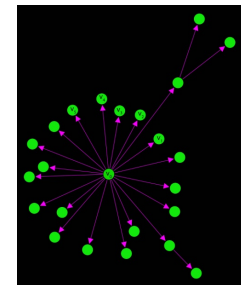
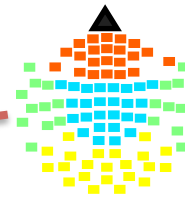
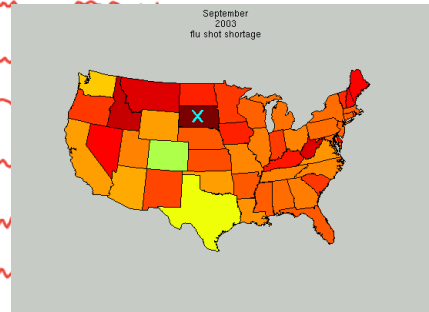
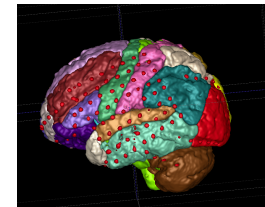
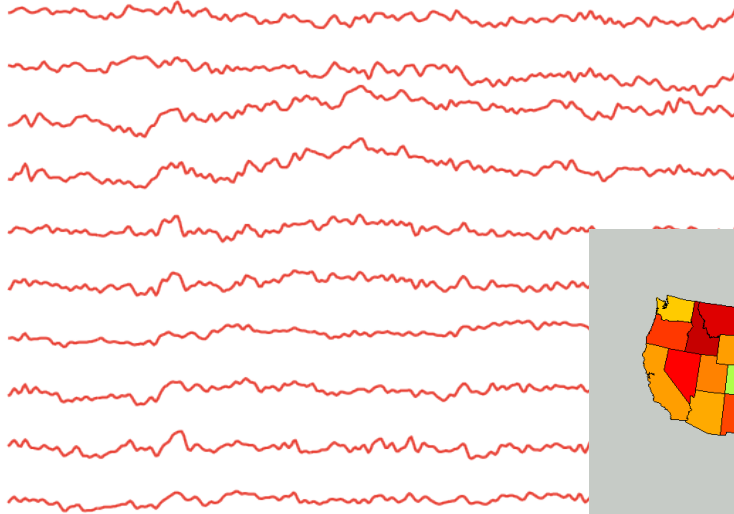


# **Bayesian Time Series: Structured Representations for Scalability**

Emily Fox

# Time Series



**HUMAN GENOME, IN NUMBERS**



**6 billion** DNA LETTERS

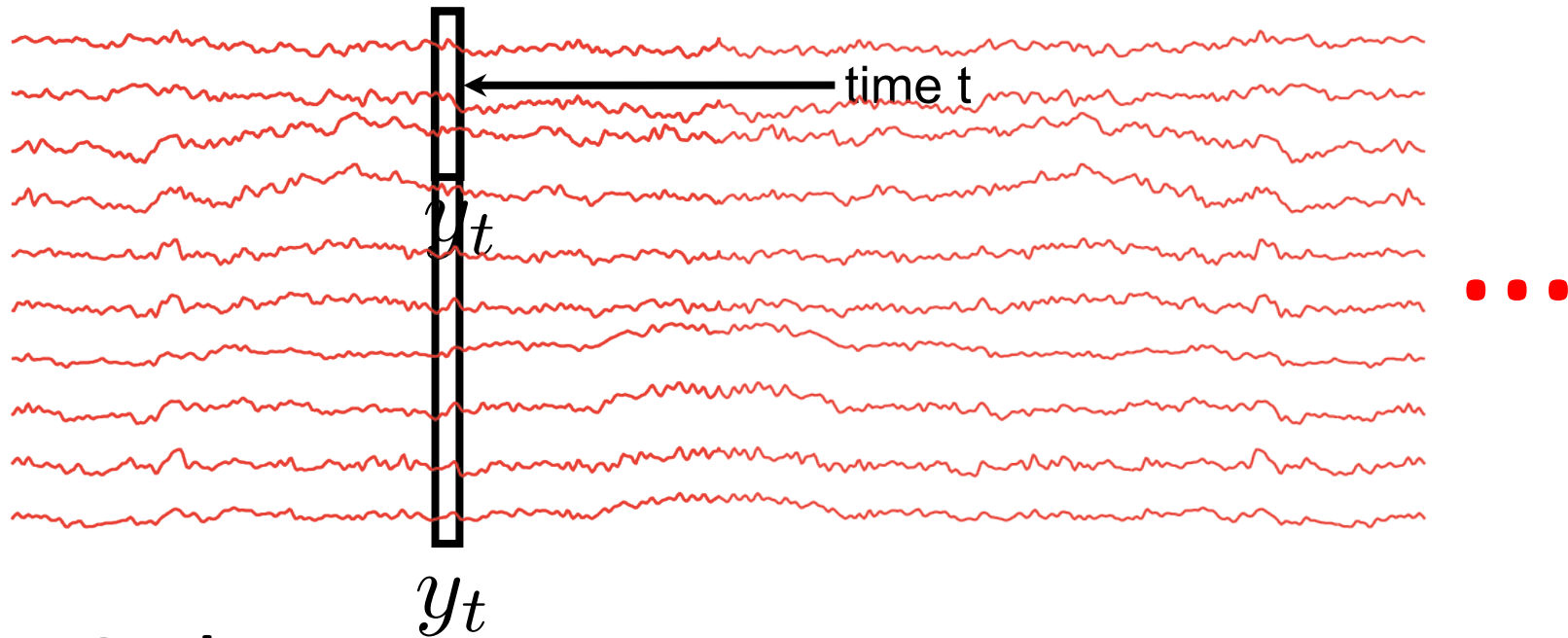
**22,000** GENES

46 CHROMOSOMES

**\$9,500** COST TO SEQUENCE

500GB SIZE ON DISK

Sources: NIH, Illumina



## Goals:

- Evolution – Dynamics across time
- Relational structure – Dependencies between series

## Modeling challenges:

- Large  $p$  – Many dimensions/series
- Irregular grid of observations
- Missing values
- Heterogeneous data sources
- ...

## Computational challenges:

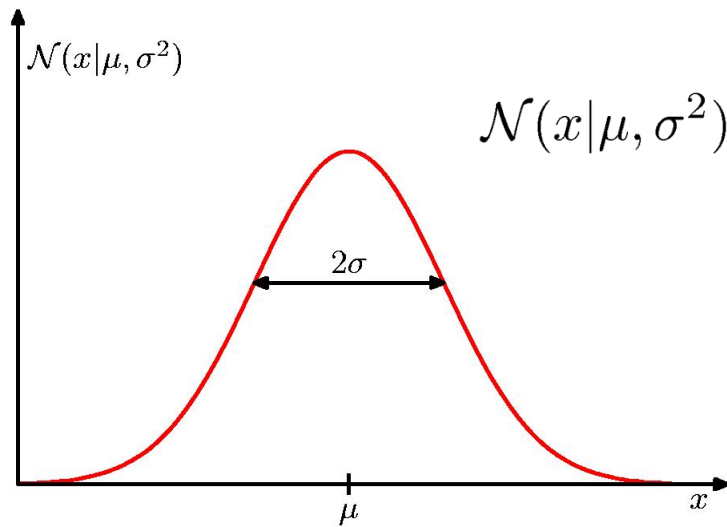
- Large  $n$  – Long time series
- Streaming data –  
Continuum of observations

# Preliminaries/Review

- Multivariate Gaussians
- Hidden Markov models (HMMs)
- Vector autoregressive (VAR) processes
  - Stability/stationarity
- Gaussian state space models
  - Identifiability

# Quick Review of Gaussians

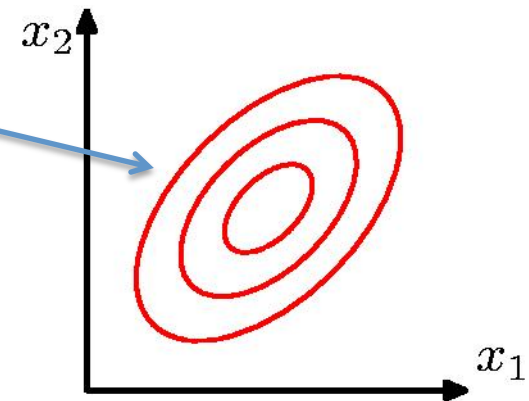
- Univariate and multivariate Gaussians



$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

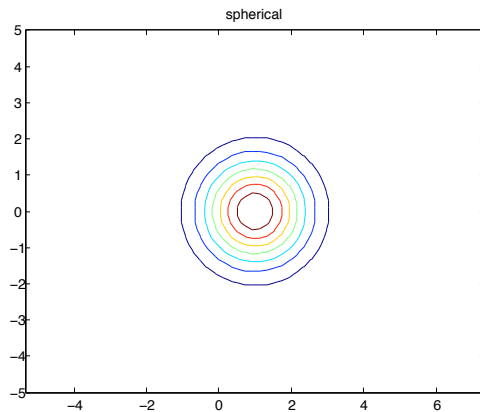
Covariance  
defines shape  
(eigvecs+eigvals)

$$\Sigma = \text{COV}(\mathbf{x})$$

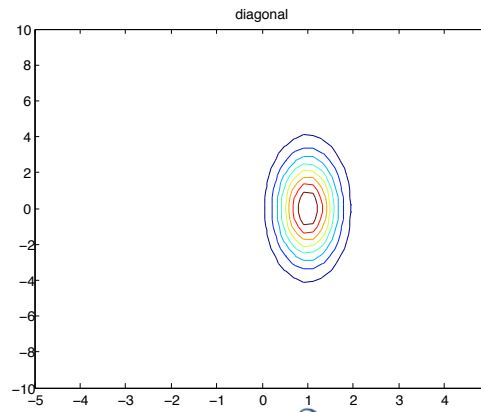


$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

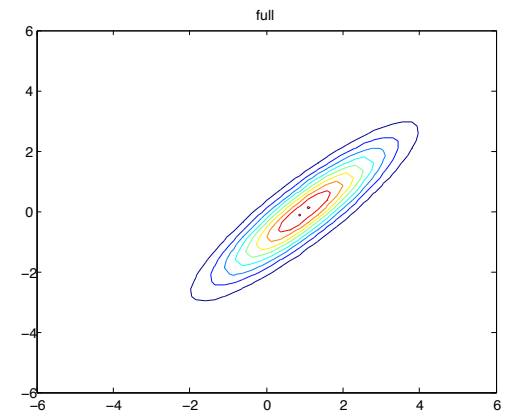
# Two-Dimensional Gaussians



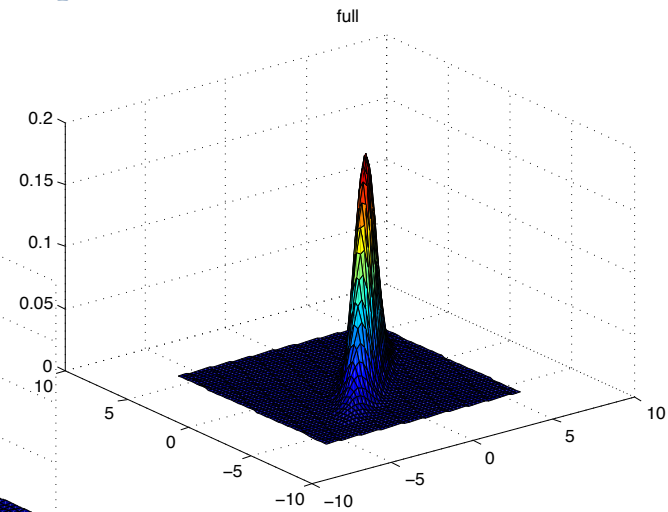
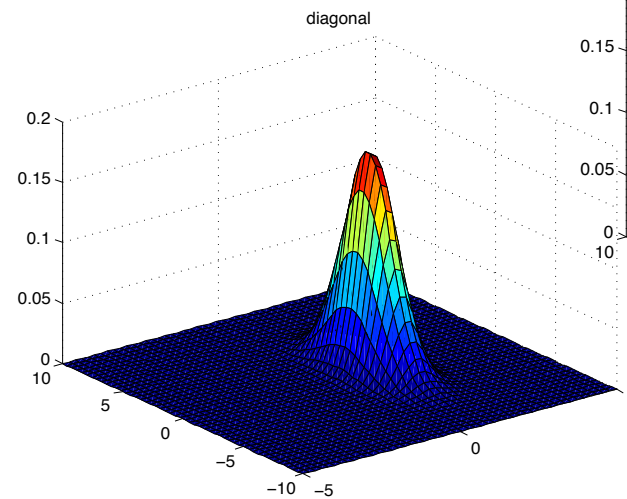
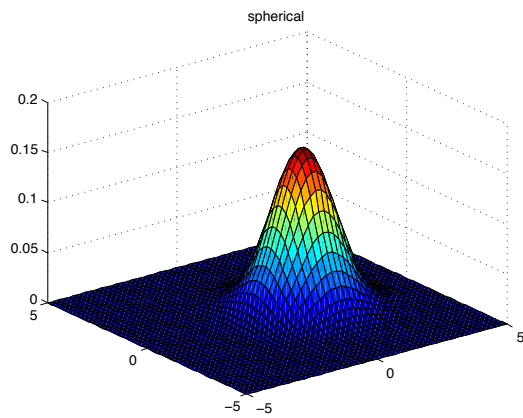
$$\Sigma = \sigma^2 I$$



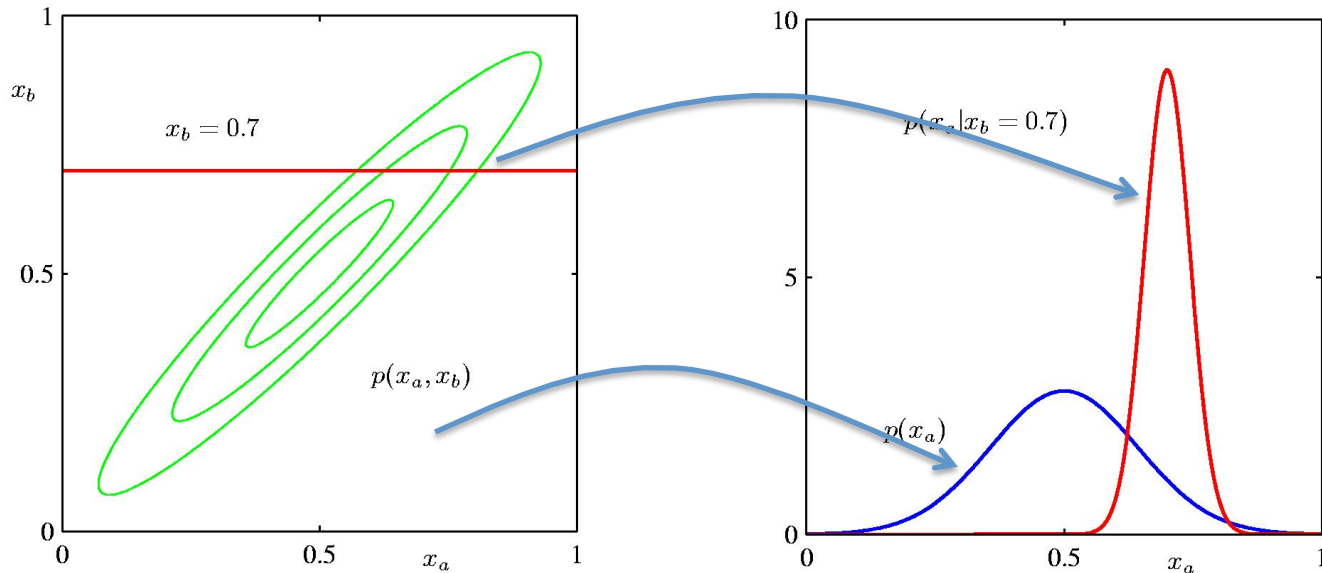
$$\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$$



$$\Sigma \text{ general}$$



# Conditional & Marginal Distributions



$$\begin{pmatrix} x_a \\ x_b \end{pmatrix} \sim N \left( \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ab}^T & \Sigma_{bb} \end{bmatrix} \right)$$

Marginally:  $x_a \sim N(\mu_a, \Sigma_{aa})$

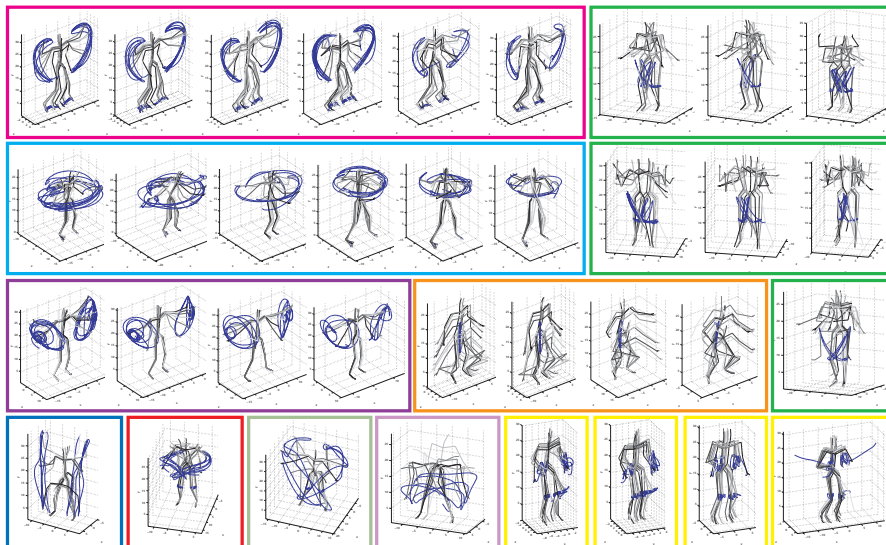
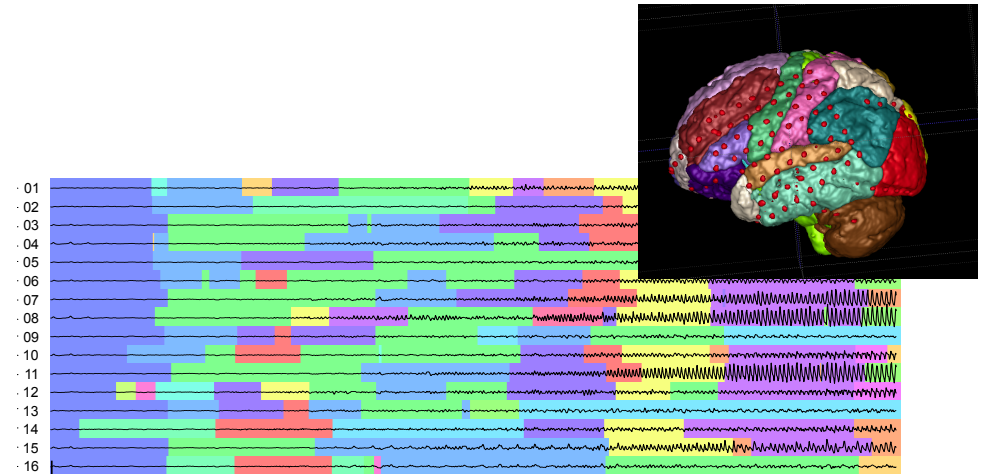
Conditionally:

$$x_a | x_b \sim N(\mu_a - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b), \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ab}^T)$$

# Hidden Markov Models

## Example applications:

- Parsing EEG recordings
- Discovering behaviors in videos
- Speech segmentation
- Volatility regimes in financial time series
- Genomics
- ...





