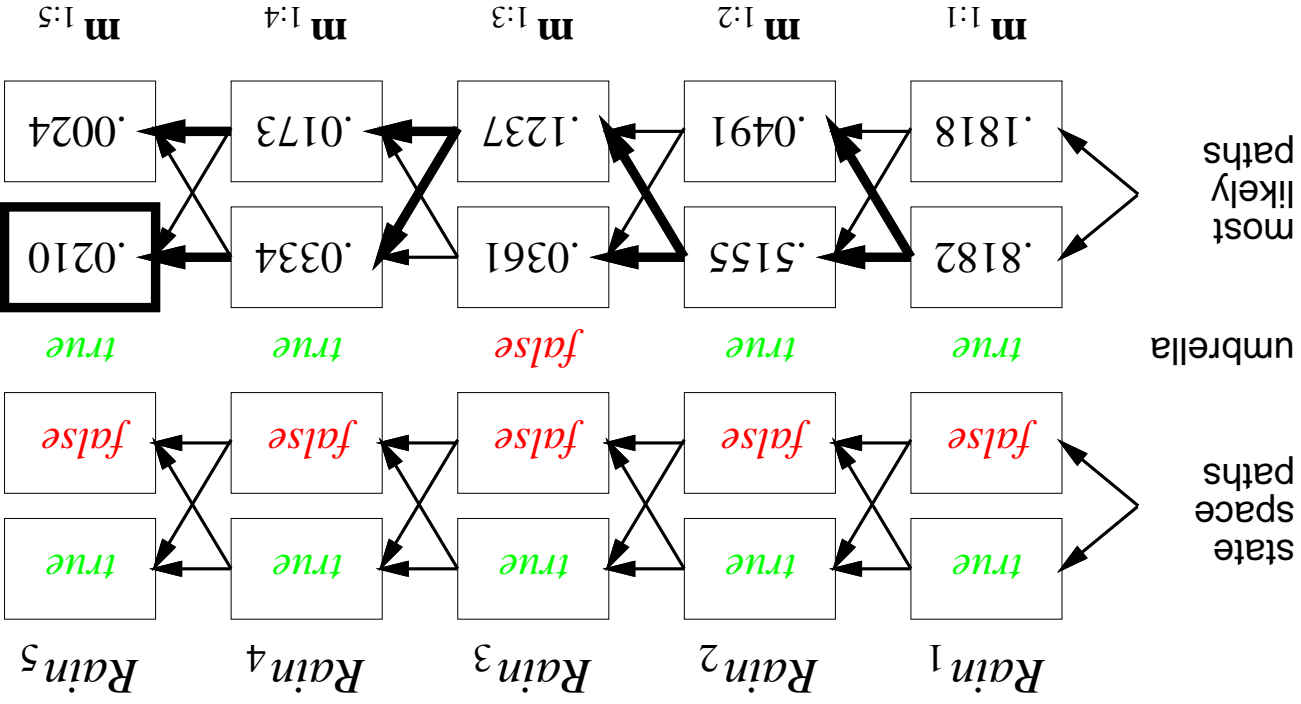
$$\mathbf{m}_{1:t} = \max_{\mathbf{x}_1 \dots \mathbf{x}_{t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{X}_t | \mathbf{e}_{1:t}),$$

i.e.,  $\mathbf{m}_{1:t} (?)$  gives the probability of the most likely path to state  $t$ .

Update has sum replaced by max, giving the Viterbi algorithm:

$$\mathbf{m}_{1:t+1} = P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \max_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t) \mathbf{m}_{1:t}$$

# Viterbi example



# Hidden Markov models

$X_t$  is a single, discrete variable (usually  $E_t$  is too)  
Domain of  $X_t$  is  $\{1, \dots, S\}$

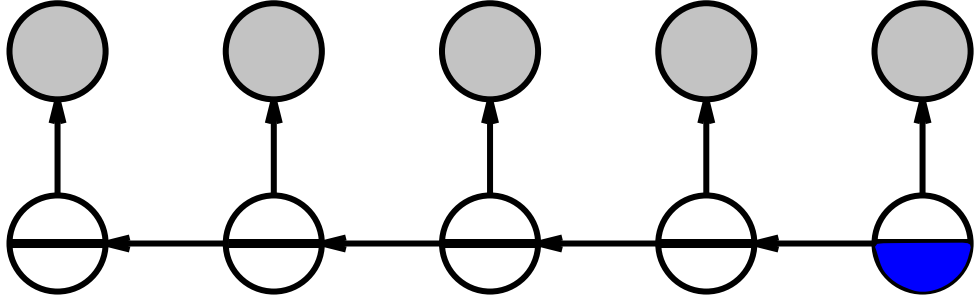
Transition matrix  $T_{ij} = P(X_t = j | X_{t-1} = i)$ , e.g.,  $\begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}$

Sensor matrix  $O_t$  for each time step, diagonal elements  $P(e_t | X_t = i)$   
e.g., with  $U_1 = true$ ,  $O_1 = \begin{pmatrix} 0 & 0.2 \\ 0.9 & 0 \end{pmatrix}$

Forward and backward messages as column vectors:

$$f_{1:t+1} = \alpha O_{t+1} T_{1:t}^T f_{1:t}$$
$$b_{k+1:t} = T O_{k+1} b_{k+2:t}$$

Forward-backward algorithm needs time  $O(S^2t)$  and space  $O(St)$



Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$

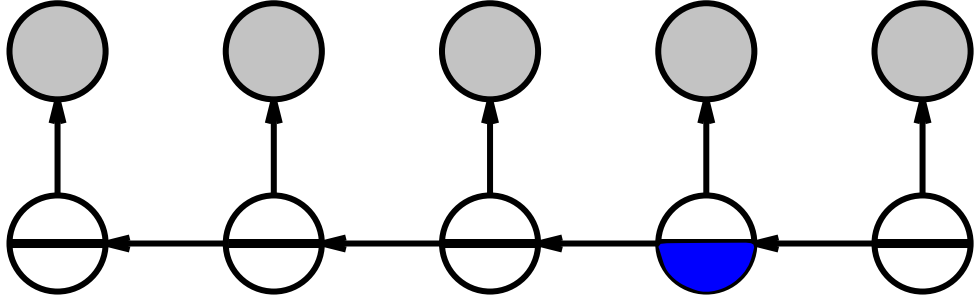
$$f_{1:t+1} = \alpha \mathbf{O}_{t+1}^T \mathbf{T}_{1:t} f_{1:t}$$

$$\mathbf{O}_{t+1}^{-1} f_{1:t+1} = \alpha \mathbf{T}_{1:t}^T f_{1:t}$$

$$\alpha' (\mathbf{T}_{1:t}^{-1} \mathbf{O}_{t+1}^{-1} f_{1:t+1}) = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm



Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$

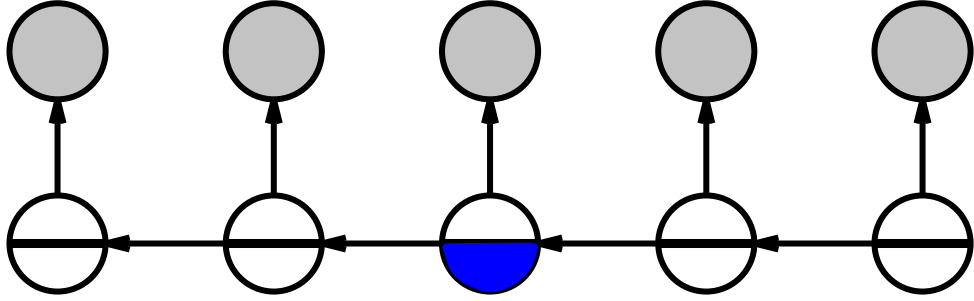
$$f_{1:t+1} = \alpha O_{t+1}^T T_{1:t} f_{1:t}$$

$$O_{t+1}^{-1} f_{1:t+1} = \alpha T_{1:t}^T f_{1:t}$$

$$\alpha (T_{1:t}^{-1} O_{t+1}^{-1} f_{1:t+1}) = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm



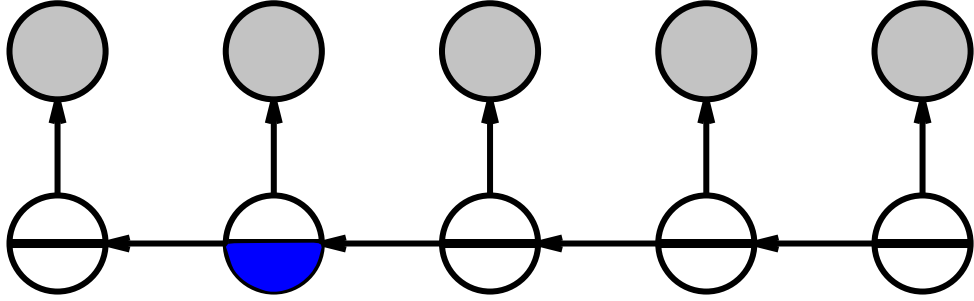
Algorithm: forward pass computes  $f_t$ , backward pass does  $b_t, b?$

$$\begin{aligned}
 f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t} \\
 \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t} \\
 \alpha'(\mathbf{T}_{1:t+1}^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1}) &= f_{1:t}
 \end{aligned}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm





Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$

$$\begin{aligned}
 f_{1:t+1} &= \alpha \mathbf{T}_{t+1}^T \mathbf{T}_{1:t} \\
 \mathbf{O}_{t+1}^{-1} f_{1:t+1} &= \alpha \mathbf{T}_{t+1}^T f_{1:t} \\
 \alpha' (\mathbf{T}_{t+1}^{-1})^{-1} \mathbf{O}_{t+1}^{-1} f_{1:t+1} &= f_{1:t}
 \end{aligned}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

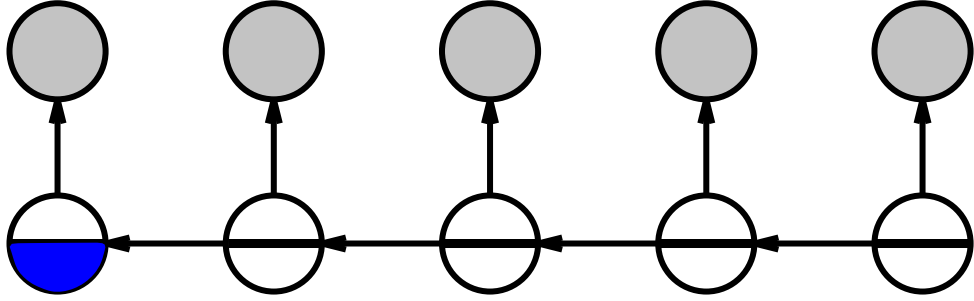
# Country dance algorithm

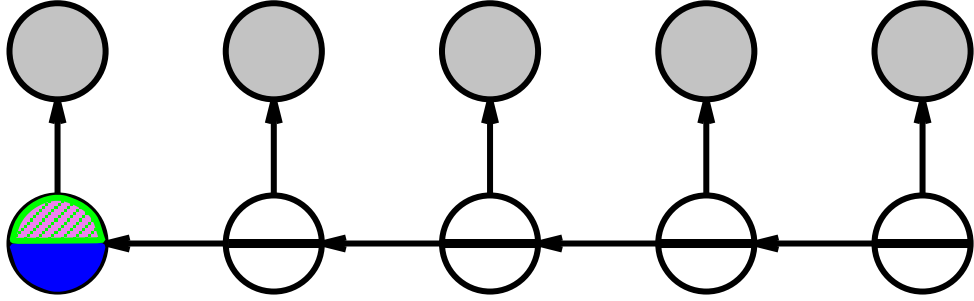
# Country dance algorithm

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\begin{aligned}
 f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^\top f_{1:t} \\
 \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^\top f_{1:t} \\
 \alpha_{t+1} (\mathbf{T}_{1:t}^\top)^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} &= f_{1:t}
 \end{aligned}$$

Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$





Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$

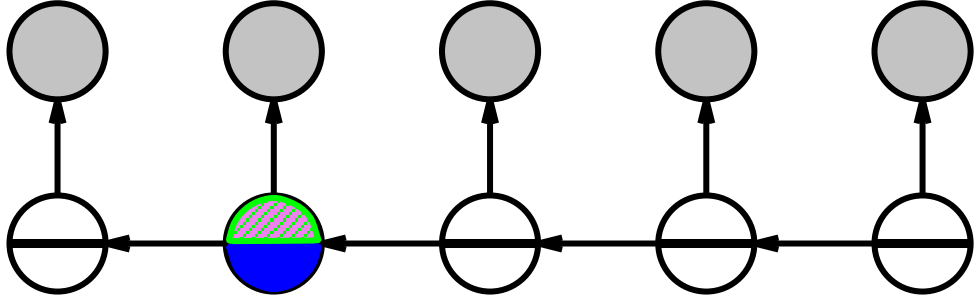
$$f_{1:t+1} = \alpha \mathbf{T}_{1:t}^T f_{1:t}$$