# Example: Motion Capture Segmentation

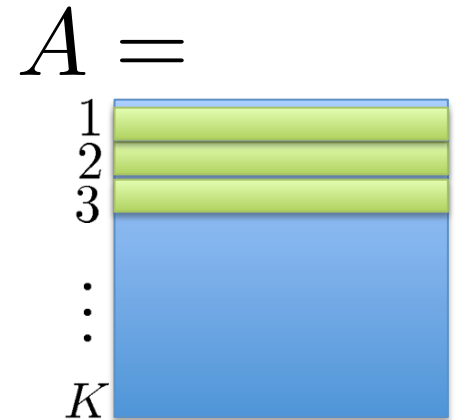
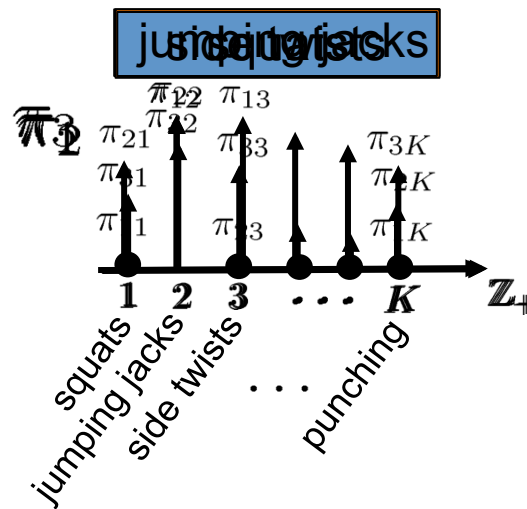
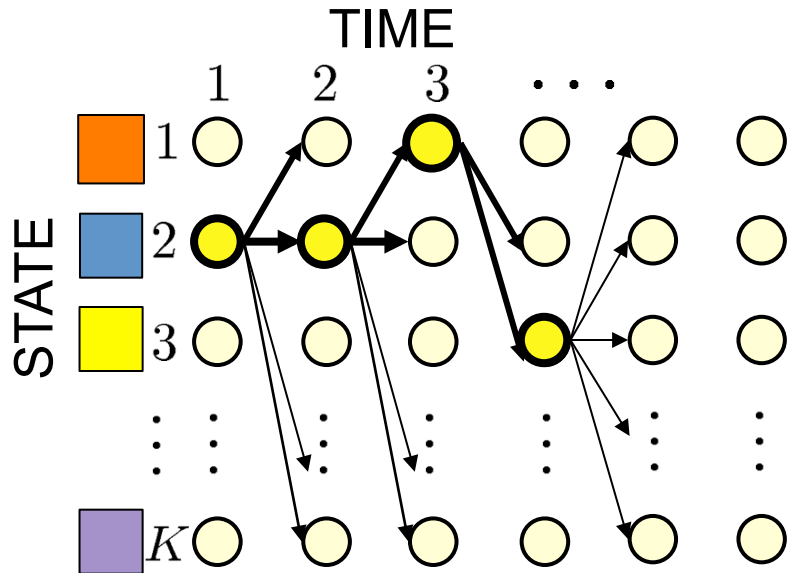
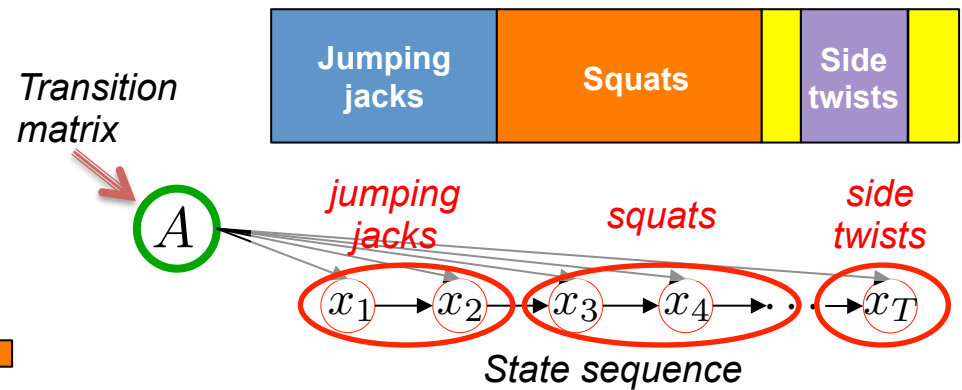


Jump- ing jacks	Side twists	Run	Squats			
-----------------------	----------------	-----	--------	--	--	--

# Hidden Markov Model

- Markov transition dynamics:

$$\Pr(x_t = \text{orange} \mid x_{t-1} = \text{blue}) = A_{\text{blue orange}}$$



Tutorial:  
Rabiner, *Proc. IEEE* 1989

# Hidden Markov Model

- Markov transition dynamics:

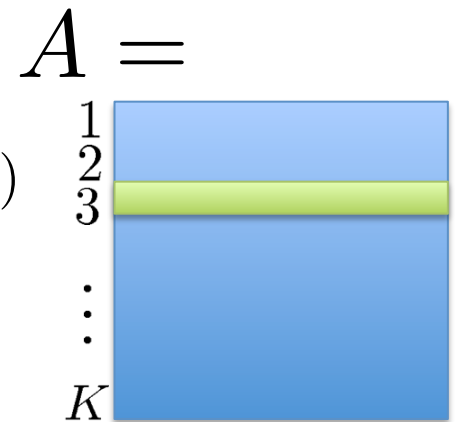
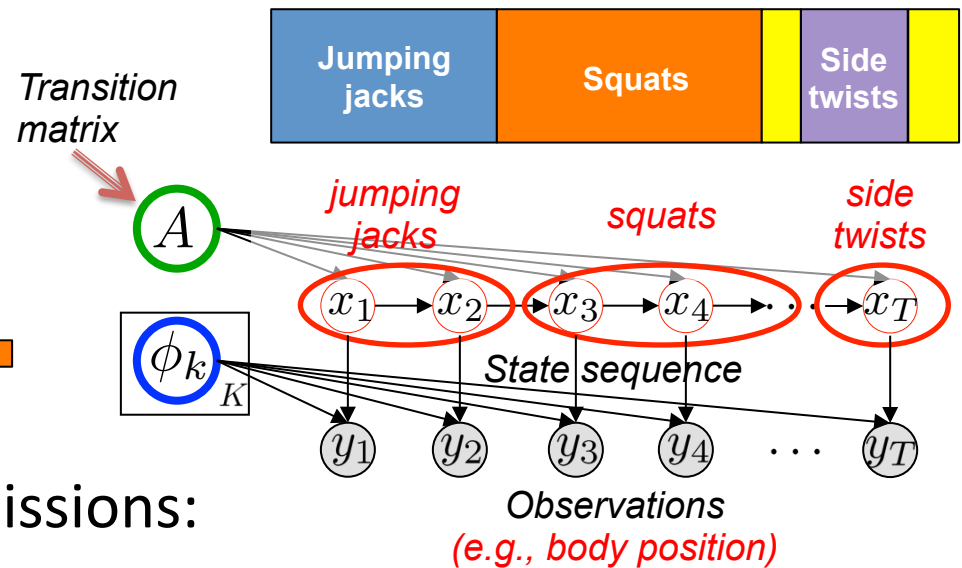
$$\Pr(x_t = \text{orange} \mid x_{t-1} = \text{blue}) = A_{\text{blue, orange}}$$

- Conditionally independent emissions:

$$y_t \mid x_t = \text{orange} \sim F(\phi_{\text{orange}})$$

- Latent Markov chain structure enables

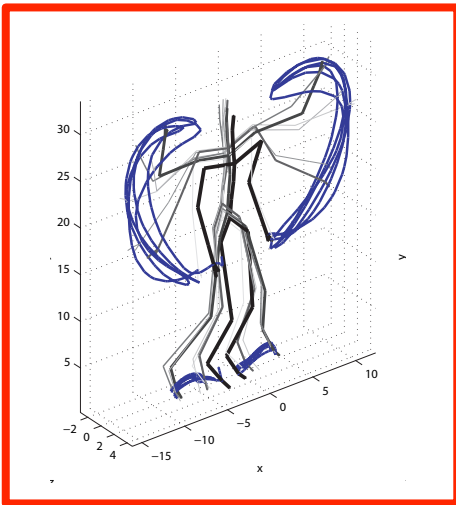
- Efficient computation of marginals  $p(x_t \mid y_1, \dots, y_T)$  using the *forward-backward* algorithm
- Most-probable sequence via *Viterbi*
- Parameter learning using *Baum-Welch* (EM for HMMs)



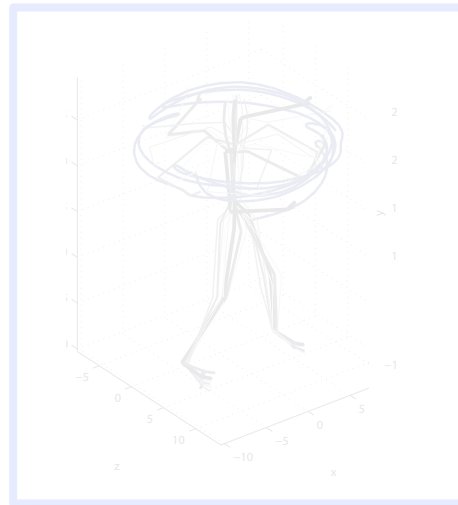
Tutorial:  
Rabiner, *Proc. IEEE* 1989

# Motivating Other Time Series Models

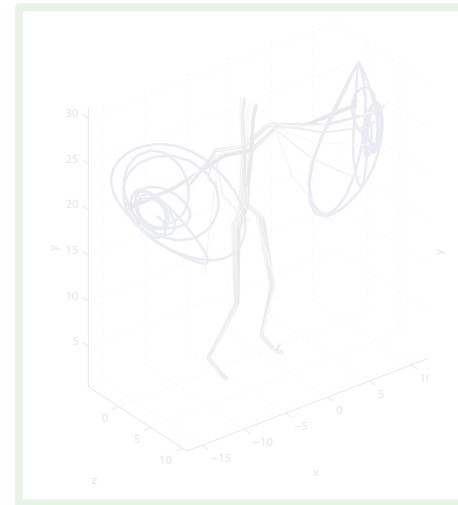
Jumping Jacks



Side Twists



Arm Circles



...

**Vector autoregressive (VAR) process:**

$$y_t = \sum_{i=1}^r A_i y_{t-i} + e_t \quad e_t \sim N(0, \Sigma)$$

# Stationary VAR Processes

$$y_t = \sum_{i=1}^r A_i y_{t-i} + e_t \quad e_t \sim N(0, \Sigma)$$

- If the *companion matrix* has eigenvalues  $\lambda$  with  $|\lambda| < 1$ , then the process is **stable**

$$\begin{bmatrix} A_1 & A_2 & \cdots & A_r \\ I & 0 & \cdots & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & \cdots & I & 0 \end{bmatrix}$$

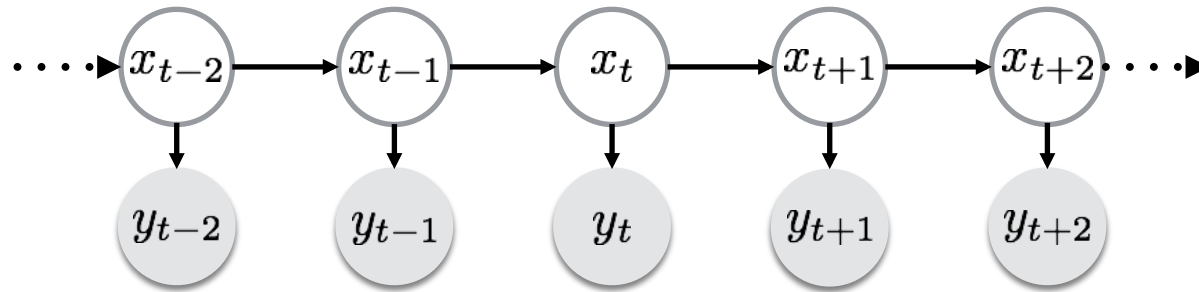
- If initialized infinitely in the past, then **stationary**

$$E[y_t] = \mu = 0 \quad \text{cov}(y_t, y_{t+h}) = \Gamma(h)$$

- For VAR(1) process, marginal covariance satisfies

$$\Gamma(0) = A_1 \Gamma(0) A_1' + \Sigma$$

# State Space Models



- Like HMMs, but continuous-valued latent state sequence

$$x_t = Ax_{t-1} + e_t \quad e_t \sim N(0, \Sigma)$$

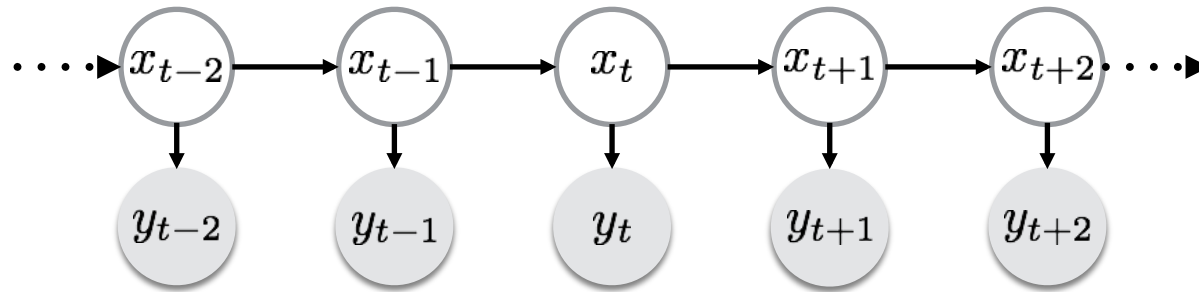
$$y_t = Cx_t + w_t \quad w_t \sim N(0, R)$$

- Entire class of equivalent systems from input/output perspective by changing latent space via  $x_t \rightarrow Tx_t$ :

$$\begin{aligned}
 x_t &= \overset{\leftarrow T^{-1}AT}{\tilde{A}}x_{t-1} + e_t & e_t &\sim N(0, \Sigma) \\
 y_t &= \overset{CT}{\tilde{C}}x_t + w_t & w_t &\sim N(0, R)
 \end{aligned}$$

Constrain  
A,  $\Sigma$ , or C

# State Space Models

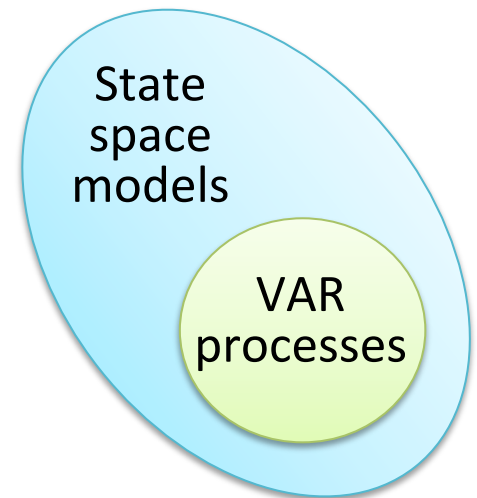


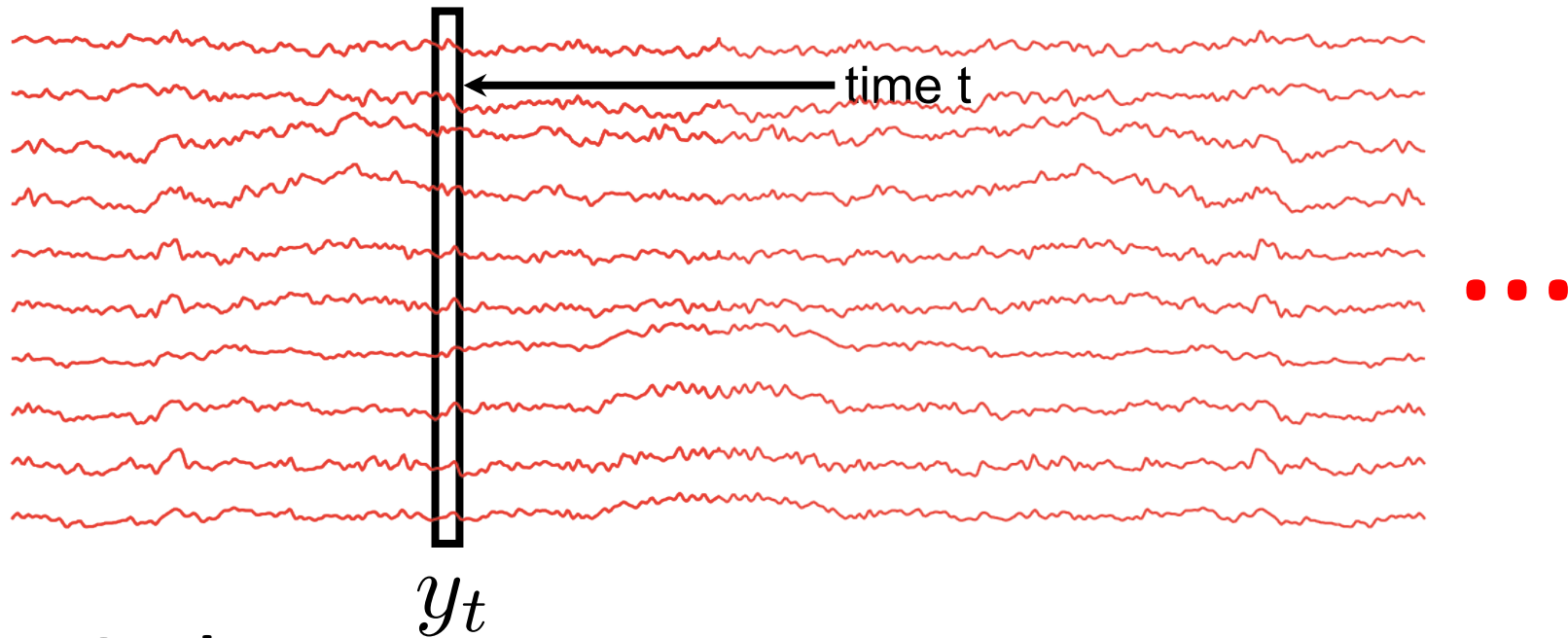
- Can write a VAR( $r$ ) process in state space form via

$$y_t = \sum_{i=1}^r A_i y_{t-i} + e_t \quad e_t \sim N(0, \Sigma)$$

$$x_t = \begin{bmatrix} A_1 & A_2 & \cdots & A_r \\ I & 0 & \cdots & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & \cdots & I & 0 \end{bmatrix} x_{t-1} + \begin{bmatrix} I \\ 0 \\ \vdots \\ 0 \end{bmatrix} e_t$$

$$y_t = [I \quad 0 \cdots 0] x_t$$





## Goals:

- Evolution – Dynamics across time
- Relational structure – Dependencies between series

## Modeling challenges:

- Large  $p$  – Many dimensions/series
- Irregular grid of observations
- Missing values
- Heterogeneous data sources
- ...

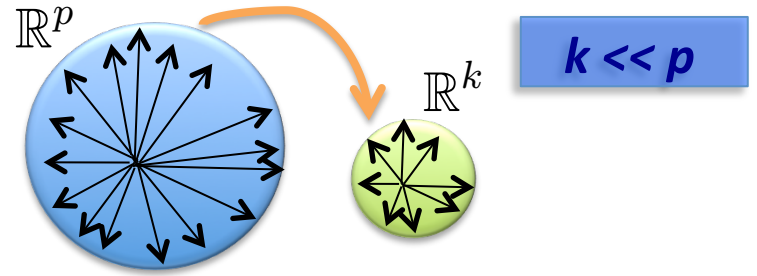
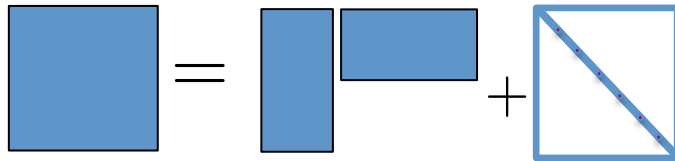
## Computational challenges:

- Large  $n$  – Long time series
- Streaming data –  
Continuum of observations



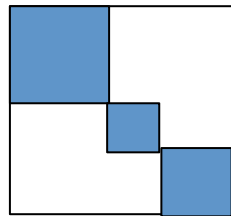
# Methods for Scaling to High Dimensions

Low Rank  $\Sigma = \Lambda\Lambda' + \Sigma_0$



*Low-dimensional embedding*

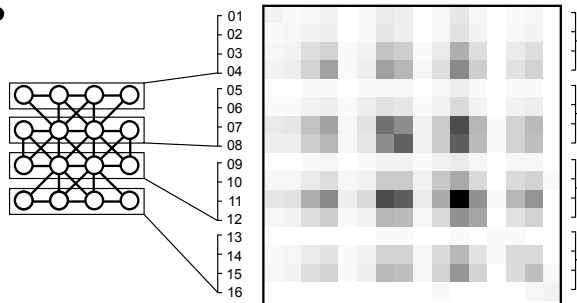
$\Sigma$  sparse



*Independent groups of nodes*

$\Sigma^{-1}$  sparse

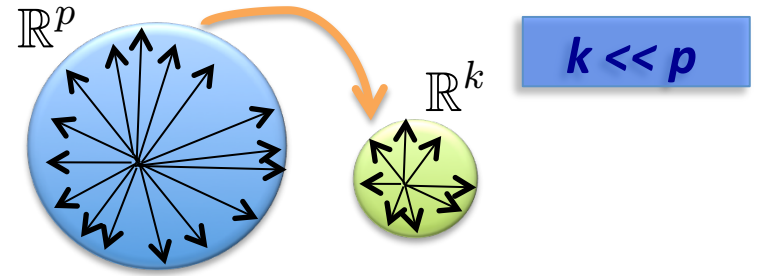
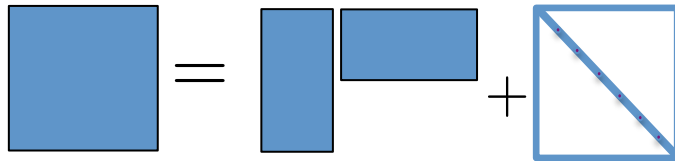
Gaussian Graphical Model



*Zeros = no edge in graph,  
Cond. ind. between nodes*

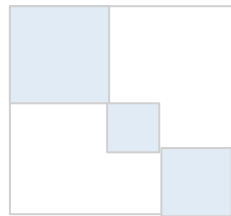
# Methods for Scaling to High Dimensions

Low Rank  $\Sigma = \Lambda\Lambda' + \Sigma_0$



Low-dimensional embedding

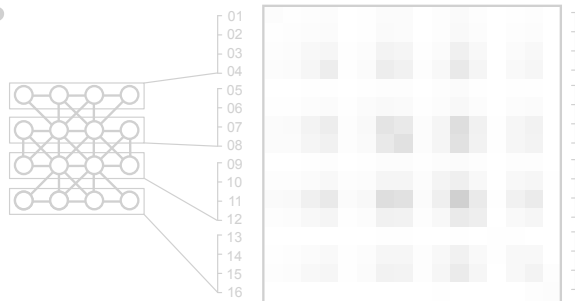
$\Sigma$  sparse



Independent groups of nodes

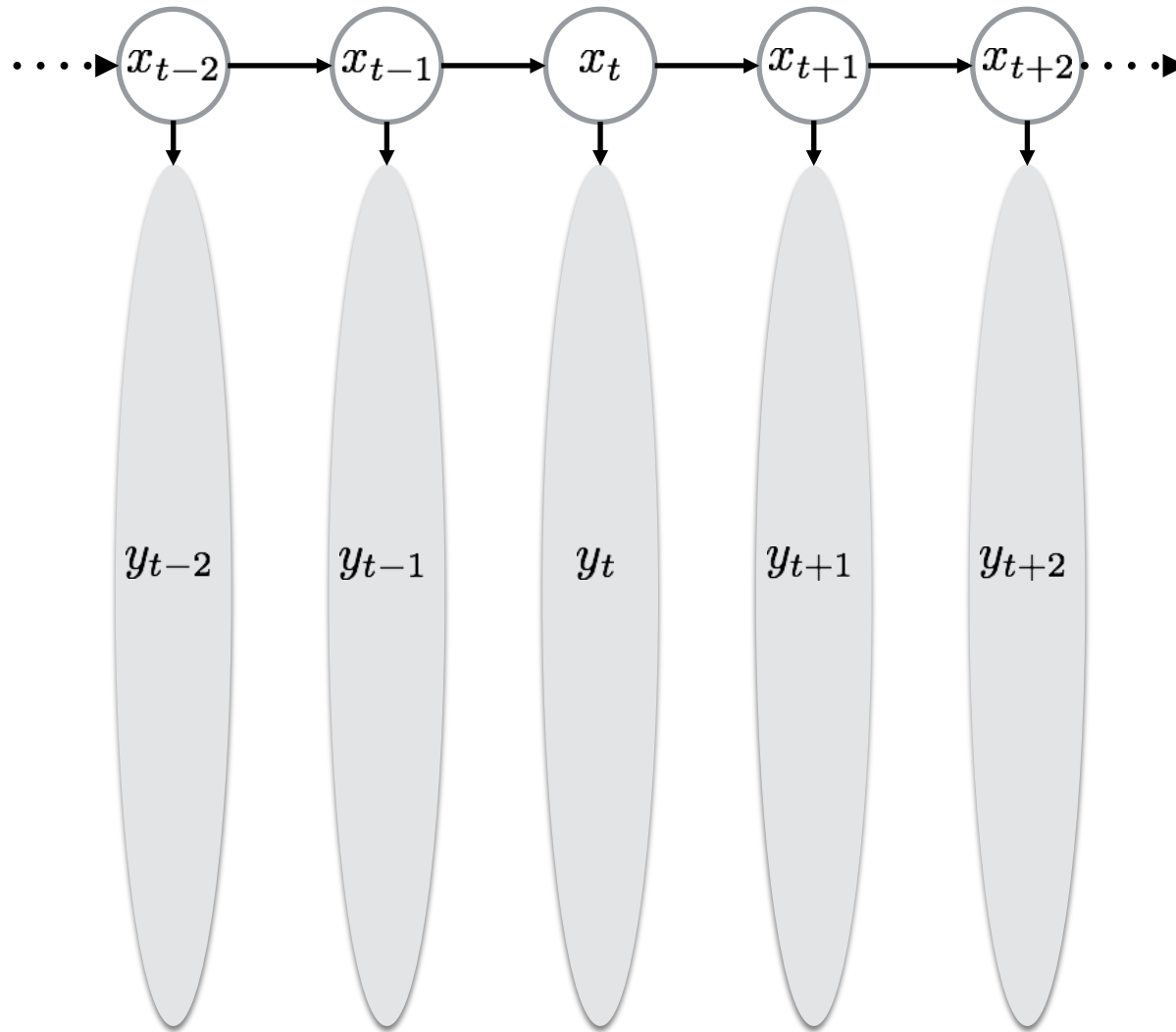
$\Sigma^{-1}$  sparse

Gaussian Graphical Model

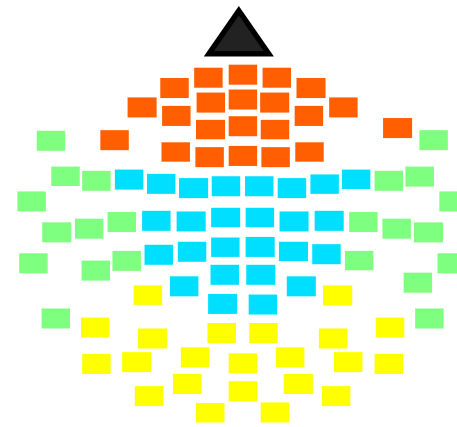


Zeros = no edge in graph,  
Cond. ind. between nodes

# Modeling High-Dimensional Time Series



# Magnetoencephalography (MEG)



Helmet with  
102 sensors

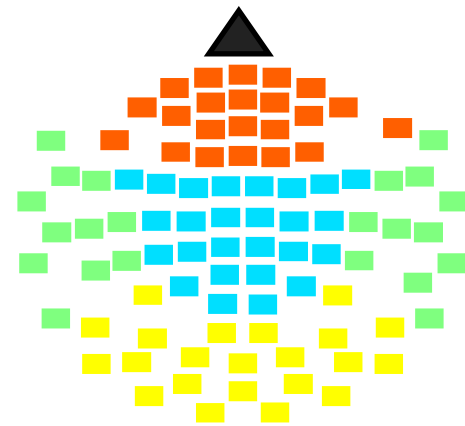
- How does the brain code concepts?
  - e.g. animals, food...

⋮

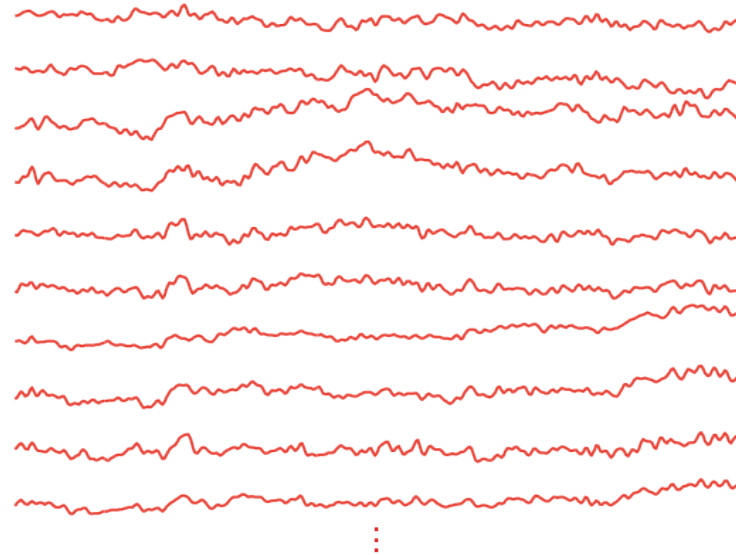
# Magnetoencephalography (MEG)