$$\mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = \alpha \mathbf{T}_{1:t}^T f_{1:t}$$

$$\alpha'(\mathbf{T}^T)^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm



Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$

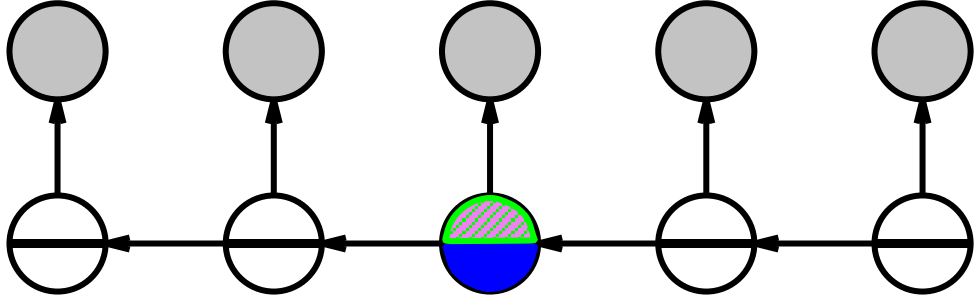
$$f_{1:t+1} = \alpha \mathbf{T}_{1:t}^T f_{1:t}$$

$$\mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = \alpha \mathbf{T}_{1:t}^T f_{1:t}$$

$$\alpha' (\mathbf{T}_{1:t}^T)^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm



Algorithm: forward pass computes  $f_t$ , backward pass does  $b_t, b_t$

$$f_{1:t+1} = \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t}$$

$$\mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t}$$

$$\alpha_{t+1} (\mathbf{T}_{1:t+1}^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1}) = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

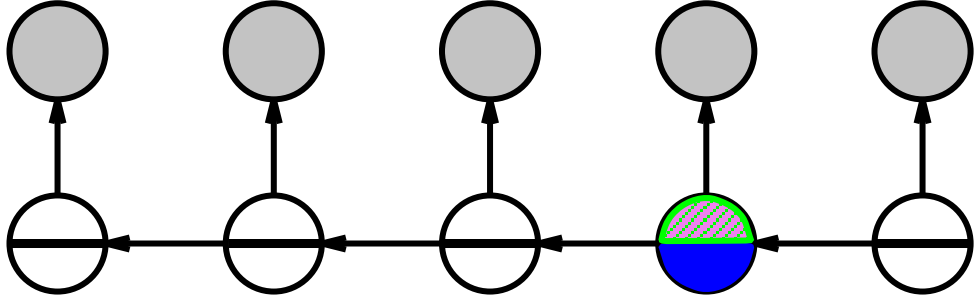
# Country dance algorithm

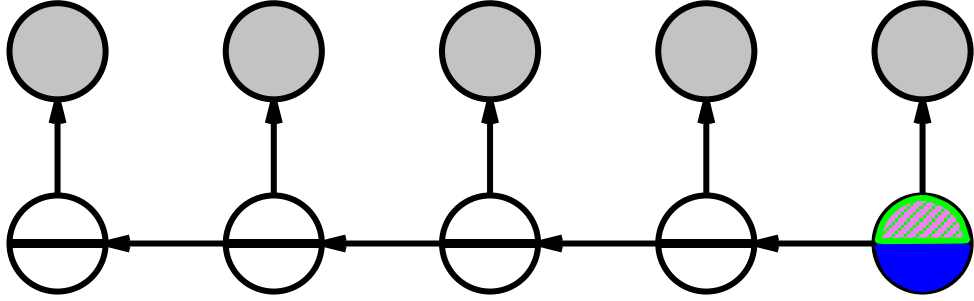
# Country dance algorithm

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\begin{aligned}
 f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^\top f_{1:t} \\
 \mathbf{O}_{t+1}^{-1} f_{1:t+1} &= \alpha_{t+1} \mathbf{T}_{1:t}^\top f_{1:t} \\
 \alpha'(\mathbf{T}^\top)^{-1} \mathbf{O}_{t+1}^{-1} f_{1:t+1} &= f_{1:t}
 \end{aligned}$$

Algorithm: forward pass computes  $f_t$ , backward pass does  $f_t, b_t$





Algorithm: forward pass computes  $f_t$ , backward pass does  $b_t, b_t$

$$f_{1:t+1} = \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t}$$

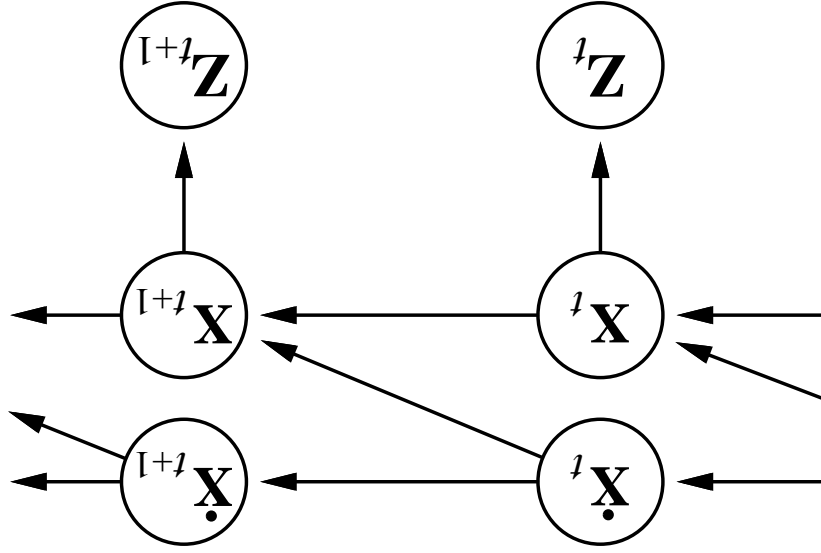
$$\mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = \alpha_{t+1} \mathbf{T}_{1:t}^T f_{1:t}$$

$$\alpha_{t+1} (\mathbf{T}_{1:t}^T)^{-1} \mathbf{O}_{1:t+1}^{-1} f_{1:t+1} = f_{1:t}$$

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

# Country dance algorithm

Gaussian prior, linear Gaussian transition model and sensor model



Airplanes, robots, ecosystems, economies, chemical plants, planets, ...

Modelling systems described by a set of continuous variables,  
 e.g., tracking a bird flying— $\mathbf{X}_t = X, \dot{X}, Y, \dot{Y}, Z, \dot{Z}$ .

# Kalman filters



## Updating Gaussian distributions

Prediction step: if  $P(\mathbf{X}_t | \mathbf{e}_{1:t})$  is Gaussian, then prediction

$$P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) = \int_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) d\mathbf{x}_t$$

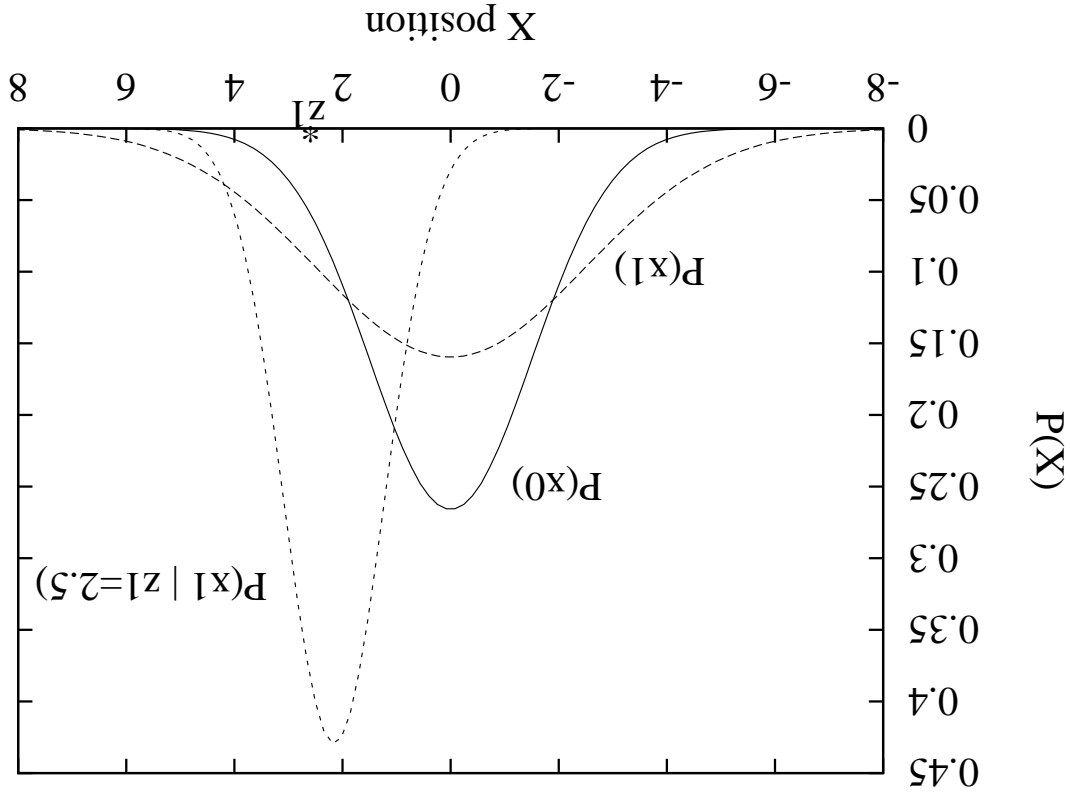
is Gaussian. If  $P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t})$  is Gaussian, then the updated distribution

$$P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t})$$

is Gaussian

Hence  $P(\mathbf{X}_t | \mathbf{e}_{1:t})$  is multivariate Gaussian  $N(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$  for all  $t$

General (nonlinear, non-Gaussian) process: description of posterior grows **unboundedly** as  $t \rightarrow \infty$



$$\mu_{t+1} = \frac{(\sigma_x^2 + \sigma_z^2)z_{t+1} + \sigma_z^2\mu_t}{\sigma_x^2 + \sigma_z^2 + \sigma_z^2}$$

$$\sigma_{t+1}^2 = \frac{(\sigma_x^2 + \sigma_z^2)\sigma_z^2}{\sigma_x^2 + \sigma_z^2 + \sigma_z^2}$$

Gaussian random walk on  $X$ -axis, s.d.  $\sigma_x$ , sensor s.d.  $\sigma_z$

# Simple 1-D example

## General Kalman update

Transition and sensor models:

$$P(\mathbf{x}_{t+1} | \mathbf{x}_t) = N(\mathbf{F}\mathbf{x}_t, \Sigma^x(\mathbf{x}_{t+1}))$$

$$P(z_t | \mathbf{x}_t) = N(\mathbf{H}\mathbf{x}_t, \Sigma^z(z_t))$$

$\mathbf{F}$  is the matrix for the transition;  $\Sigma^x$  the transition noise covariance  
 $\mathbf{H}$  is the matrix for the sensors;  $\Sigma^z$  the sensor noise covariance

Filter computes the following update:

$$\mu_{t+1} = \mathbf{F}\mu_t + \mathbf{K}_{t+1}(z_{t+1} - \mathbf{H}\mathbf{F}\mu_t)$$