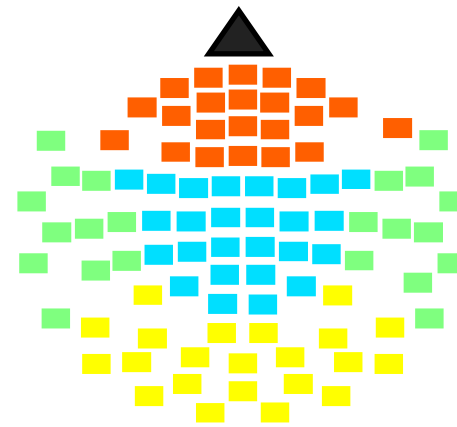
COW



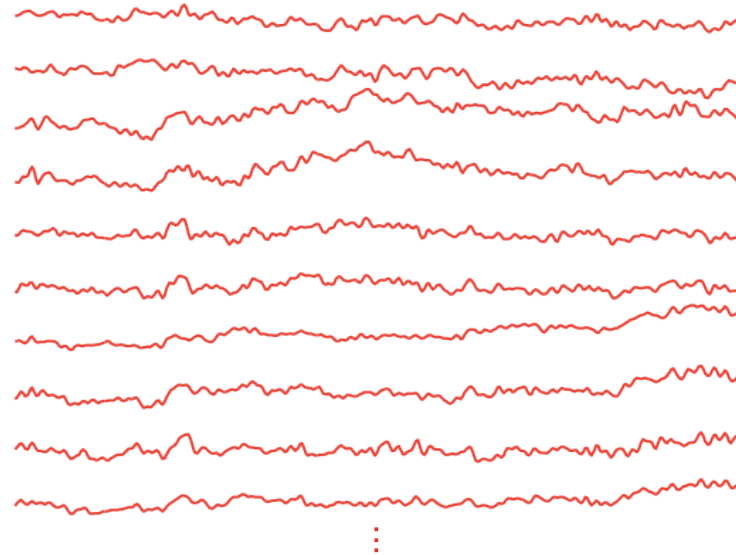
Helmet with  
102 sensors



# Magnetoencephalography (MEG)



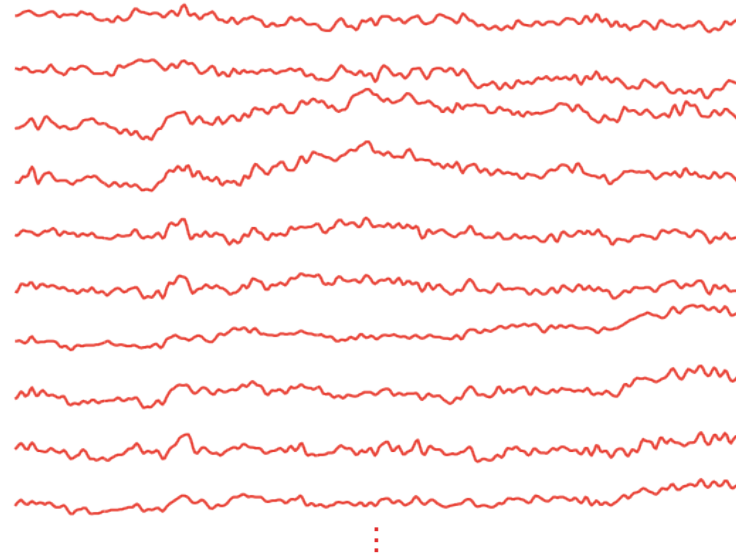
Helmet with  
102 sensors



# Magnetoencephalography (MEG)

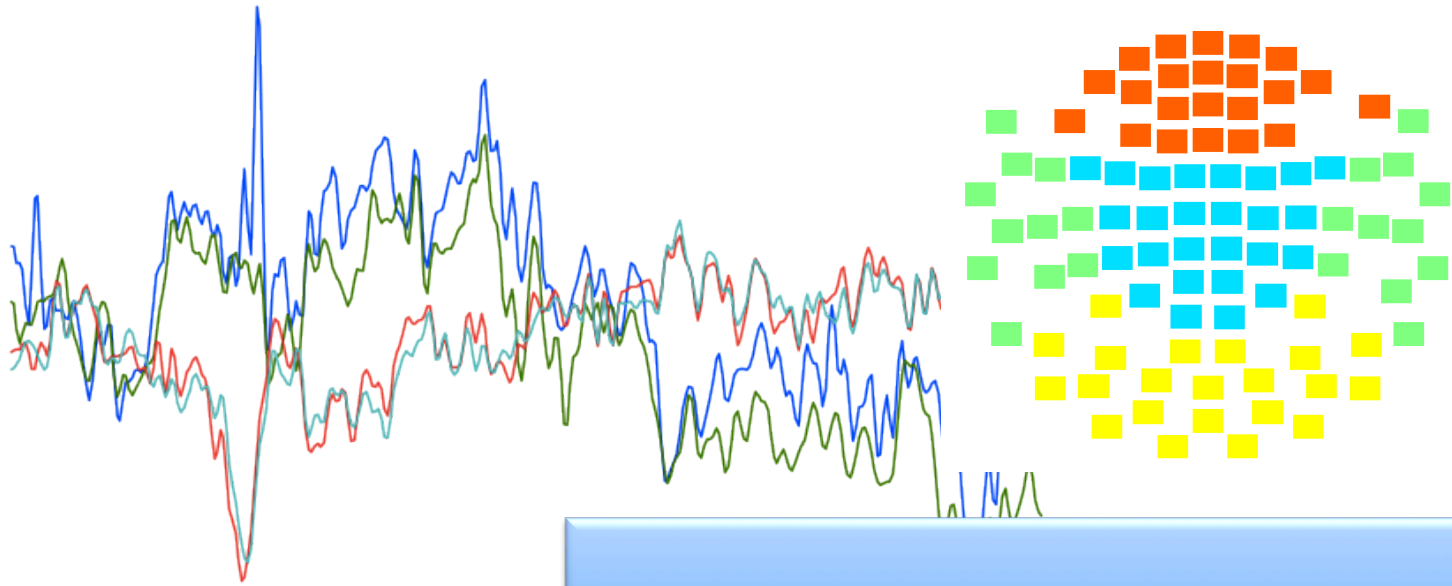


- High dimensional
- Time-varying correlations  
→ *Functional connectivity*



# Coping with Dimensionality

- Observation: Sensors are *redundant*

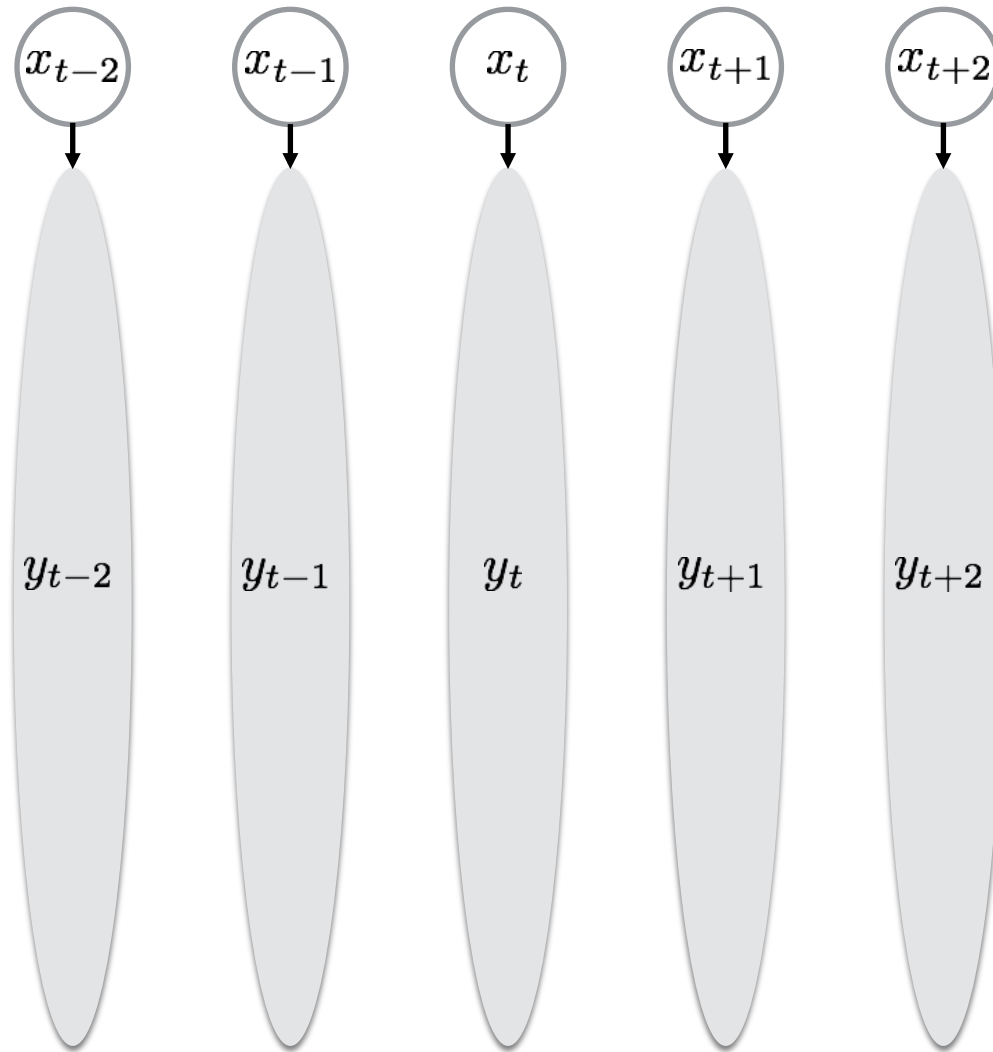


Dynamic Latent Factor Models

- Goal:
  - Harness *low-dimensional embedding* of dynamics



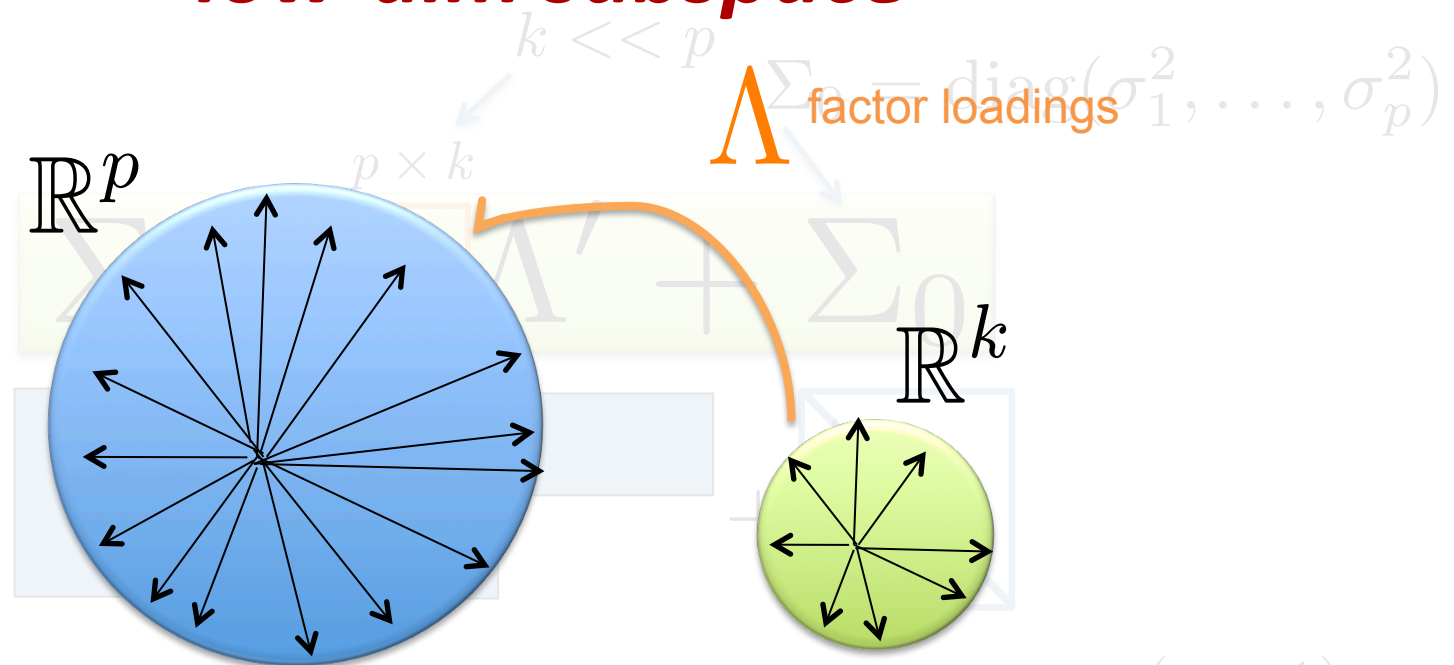
# High-Dim i.i.d. Data



# Latent Factor Model

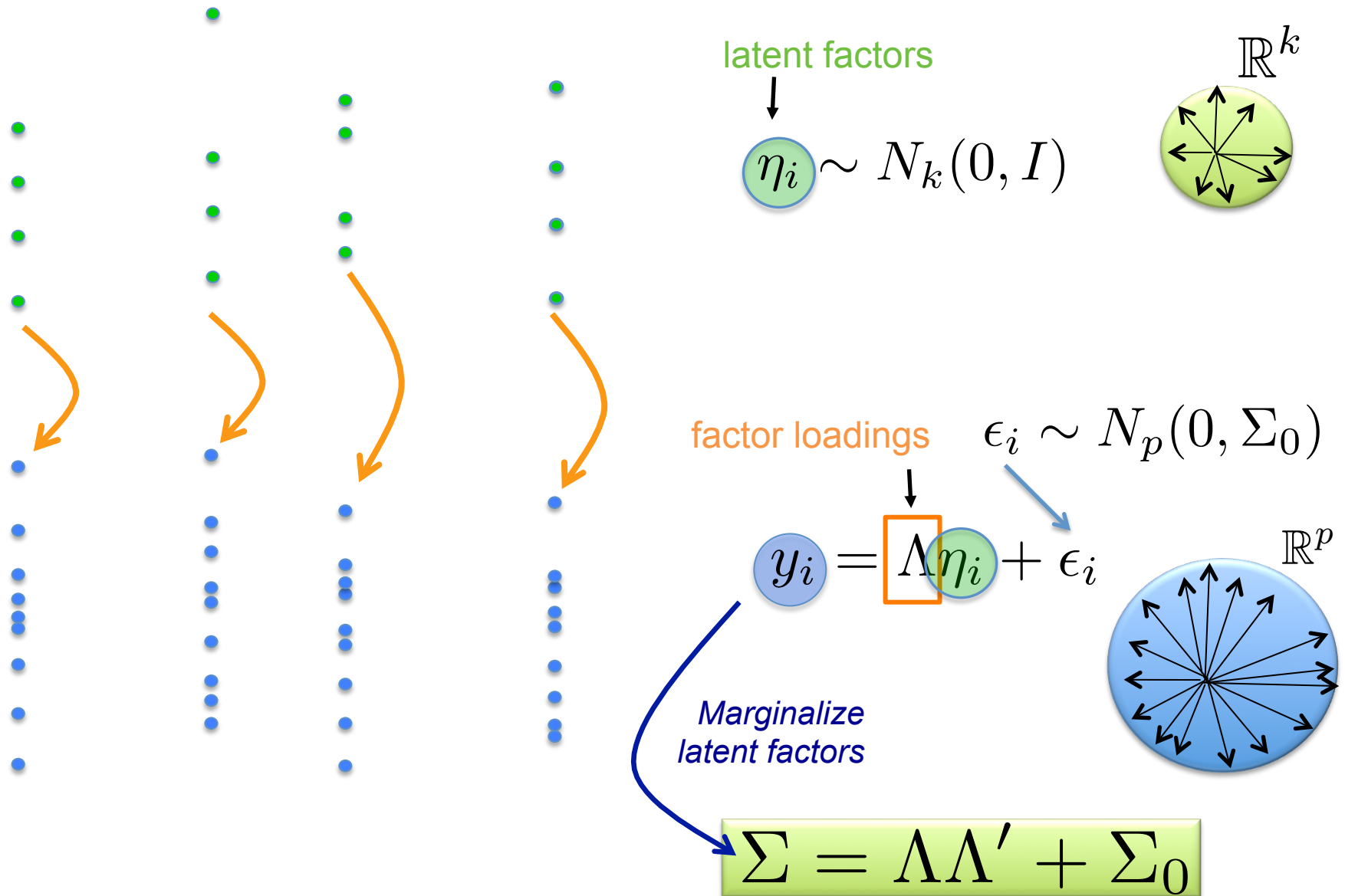
- Assume normally distributed data

## Modeling statistical uncertainty in *low-dim subspace*



- Number of parameters:  $pk + p = p(k + 1) \ll \frac{p(p + 1)}{2}$

# Latent Factor Model

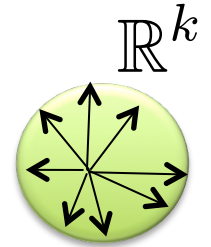


# Derivation of Marginal Distribution

- Marginal mean:

$$\begin{aligned} E[y_i] &= E[\Lambda\eta_i + \epsilon_i] \\ &= \Lambda E[\eta_i] + E[\epsilon_i] = 0 \end{aligned}$$

$$\eta_i \sim N_k(0, I)$$

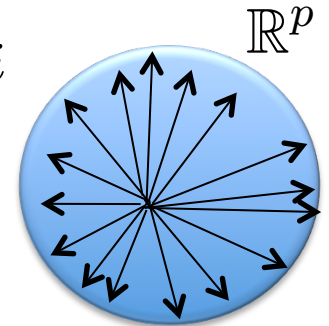


- Marginal covariance:

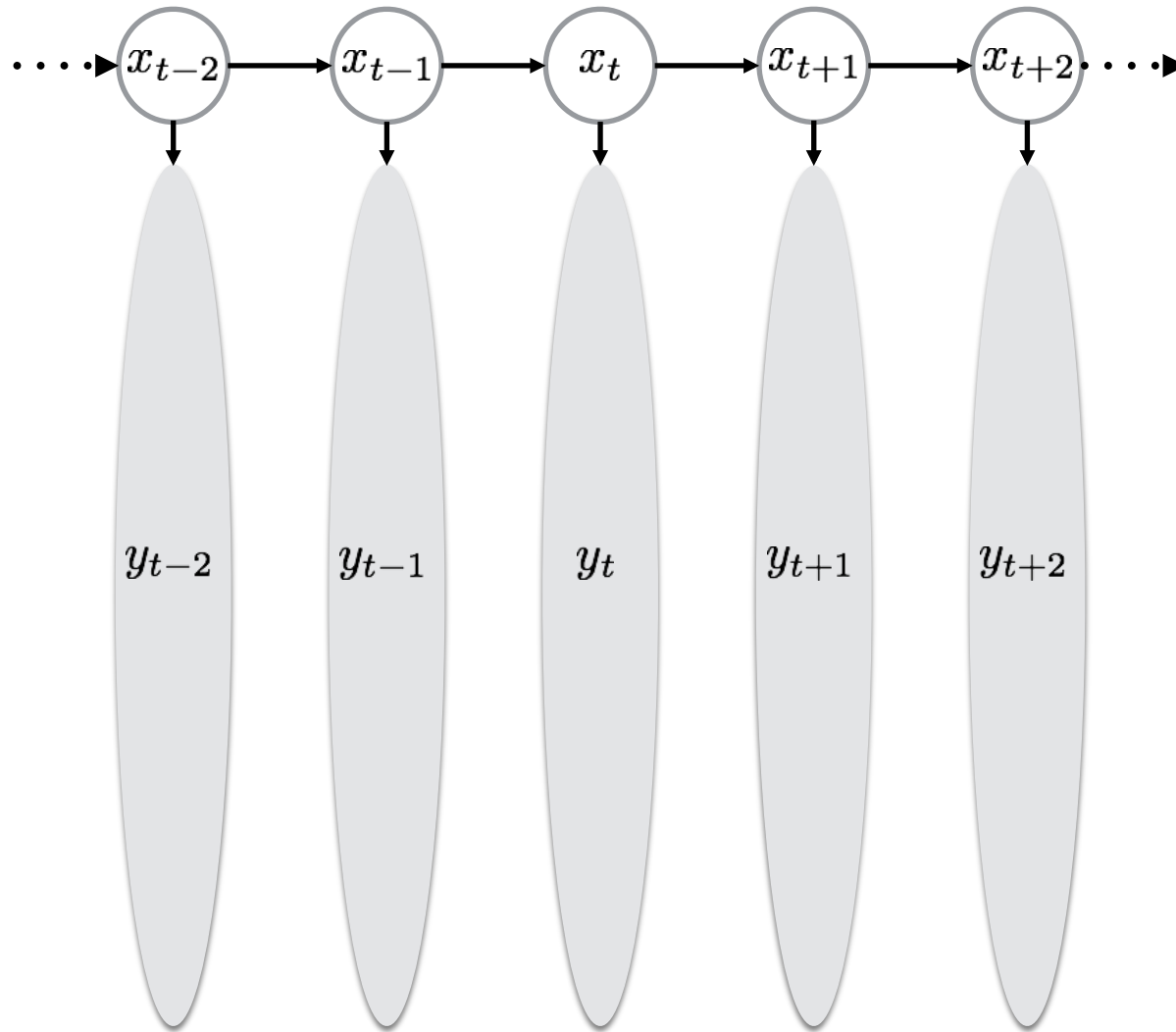
$$\begin{aligned} \text{cov}(y_i) &= E[(y_i - E[y_i])(y_i - E[y_i])'] \\ &= E[y_i y_i'] \\ &= E[(\Lambda\eta_i + \epsilon_i)(\Lambda\eta_i + \epsilon_i)'] \\ &= \Lambda E[\eta_i \eta_i'] \Lambda' + 2\Lambda E[\eta_i \epsilon_i] + E[\epsilon_i \epsilon_i'] \\ &= \Lambda I \Lambda' + 0 + \Sigma_0 \\ &= \Lambda \Lambda' + \Sigma_0 \end{aligned}$$