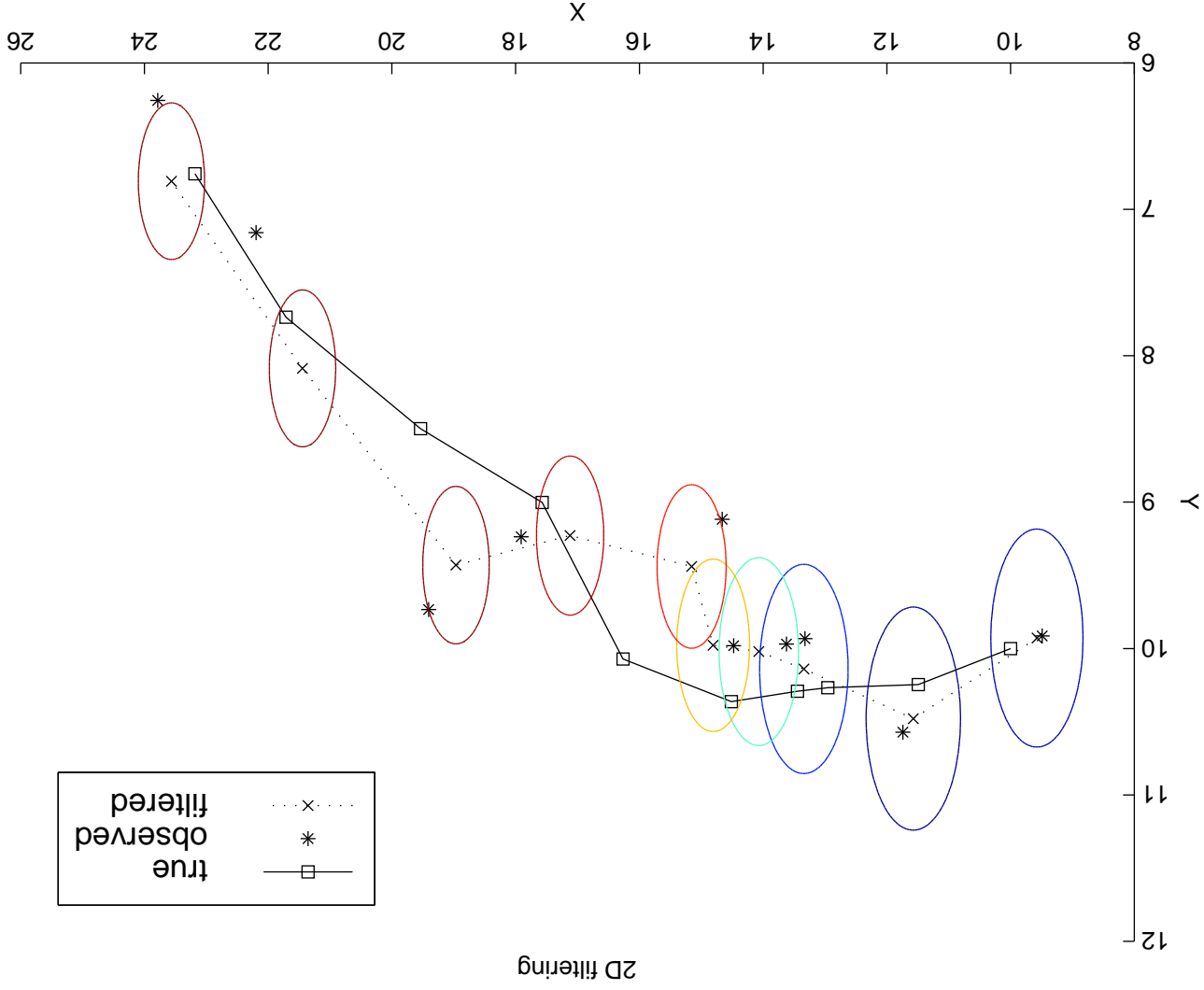
$$\Sigma_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1})(\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma^x)$$

where  $\mathbf{K}_{t+1} = (\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma^x)\mathbf{H}^\top(\mathbf{H}(\mathbf{F}\Sigma_t\mathbf{F}^\top + \Sigma^x)\mathbf{H}^\top + \Sigma^z)^{-1}$

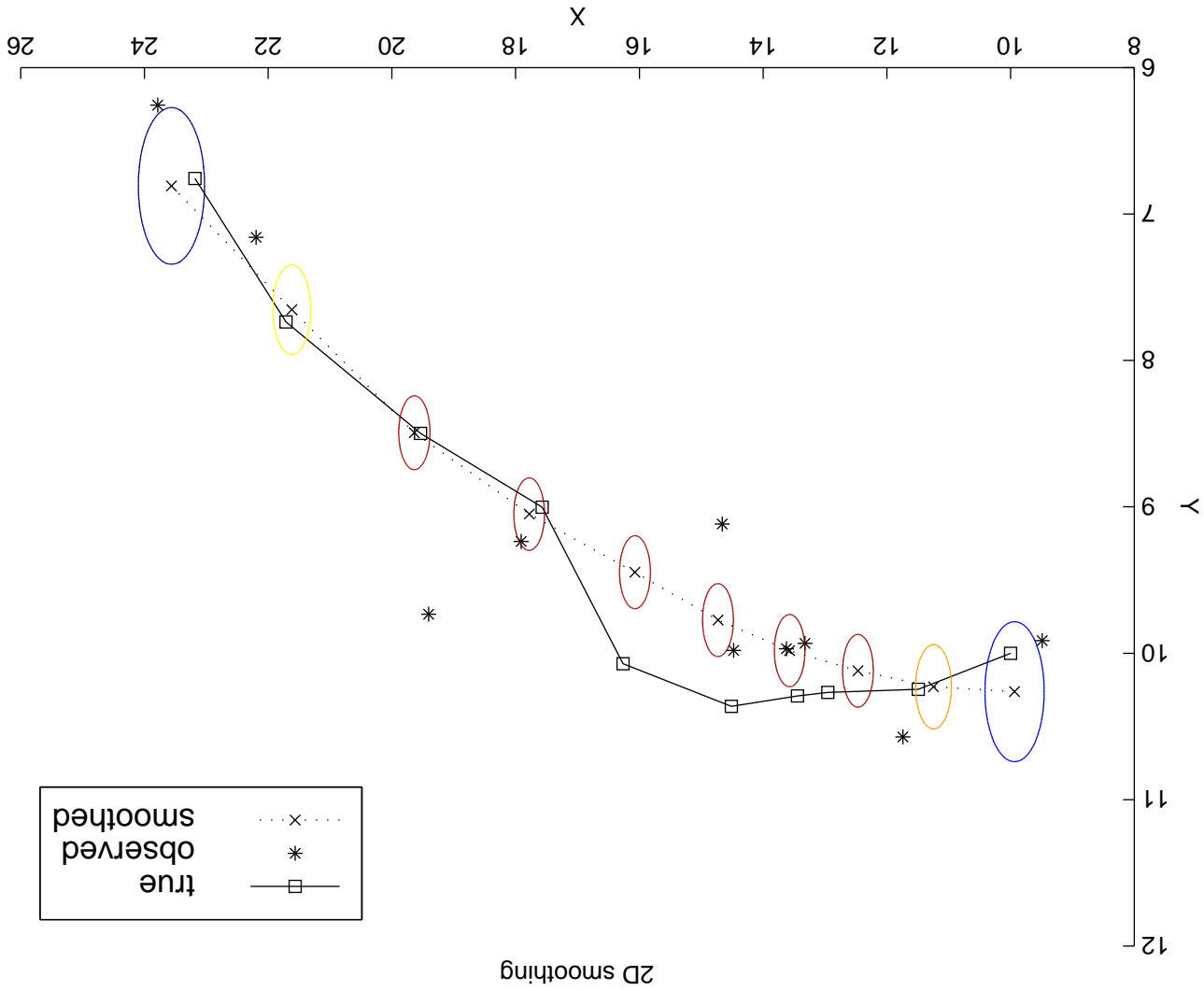
is the Kalman gain matrix

$\Sigma_t$  and  $\mathbf{K}_t$  are independent of observation sequence, so compute offline