$$\epsilon_i \sim N_p(0, \Sigma_0)$$

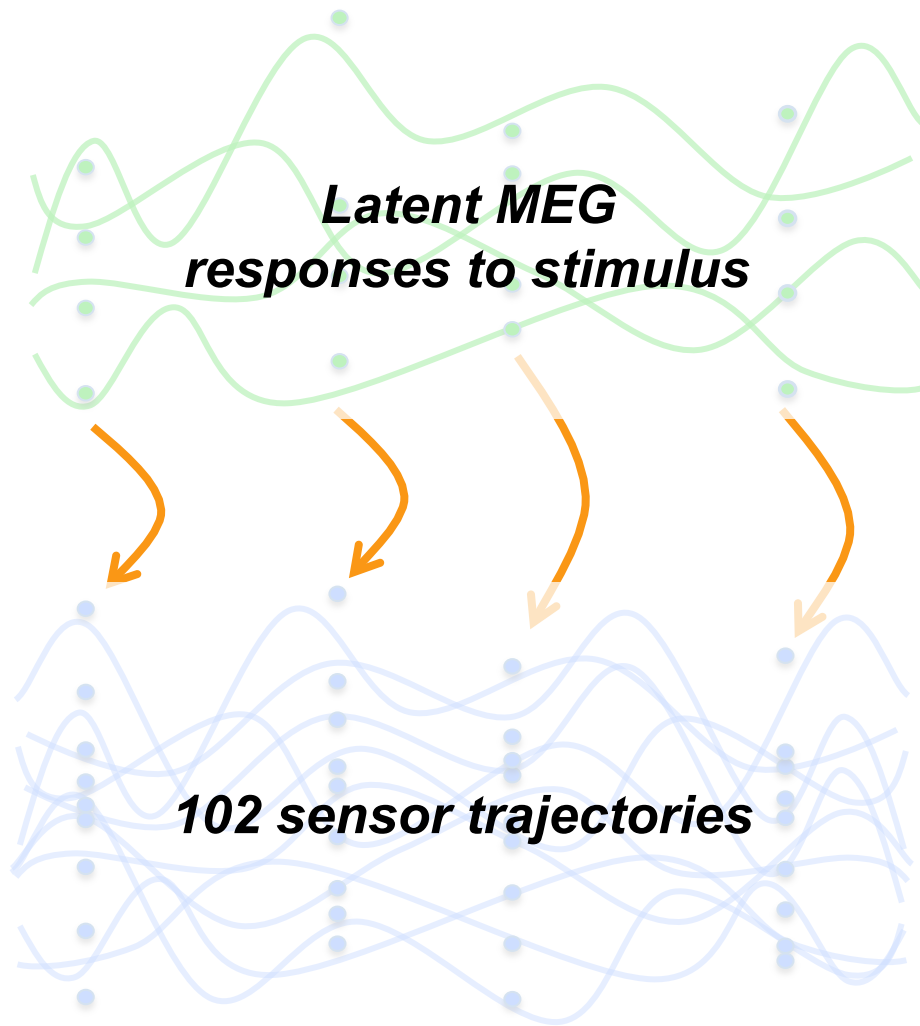
$$y_i = \Lambda \eta_i + \epsilon_i$$



# Adding Dynamics



# Dynamic Latent Factor Model



$$\eta_t = \Phi \eta_{t-1} + \nu_t$$

Evolution of latent factors

$\nu_t \sim N_k(0, I)$

$$y_t = \Lambda \eta_t + \epsilon_t$$

$\epsilon_t \sim N_p(0, \Sigma_0)$

# Dynamic Latent Factor Model

- State-space model with low-dim state and high-dim observations
- Originally developed by Geweke (1977)
  - Other early work:
    - Sargent and Sims (1977)
    - Watson and Engle (1983)
- Very popular in econometrics
- Most foundational dynamic model of high-dimensional time series

$$\eta_t = \Phi \eta_{t-1} + \nu_t$$

$N_k(0, I)$

Evolution of latent factors

$$y_t = \Lambda \eta_t + \epsilon_t$$

$N_p(0, \Sigma_0)$

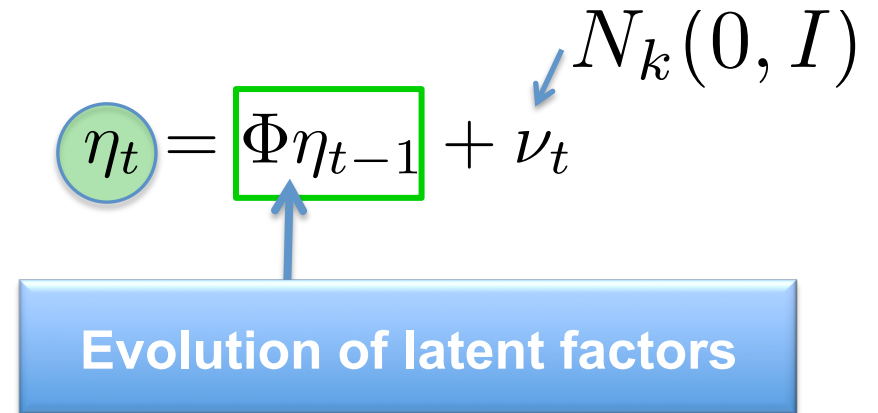
# Dynamic Latent Factor Model

- Assuming latent process is *stable*, marginally

$$y_t \sim N(0, \Sigma)$$

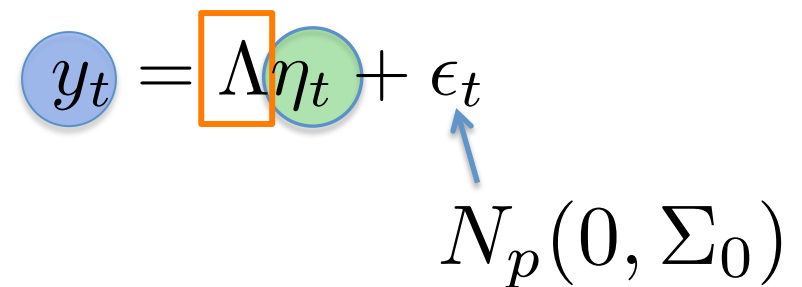
$$\Sigma = \Lambda \Sigma_{\eta} \Lambda' + \Sigma_0$$

$\Gamma_{\eta}(0)$



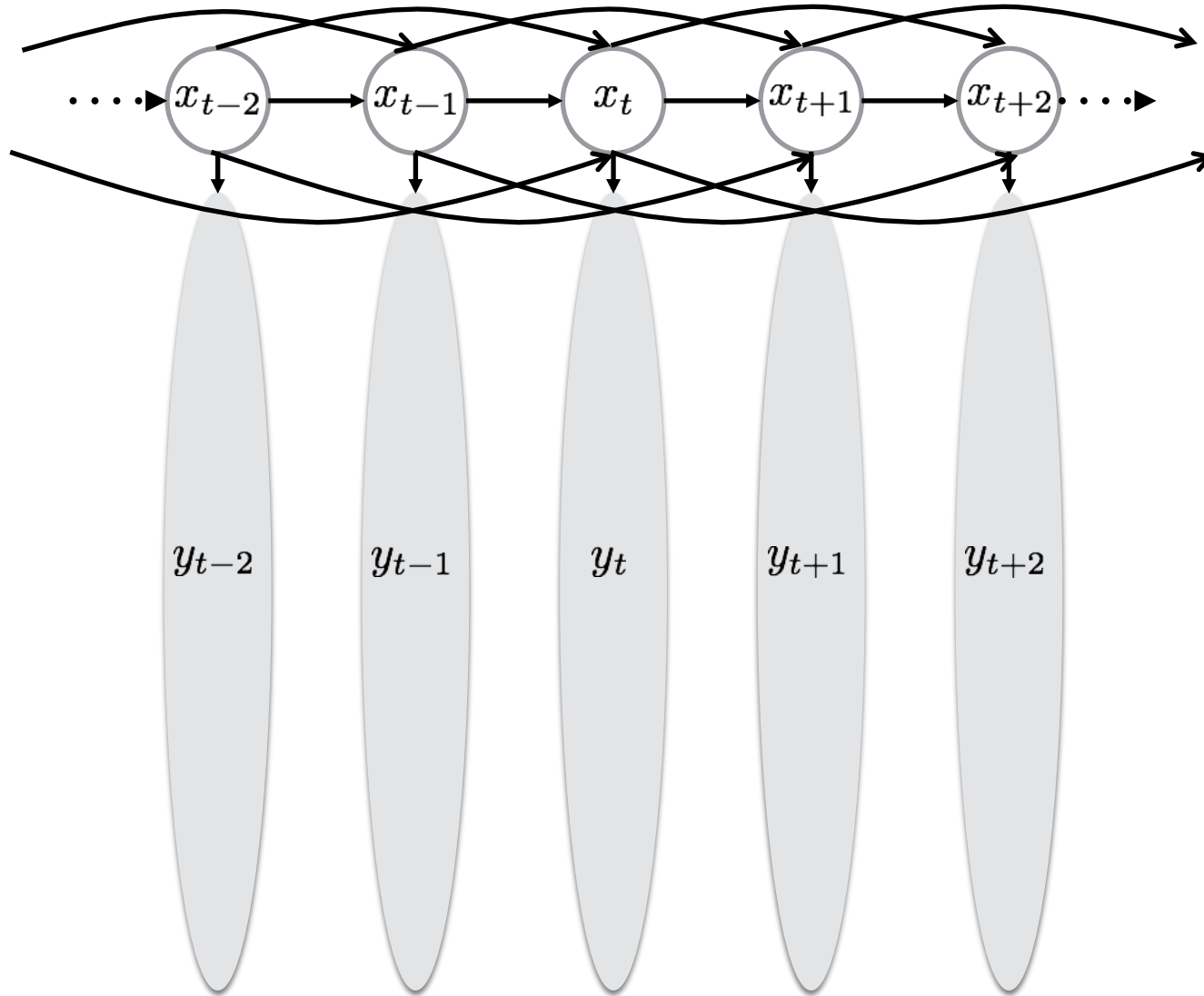
- Though, still a dynamic process with lag covariance

$$\begin{aligned} \Gamma_y(h) &= \text{cov}(y_t, y_{t+h}) \\ &= \Lambda \Gamma_{\eta}(h) \Lambda' \quad h > 0 \end{aligned}$$



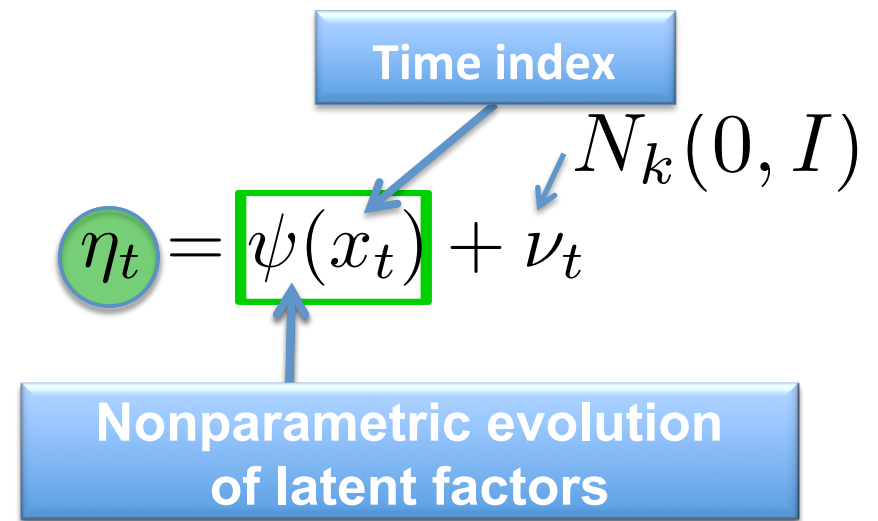


# Adding *Complex* Dynamics



# Semiparametric Factor Model

- Can consider a *nonparametric latent factor process*
  - Gaussian processes
  - More on next slides...
- For a regression setting, looks very similar to  
**Teh, Seeger, & Jordan 2004**
  - $x$  an arbitrary covariate, not necessarily time



The diagram illustrates the regression equation:  $y_t = \Lambda \eta_t + \epsilon_t$ . A blue box labeled "Time index" has an arrow pointing to the  $\eta_t$  term, which is enclosed in a blue box. A blue arrow points from the  $\epsilon_t$  term to the notation  $N_p(0, \Sigma_0)$ .

# Gaussian Processes

**REVIEW**

- Distribution on functions

- $f \sim \text{GP}(m, k)$

- $m$ : mean function
    - $k$ : covariance function



- $p(f(x_1), \dots, f(x_n)) \sim N_n(\mu, K)$

- $\mu = [m(x_1), \dots, m(x_n)]$
    - $K_{ij} = k(x_i, x_j)$

- Idea: If  $x_i, x_j$  are similar according to the kernel, then  $f(x_i)$  is similar to  $f(x_j)$

# Gaussian Processes

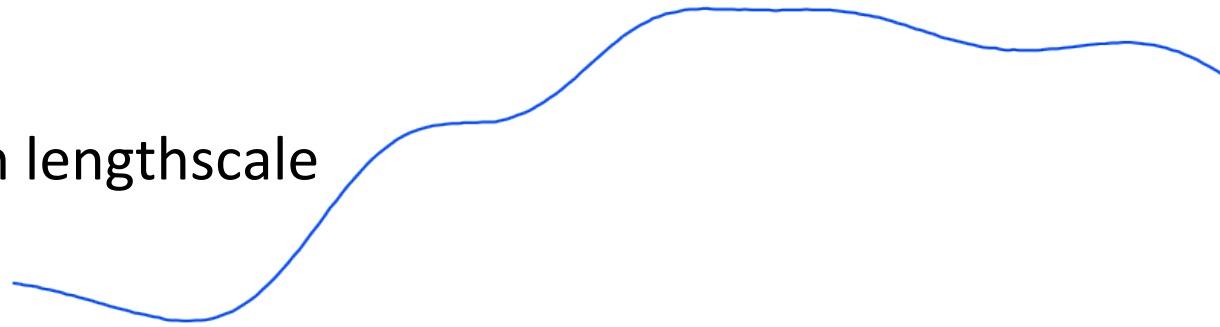
**REVIEW**

**$\kappa$** : covariance function

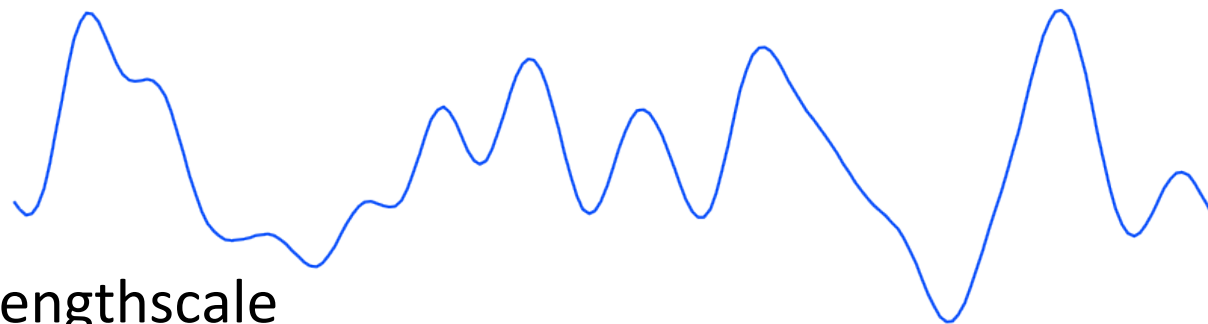
$f \sim \text{GP}(\mathbf{m}, \mathbf{\kappa})$

$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$

High lengthscale



Low lengthscale

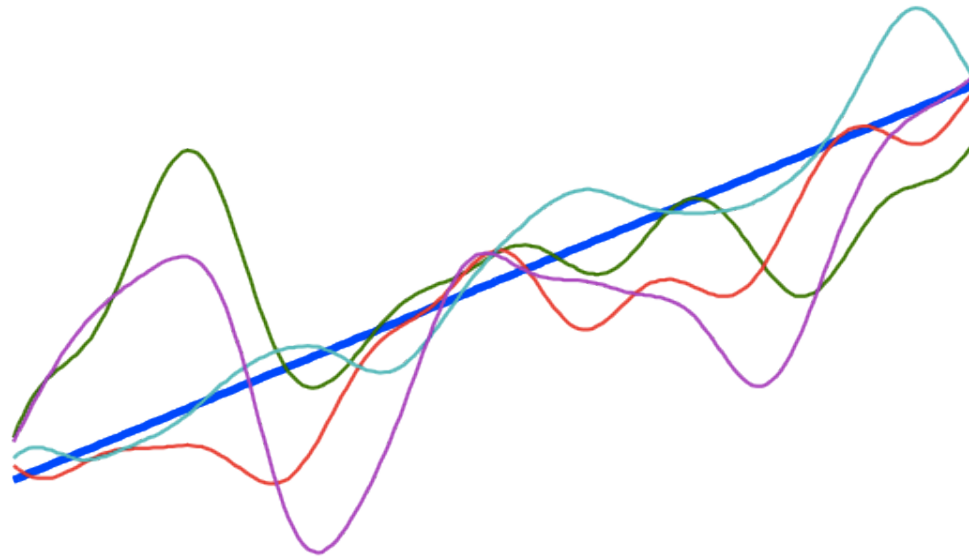


# Gaussian Processes

**REVIEW**

$m$ : mean function

$$f \sim \text{GP}(m, k)$$

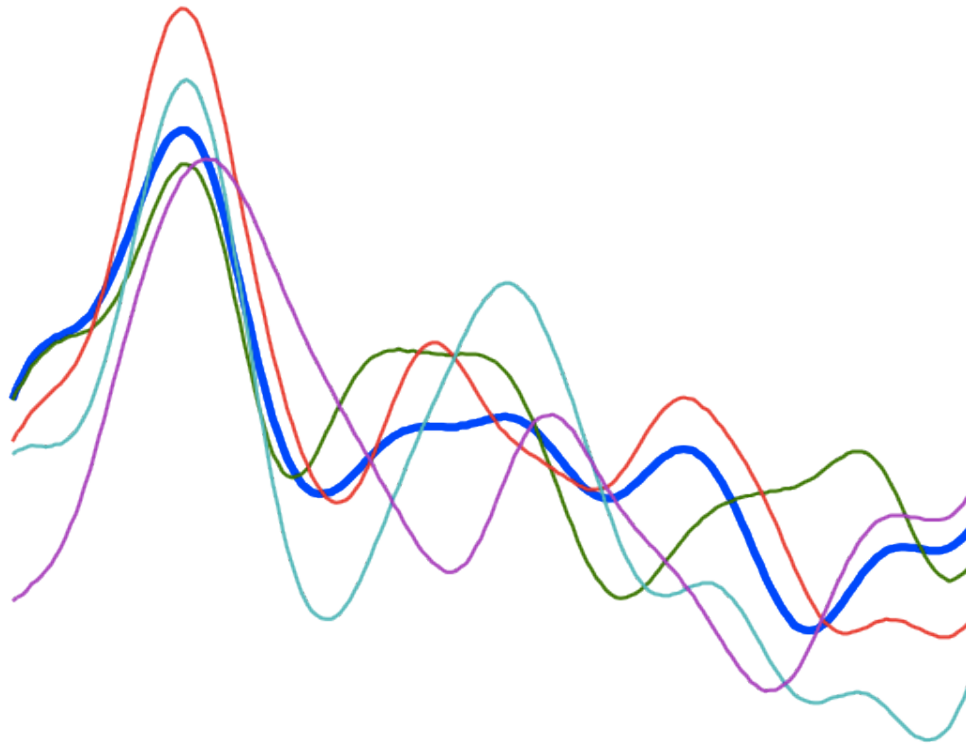


# Gaussian Processes

**REVIEW**

$m$ : mean function

$$f \sim \text{GP}(m, k)$$

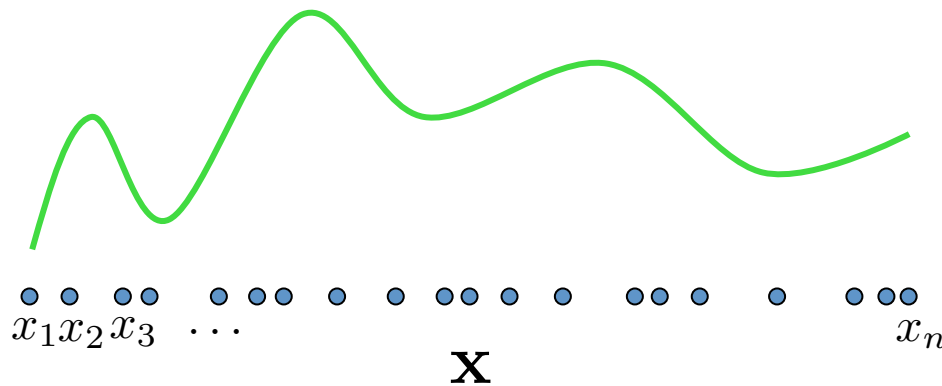


# Induced Multivariate Gaussian

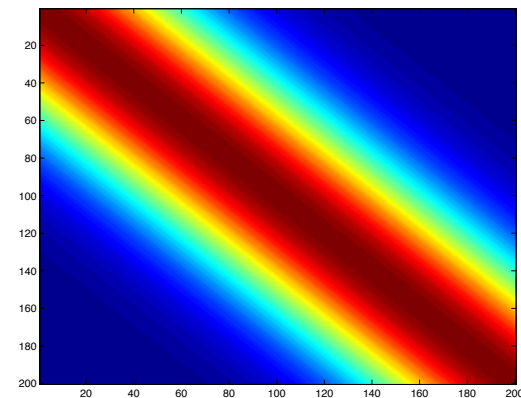
**REVIEW**

- Evaluating the GP-distributed function at any set of locations  $(x_1, \dots, x_n)$ , we have

$$\begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix} \sim N(\mu, K)$$



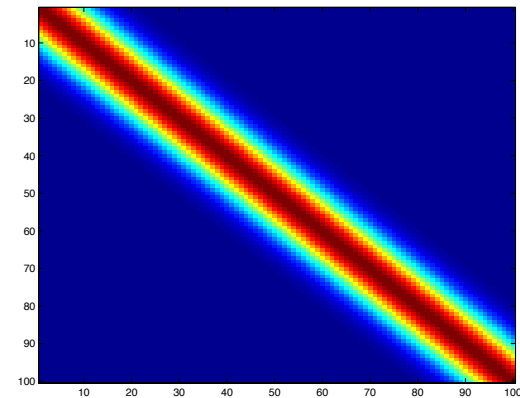
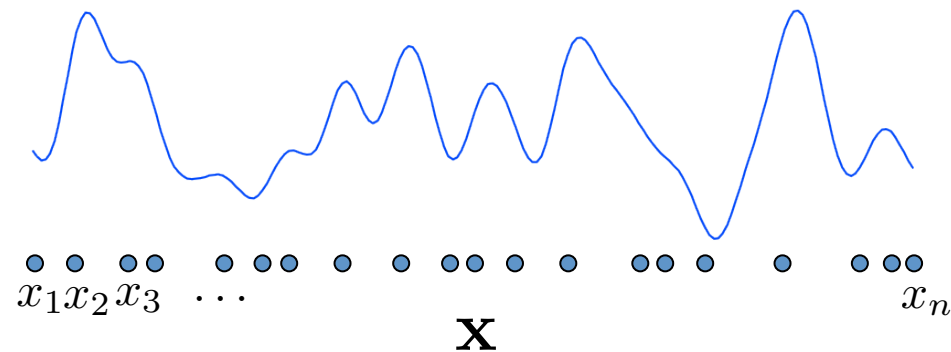
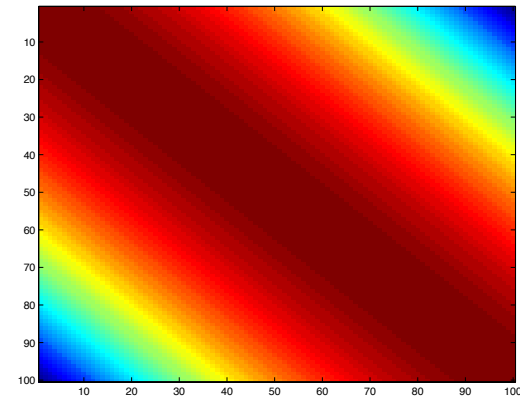
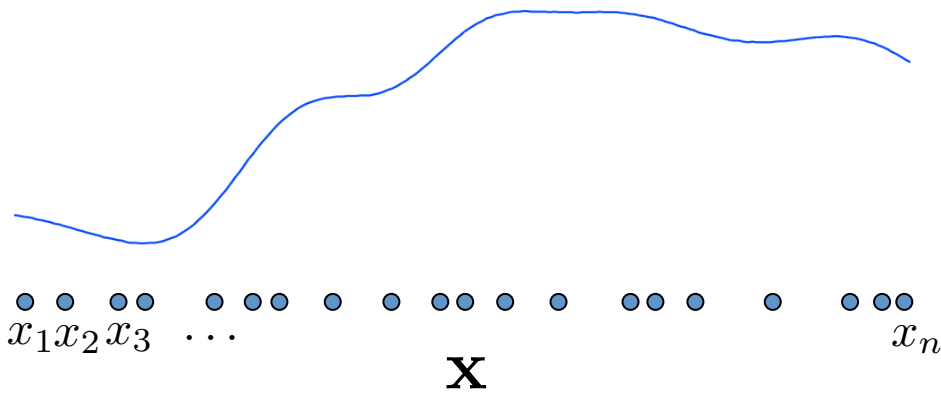
$$K =$$



# Induced Multivariate Gaussian

**REVIEW**

- Comparing length-scales:

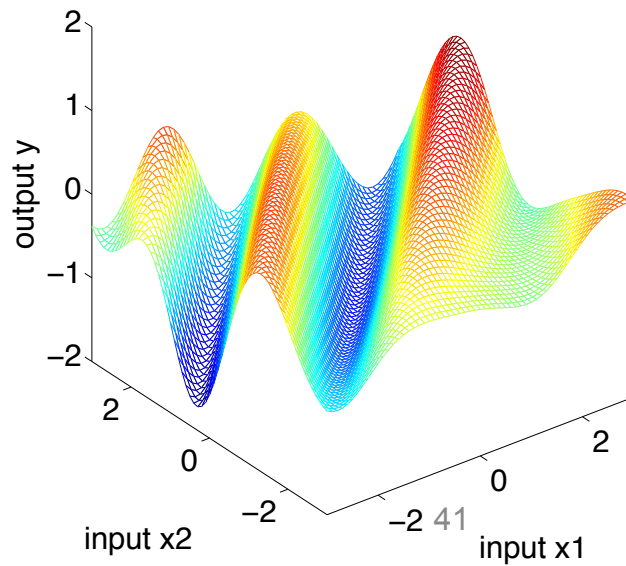
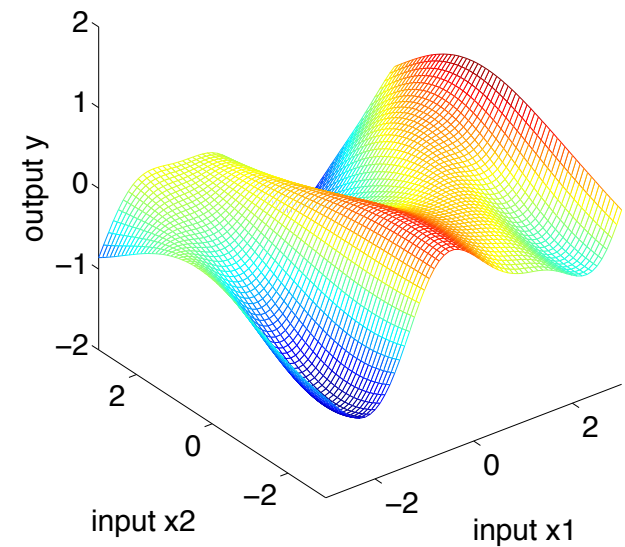
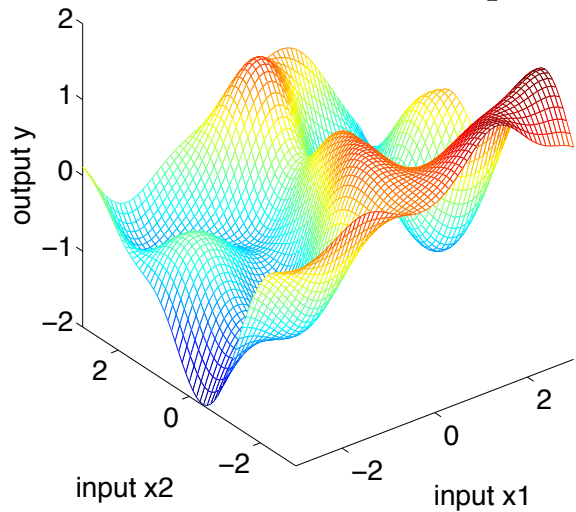




# 2D Gaussian Processes

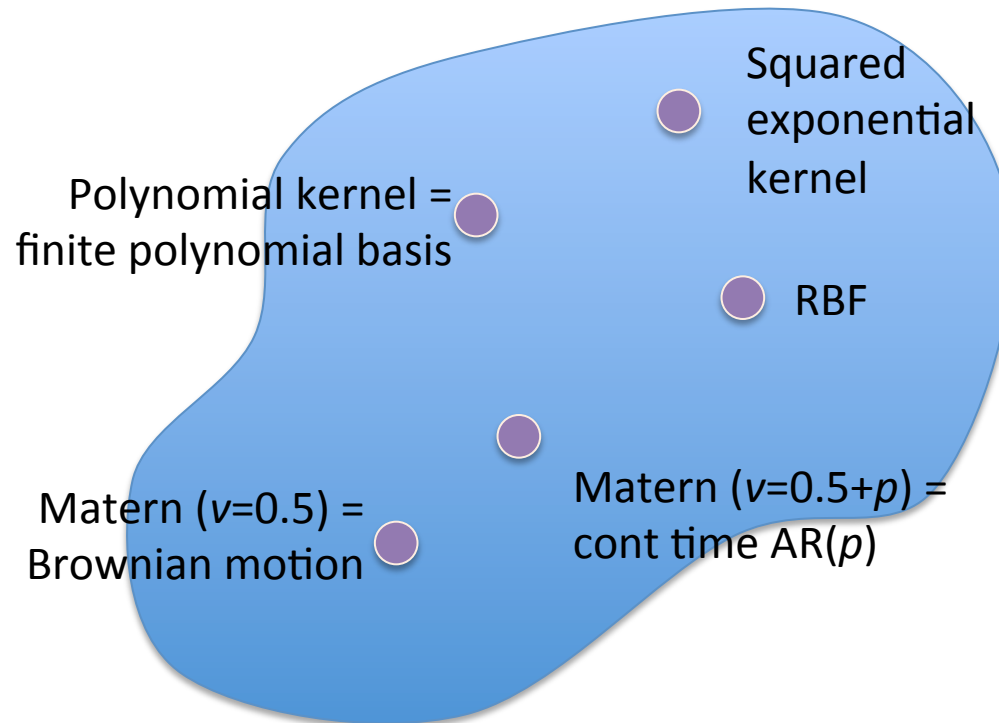
**REVIEW**

$$\kappa(x_p, x'_q) = \frac{2}{f} \exp\left(-\frac{1}{2}(x_p - x'_q)^T M(x_p - x'_q)\right)$$



# Family of Gaussian Processes

**REVIEW**



# GPs for Regression


**REVIEW**

- Start with noise-free scenario: directly observe the function
- Training data  $\mathcal{D} = \{(x_i, f_i), i = 1, \dots, n\}$
- Test data locations  $X^* \rightarrow$  predict  $f^*$


- Jointly, we have

$$\begin{pmatrix} f \\ f^* \end{pmatrix} \sim N \left( \begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} K & K_* \\ K_*^T & K_{**} \end{pmatrix} \right)$$

$\kappa(X, X^*)$



$\kappa(X^*, X^*)$

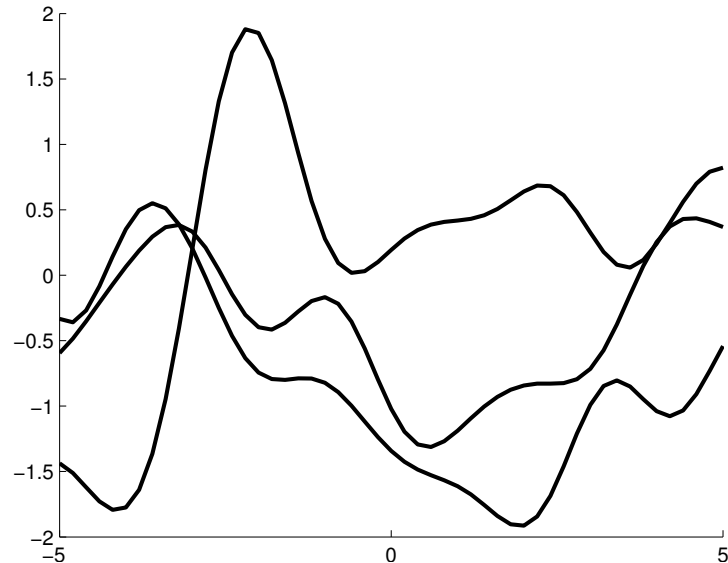


- Therefore,

$$p(f^* | X^*, X, f) = N(f^* | \mu_* + K_*' K^{-1} (f - \mu), K_{**} - K_*' K^{-1} K_*)$$

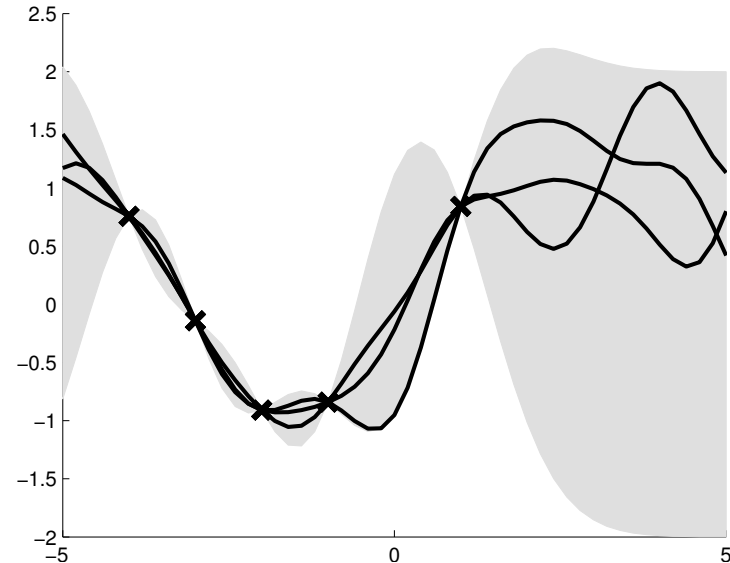
# 1D Noise-Free Example

**REVIEW**



*Samples from Prior*

$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$



*Posterior Given 5  
Noise-Free Observations*

- Interpolator, where uncertainty increases with distance
- Useful as a computationally cheap proxy for a complex simulator
  - Examine effect of simulator params on GP predictions instead of doing expensive runs of the simulator

# GPs for Regression



- Noisy scenario: observe a noisy version of underlying function

$$y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma_y^2)$$

- Not required to interpolate, just come “close” to observed data

$$\text{cov}(y|X) = \text{cov}(f) + \text{cov}(e) = K + \sigma_y^2 I \triangleq K_y$$