# 2-D tracking example: filtering

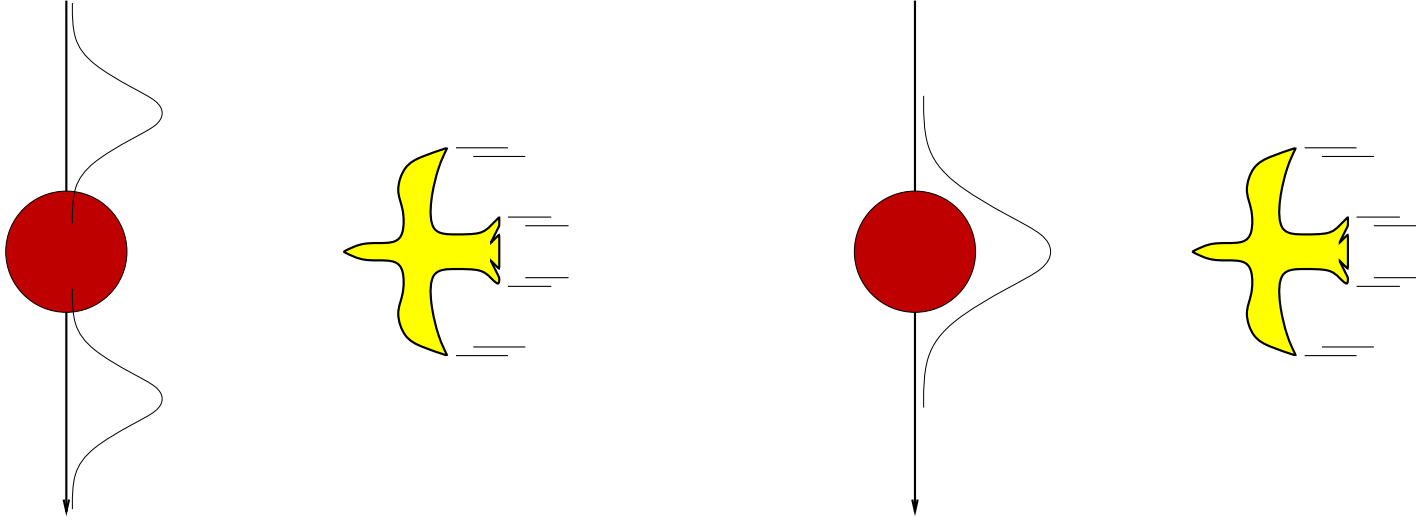


# 2-D tracking example: smoothing



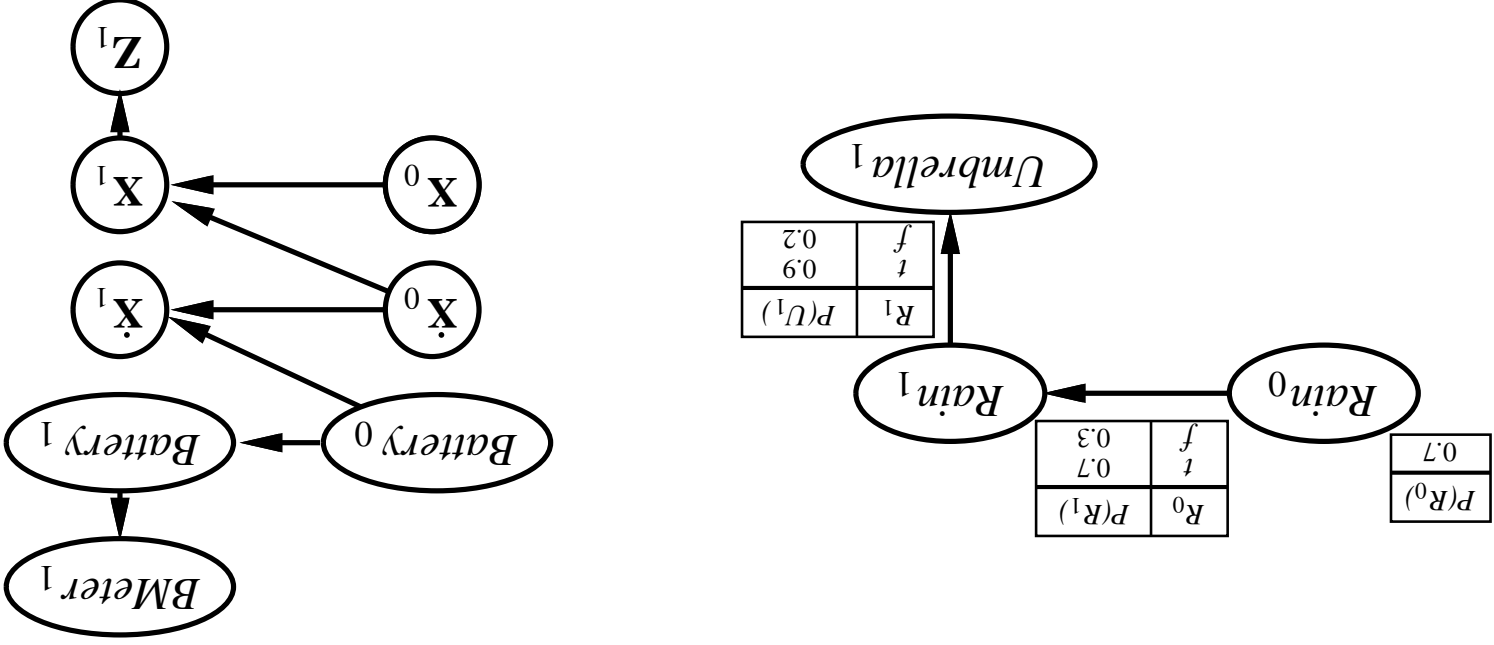
Cannot be applied if the transition model is nonlinear  
Extended Kalman Filter models transition as **locally linear** around  $\mathbf{x}_t = \mu_t$   
Fails if systems is locally unsmooth

**Where it breaks**



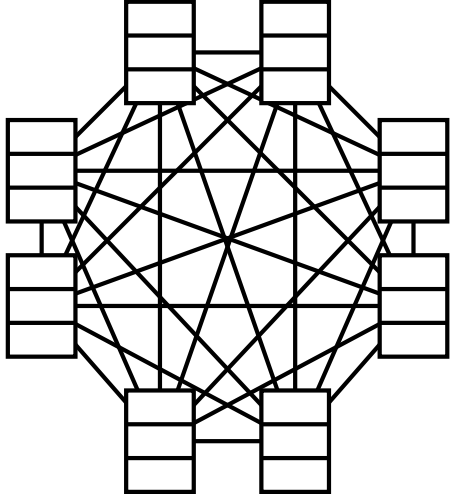
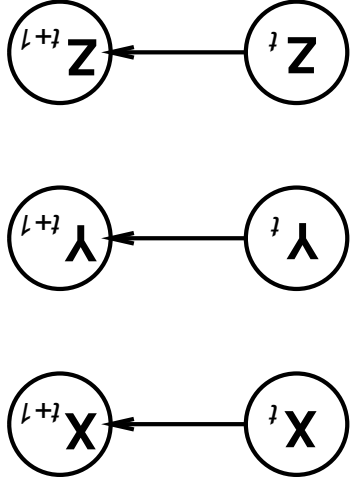
# Dynamic Bayesian networks

$X_t, E_t$  contain arbitrarily many variables in a replicated Bayes net



# DBNs vs. HMMs

Every HMM is a single-variable DBN; every discrete DBN is an HMM



Sparse dependencies  $\Rightarrow$  exponentially fewer parameters;

e.g., 20 state variables, three parents each

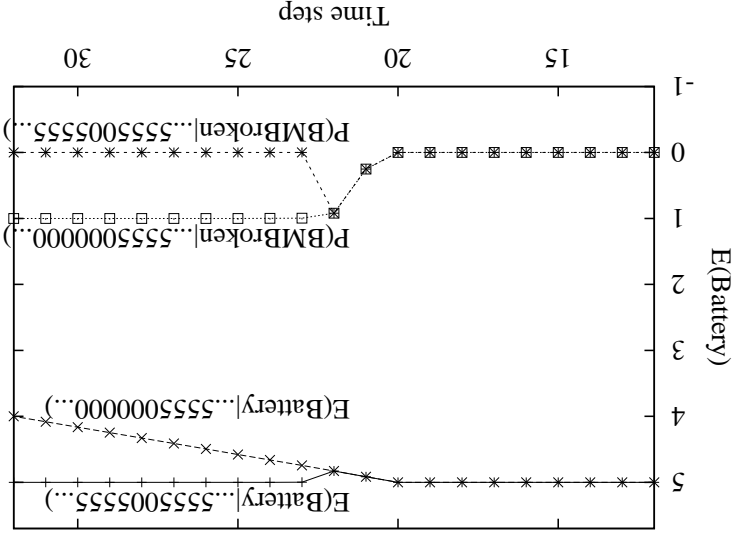
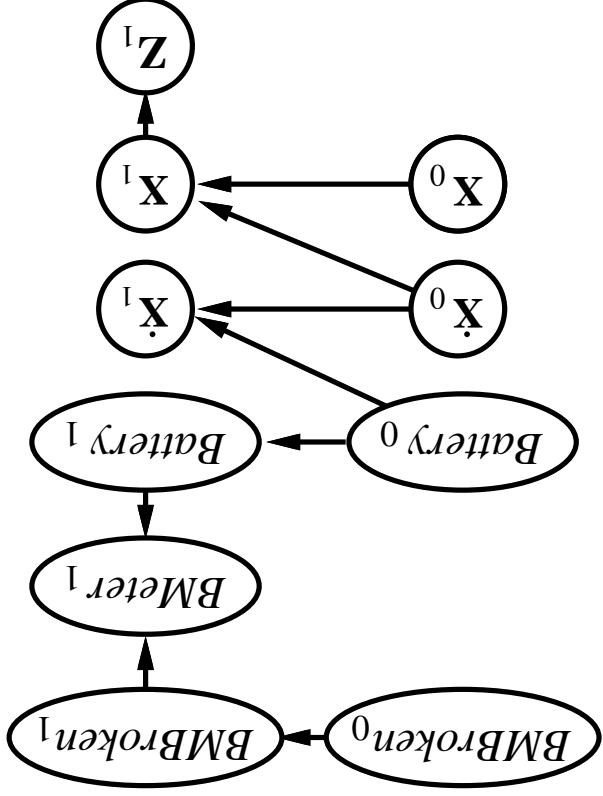
DBN has  $20 \times 2^3 = 160$  parameters, HMM has  $2^{20} \times 2^{20} \approx 10^{12}$



# DBNs vs Kalman filters

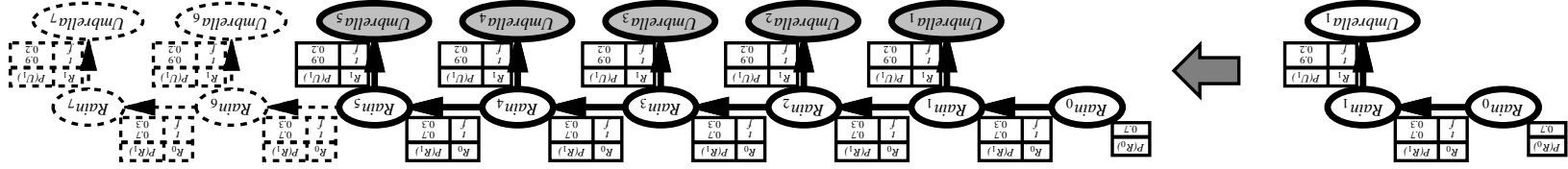
Every Kalman filter model is a DBN, but few DBNs are KFs; real world requires non-Gaussian posteriors

E.g., where are bin Laden and my keys? What's the battery charge?



# Exact inference in DBNs

Naive method: **unroll** the network and run any exact algorithm



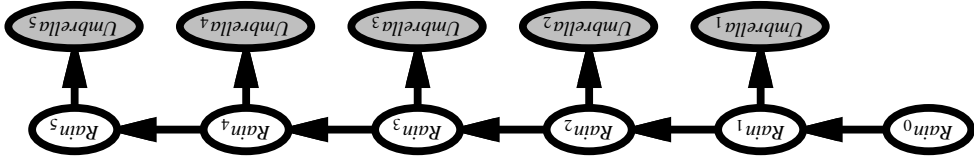
Problem: inference cost for each update grows with  $t$

**Rollup filtering**: add slice  $t + 1$ , "sum out" slice  $t$  using variable elimination

Largest factor is  $O(p^{n+1})$ , update cost  $O(p^{n+2})$  (cf. HMM update cost  $O(p^{2n})$ )

# Likelihood weighting for DBNs

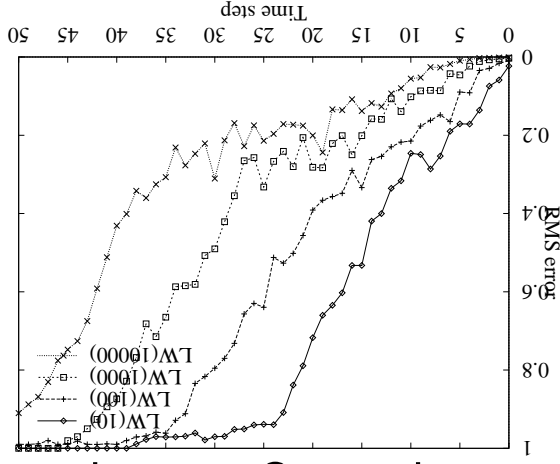
Set of weighted samples approximates the belief state



LW samples pay no attention to the evidence!