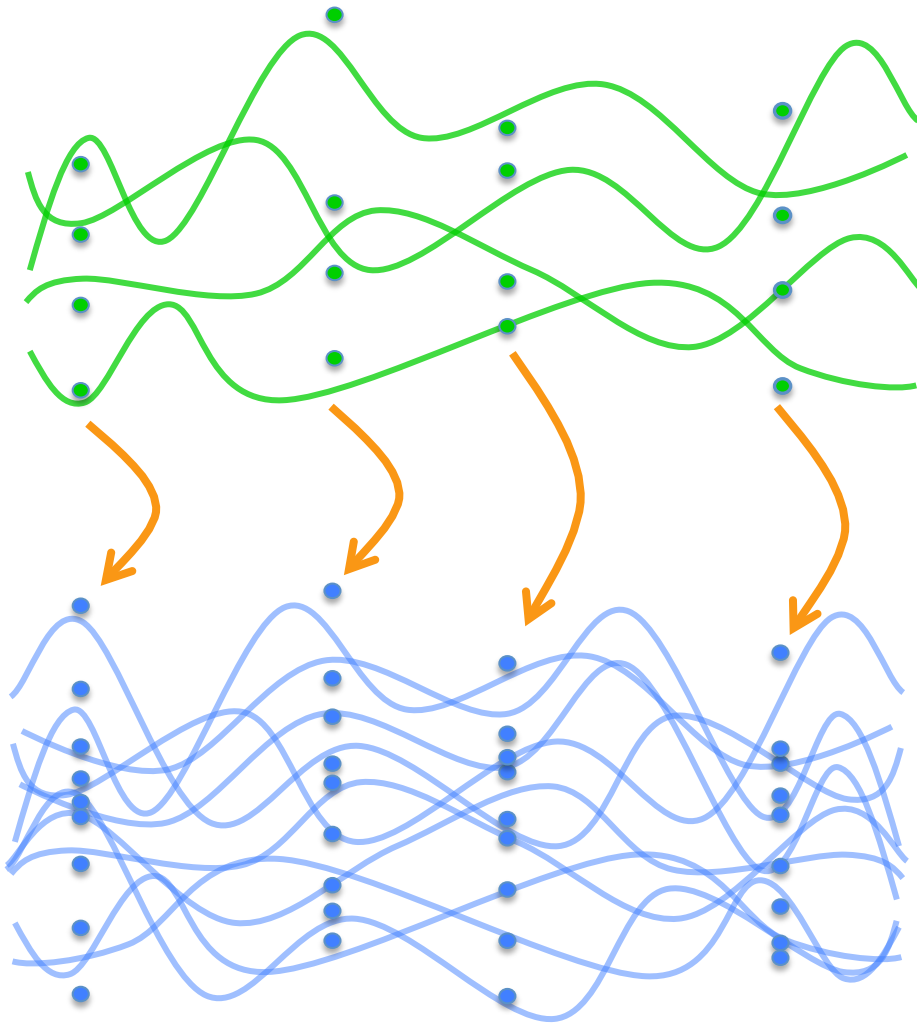
- Training data  $\mathcal{D} = \{(x_i, y_i), i = 1, \dots, n\}$
- Test data locations  $X^* \rightarrow$  predict  $f^*$

- Jointly, we have  $\begin{pmatrix} y \\ f^* \end{pmatrix} \sim N \left( 0, \begin{pmatrix} K_y & K_* \\ K_*^T & K_{**} \end{pmatrix} \right)$

- Therefore,

$$p(f^* | X^*, X, y) = N(f^* | K_*' K_y^{-1} y, K_{**} - K_*' K_y^{-1} K_*)$$

# Dynamic Latent Factor Model



$$\eta_t = f(\eta_{1:t-1}) + \nu_t$$

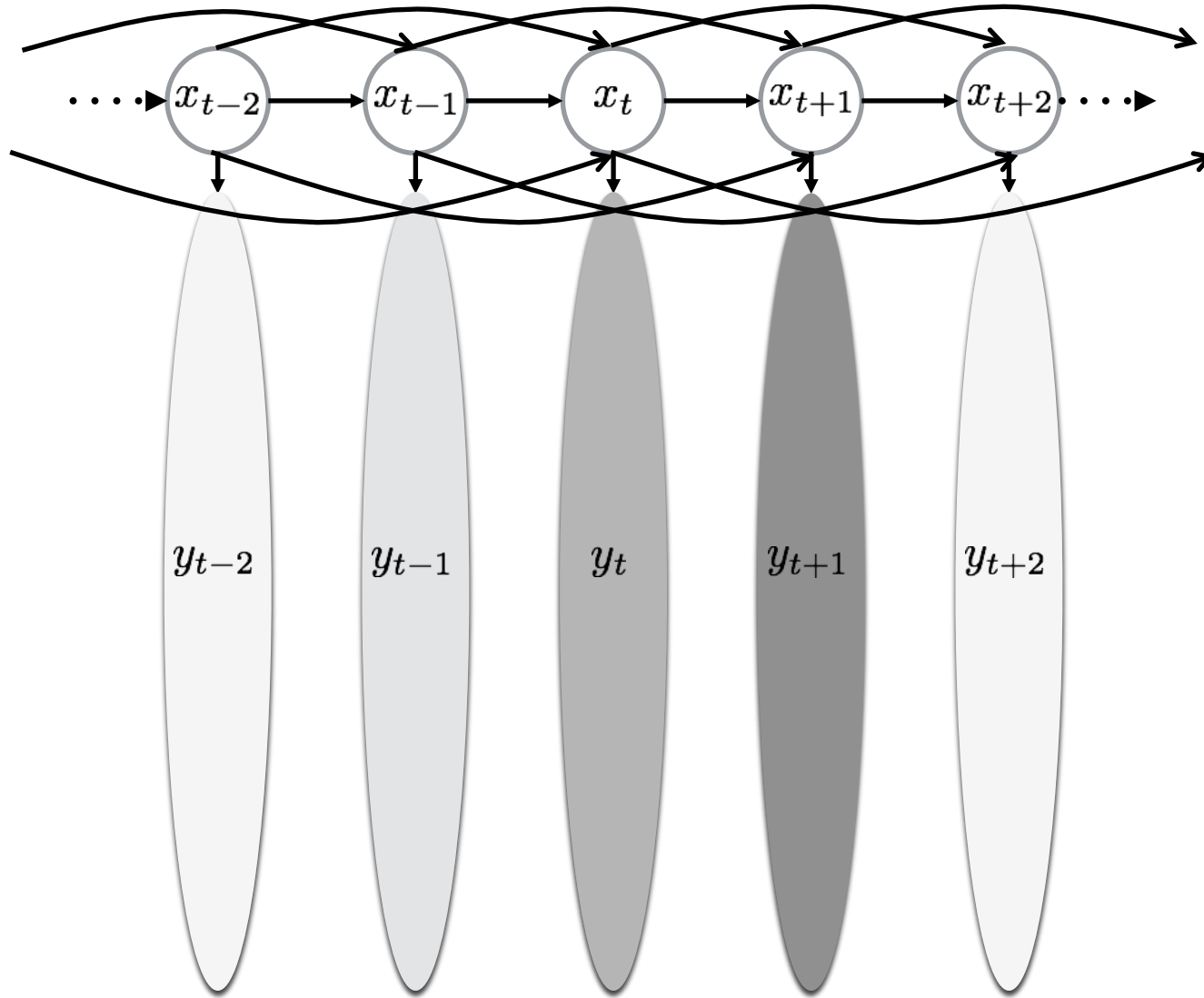
$N_k(0, I)$

Evolution of latent factors

$$y_t = \Lambda \eta_t + \epsilon_t$$

$N_p(0, \Sigma_0)$

# Capturing Changing Correlations



# Capturing Changing Correlations

Observation:

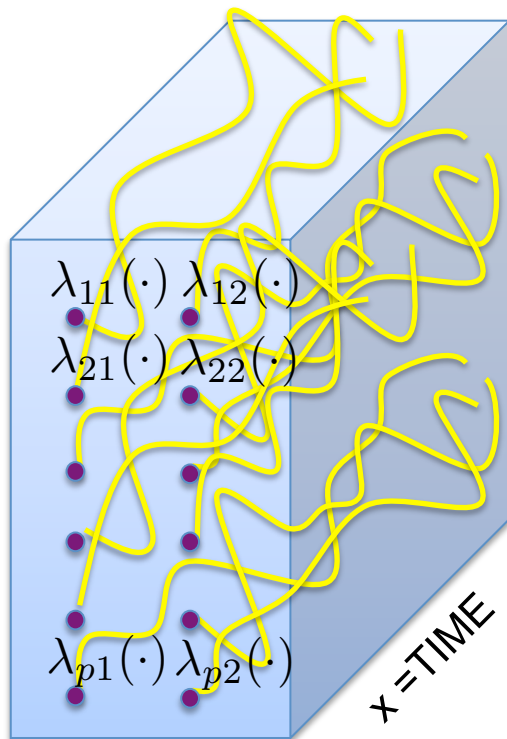
1. Sensors are *redundant*
2. Correlation pattern *changes* with time

$$\Sigma(x) = \overset{p \times k}{\Lambda(x)\Lambda(x)'} + \Sigma_0$$

The diagram illustrates the equation  $\Sigma(x) = \Lambda(x)\Lambda(x)' + \Sigma_0$ . A blue box labeled "Time index" points to the  $\Sigma(x)$  term. The  $\Lambda(x)\Lambda(x)'$  term is circled in yellow. Below the equation, a blue square is shown to be equal to a blue rectangle plus a blue square with a diagonal line and red dots.



# Low-Rank Covariance Evolution



$p \times k$  array of  
processes over time

$$k \ll p$$

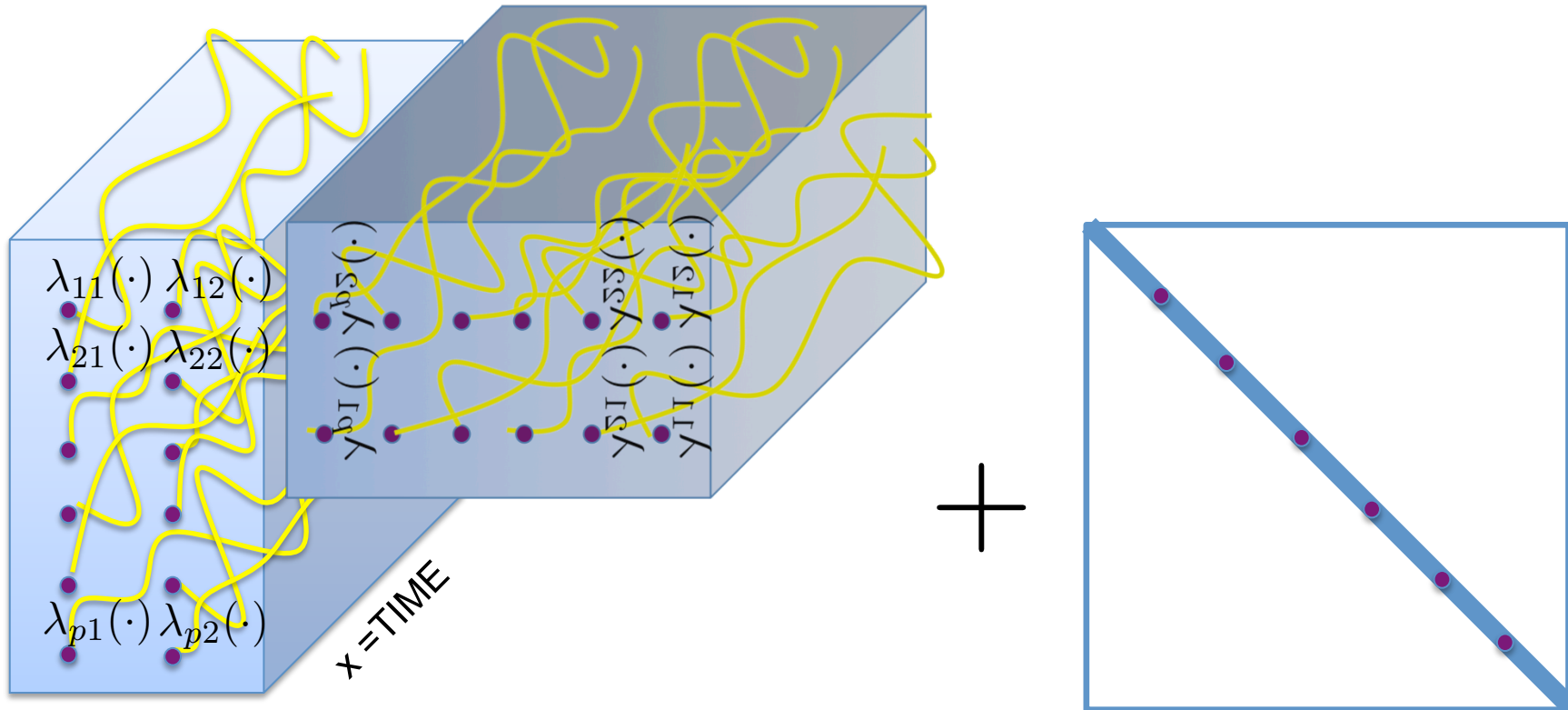
$$p \times k$$

$$\Sigma(x) = \boxed{\Lambda(x)} \Lambda(x)' + \Sigma_0$$

Fox and Dunson, arXiv 2011.

Related model without low-rank structure: Wilson and Ghahramani, UAI 2011.

# Low-Rank Covariance Evolution

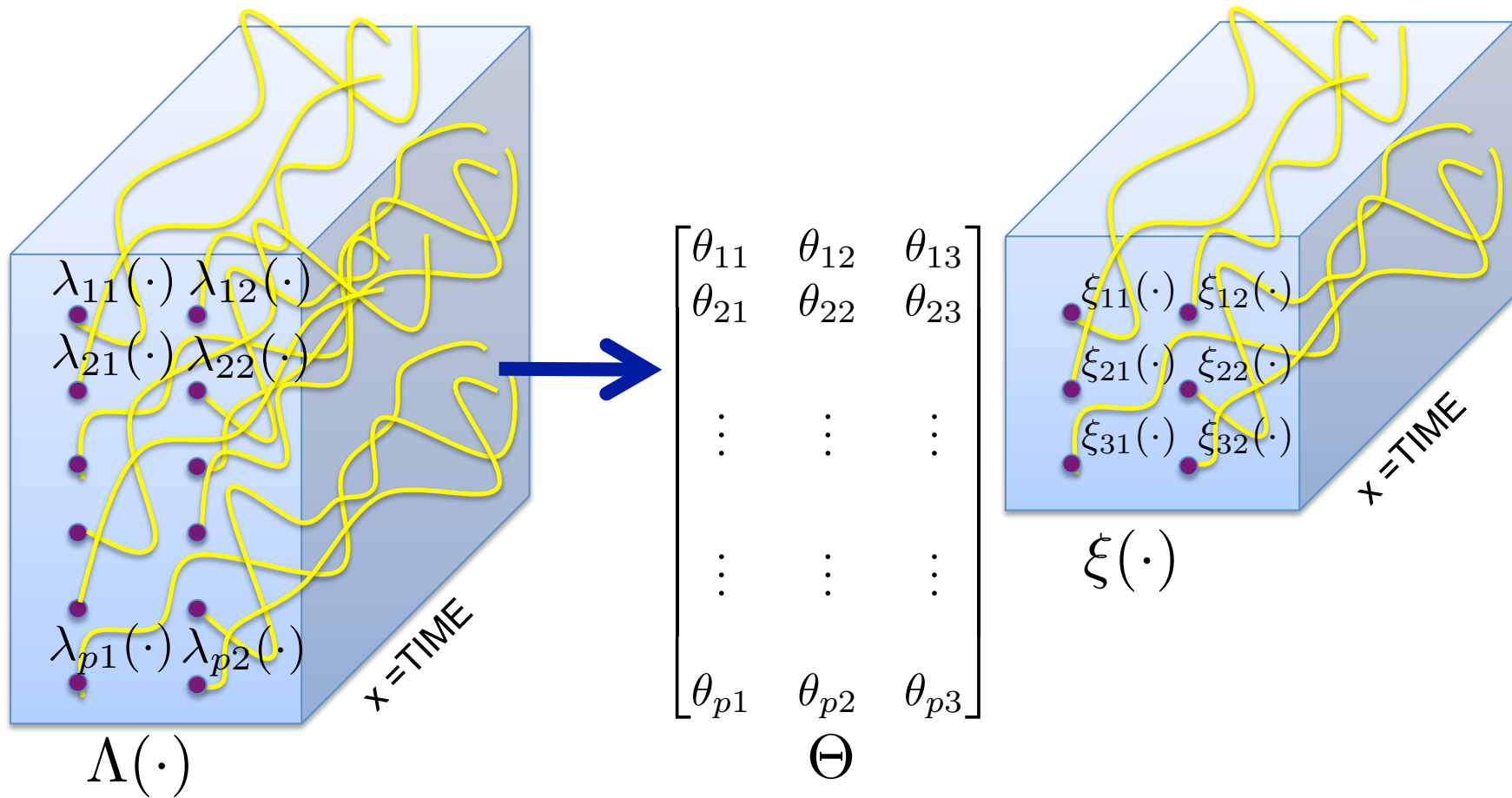


$$\Sigma(x) = \Lambda(x)\Lambda(x)' + \Sigma_0$$

Fox and Dunson, arXiv 2011.

Related model without low-rank structure: Wilson and Ghahramani, UAI 2011.