⇒ fraction “agreeing” falls exponentially with  $t$

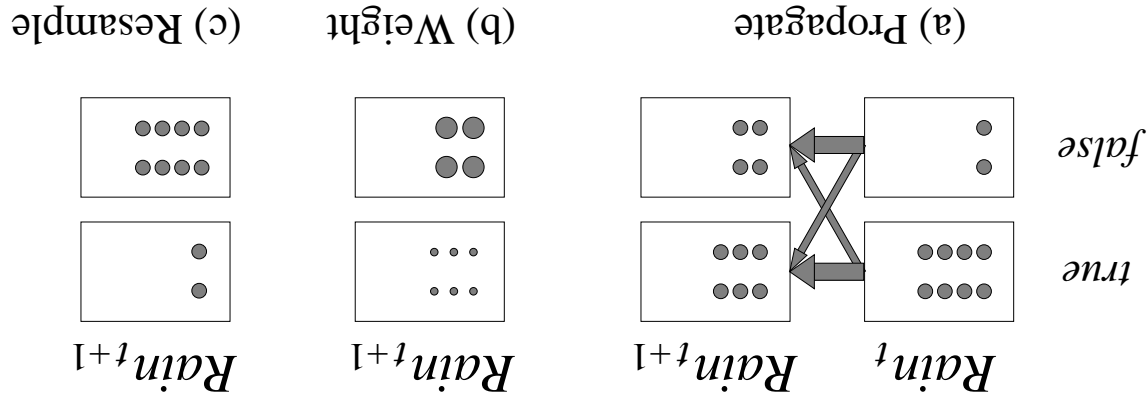
⇒ number of samples required grows exponentially with  $t$



# Particle filtering

Basic idea: ensure that the population of samples ("particles") tracks the high-likelihood regions of the state-space

Replicate particles proportional to likelihood for  $e_t$



Widely used for tracking nonlinear systems, esp. in vision

Also used for simultaneous localization and mapping in mobile robots  
 $10^5$ -dimensional state space

## Particle filtering contd.

Assume consistent at time  $t$ :  $N(\mathbf{x}_t | \mathbf{e}_{1:t}) / N = P(\mathbf{x}_t | \mathbf{e}_{1:t})$

Propagate forward: populations of  $\mathbf{x}_{t+1}$  are

$$N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t}) = \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) N(\mathbf{x}_t | \mathbf{e}_{1:t})$$

Weight samples by their likelihood for  $\mathbf{e}_{t+1}$ :

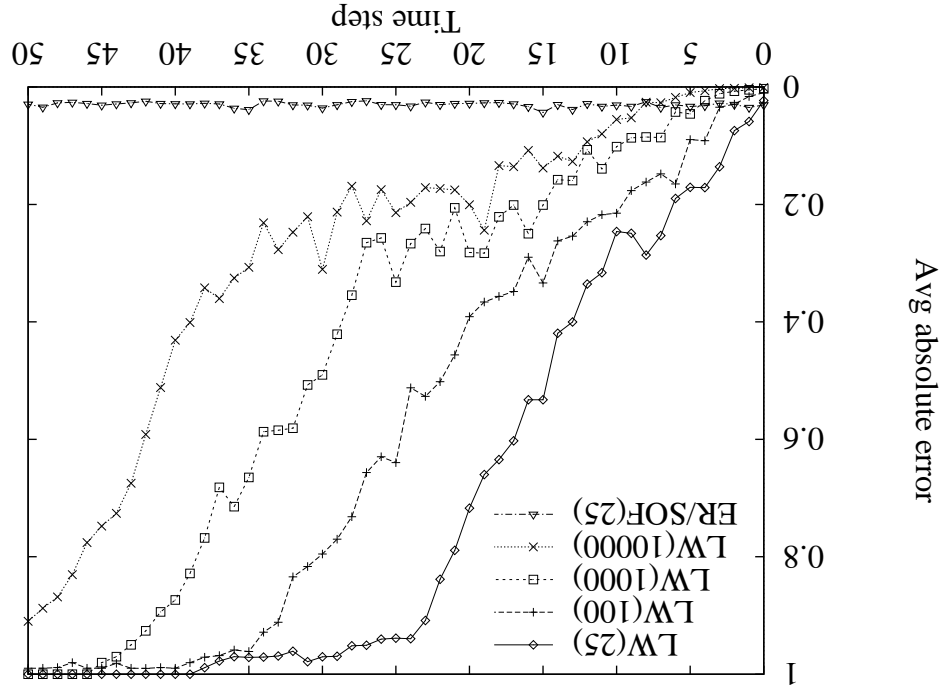
$$W(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) = P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t})$$

Resample to obtain populations proportional to  $W$ :

$$\begin{aligned} N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) / N &= \alpha W(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) = \alpha P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) N(\mathbf{x}_{t+1} | \mathbf{e}_{1:t}) \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) N(\mathbf{x}_t | \mathbf{e}_{1:t}) \\ &= \alpha' P(\mathbf{e}_{t+1} | \mathbf{x}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) \\ &= P(\mathbf{x}_{t+1} | \mathbf{e}_{1:t+1}) \end{aligned}$$

# Particle filtering performance

Approximation error of particle filtering remains bounded over time, at least empirically—theoretical analysis is difficult



## Summary

Temporal models use state and sensor variables replicated over time

Markov assumptions and stationarity assumption, so we need

- transition model  $P(X_t | X_{t-1})$
- sensor model  $P(E_t | X_t)$

Tasks are filtering, prediction, smoothing, most likely sequence;

**all done recursively with constant cost per time step**

Hidden Markov models have a single discrete state variable; used

for speech recognition

Kalman filters allow  $n$  state variables, linear Gaussian,  $O(n^3)$  update

Dynamic Bayes nets subsume HMMs, Kalman filters; exact update intractable

Particle filtering is a good approximate filtering algorithm for DBNs