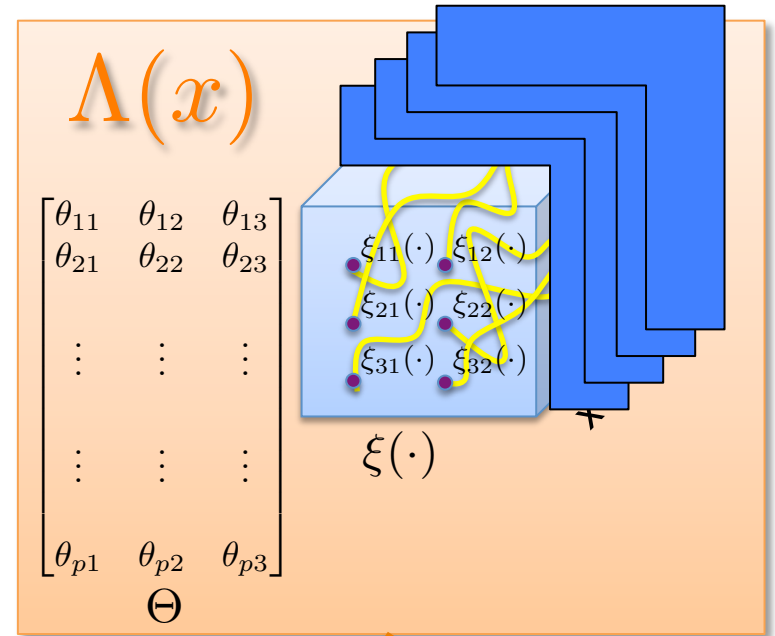
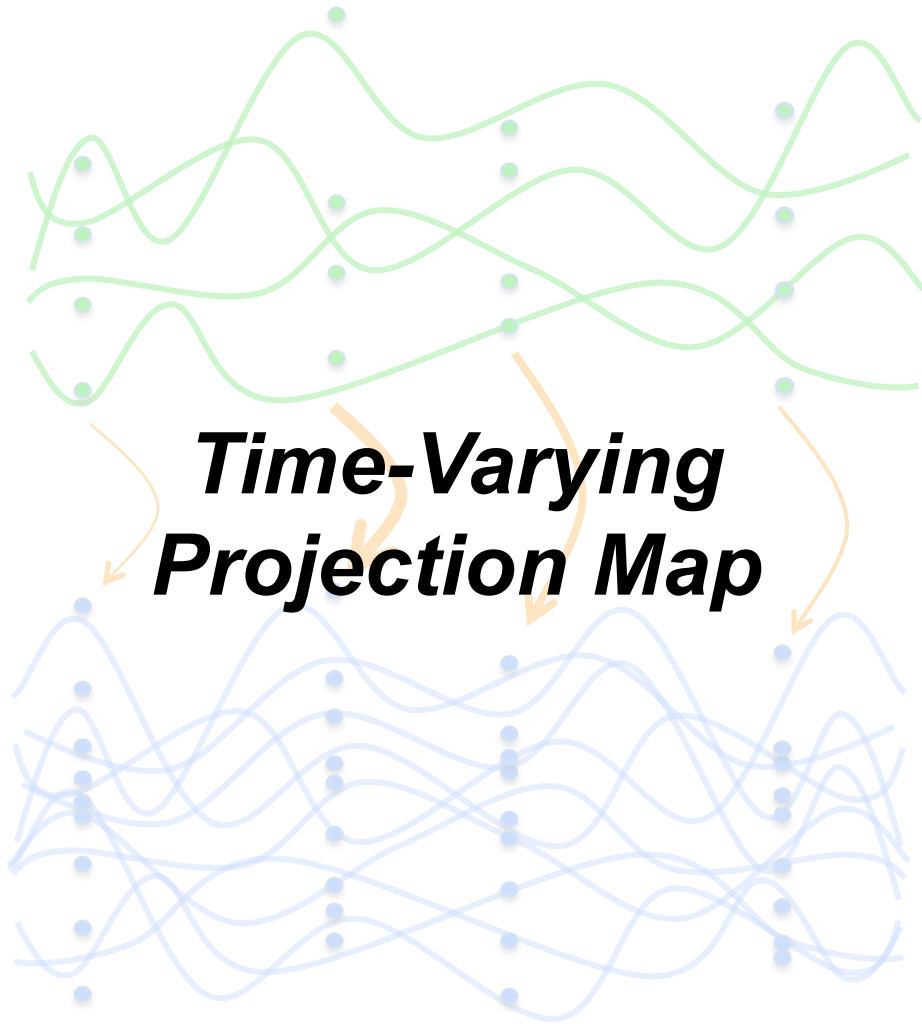
# One Step Further...



$$\Sigma(x) = \Theta \xi(x) \xi(x)' \Theta' + \Sigma_0$$

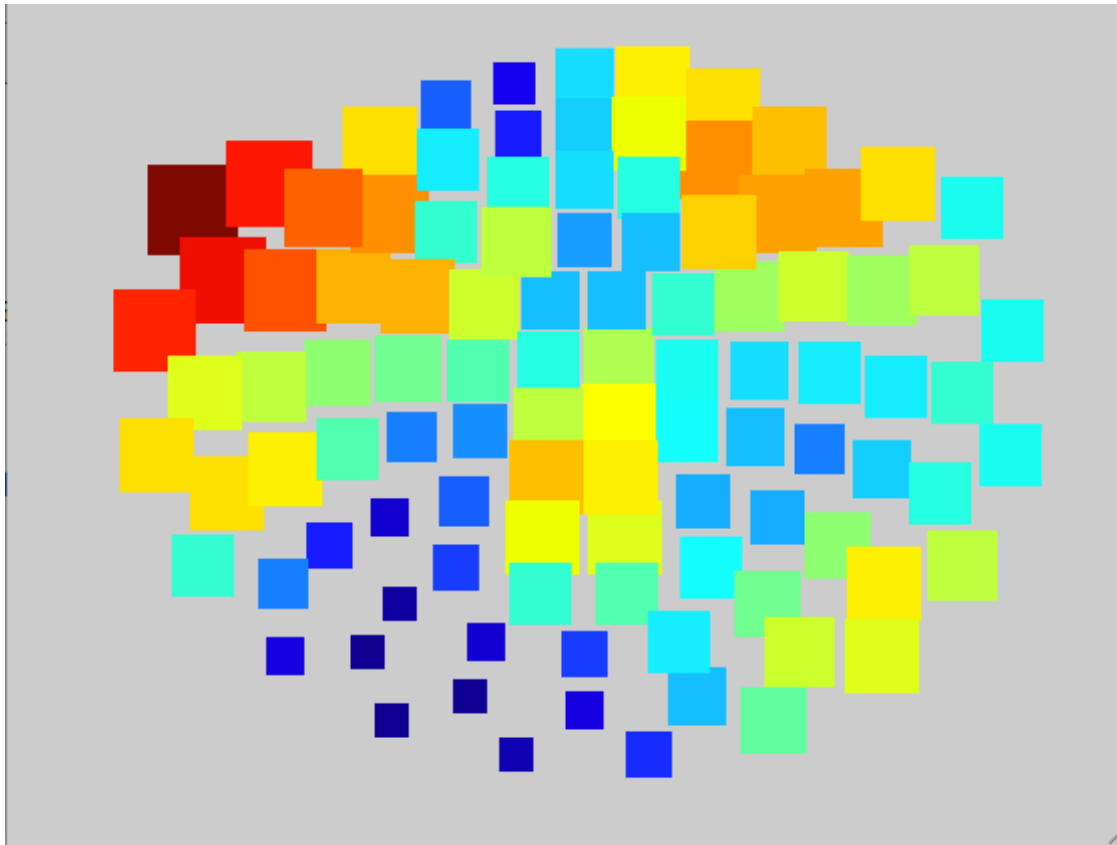
Fox and Dunson, arXiv 2011.

# Interpretation as Dynamic LFM

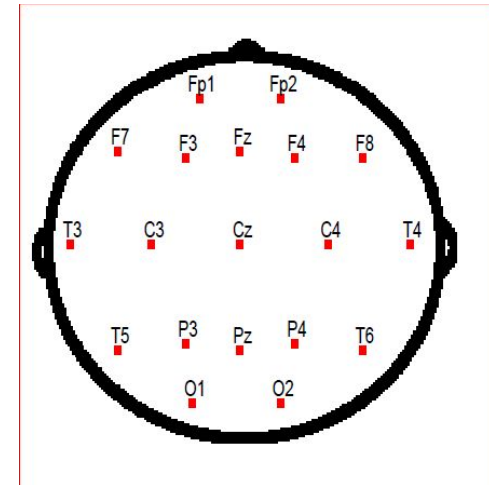


$$y_t = \Theta \xi(x_t) \eta_t + N_p(0, \Sigma_0)$$

# Changing Correlations – MEG



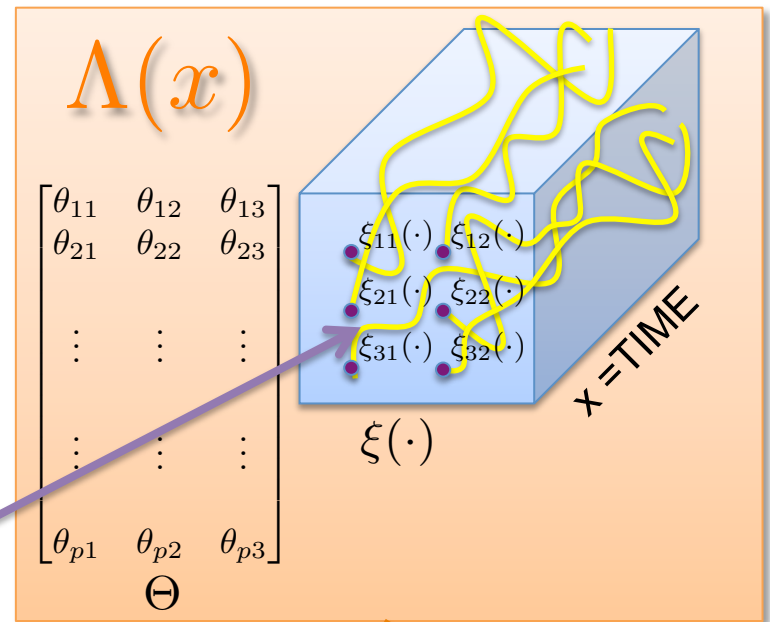
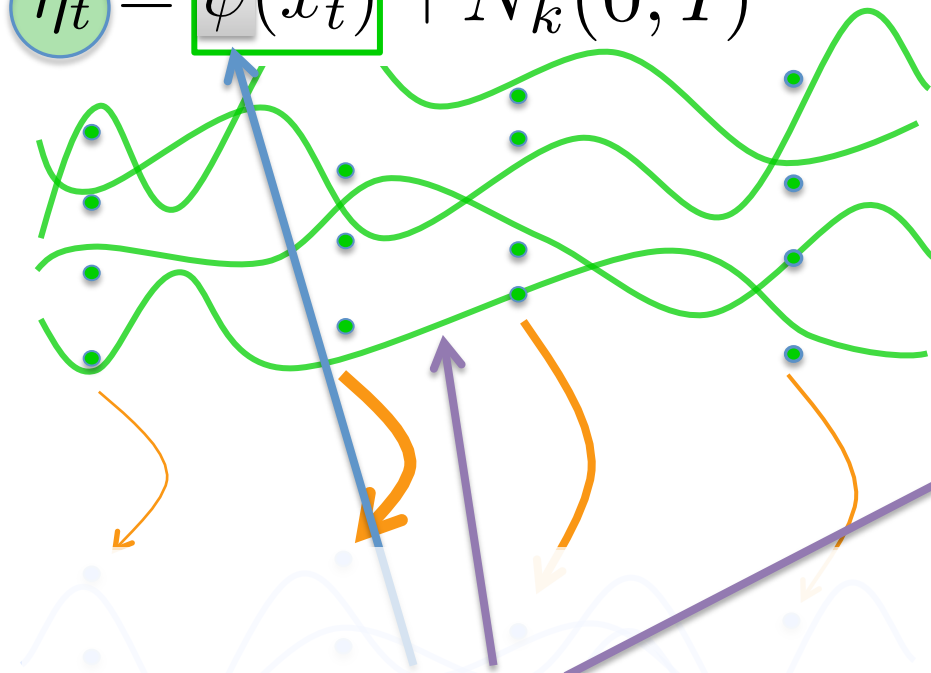
102 sensors:



*Correlations between sensors change with processing of word "kick"*

# Prior Specification

$$\eta_t = \psi(x_t) + N_k(0, I)$$



$$y_t = \Theta \xi(x_t) \eta_t$$

$$+ N_p(0, \Sigma_0)$$

**Gaussian  
Processes**

# Data Collection

- 4 word categories, 5 words per category



Animals



Food

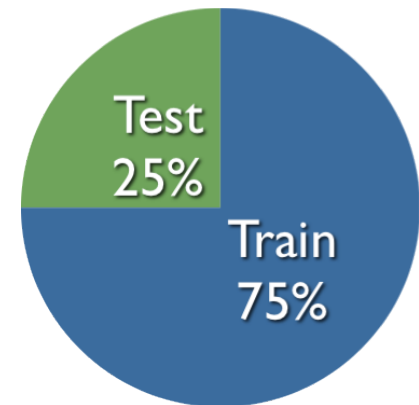


Tools

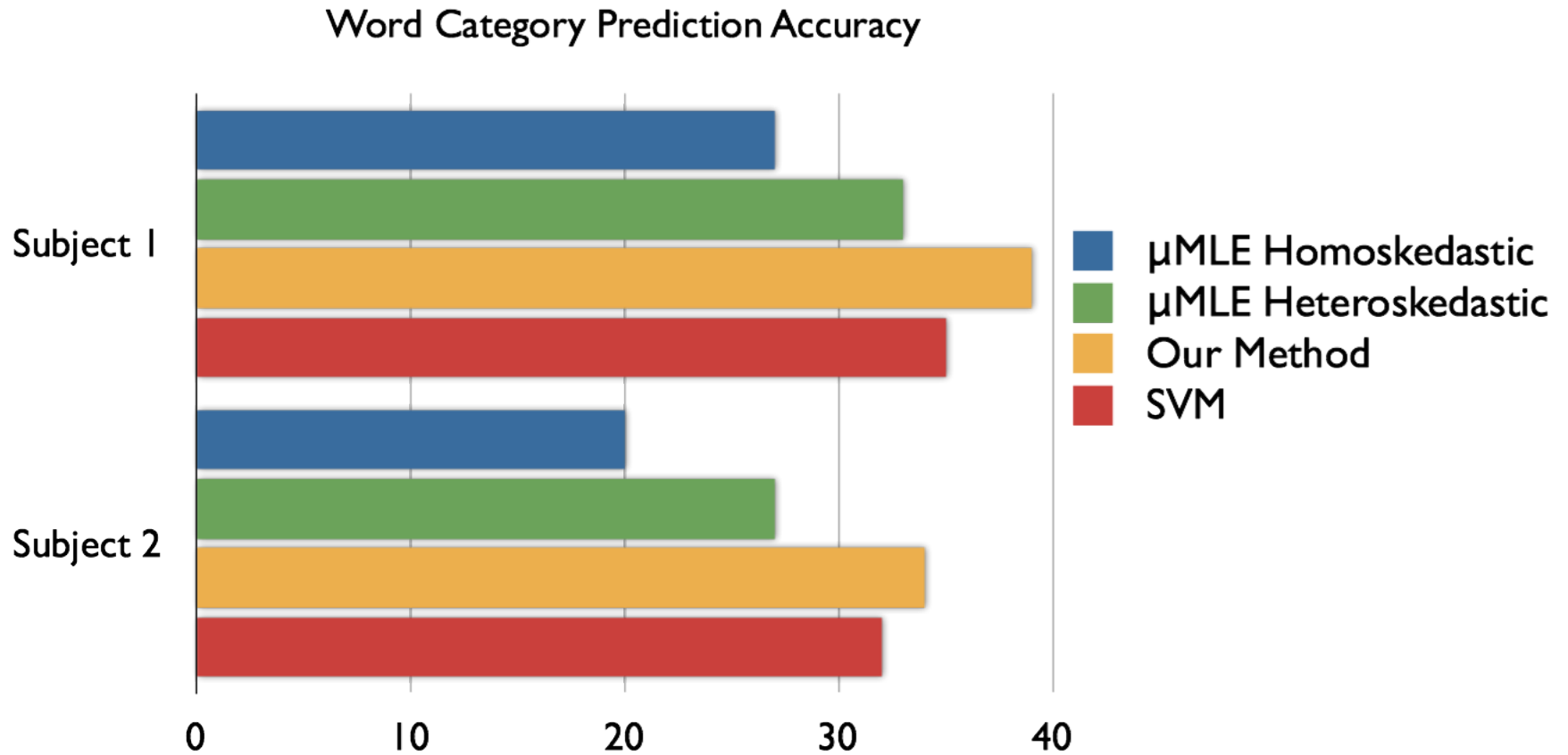


Buildings

- 20 repetitions per word (400 total)
  - 15 train/word (300 total)
  - 5 test/word (100 total)



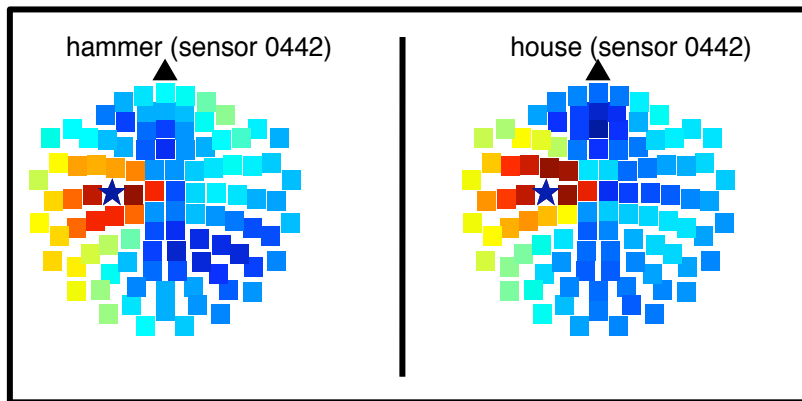
# Classification Performance





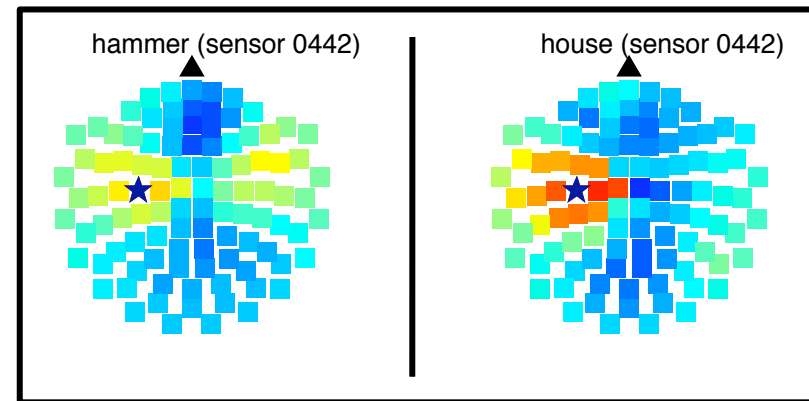
# Perceptual vs. Semantic Correlations

## *Perceptual*

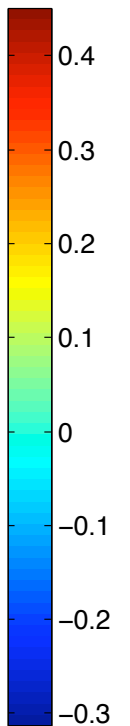


**t = 101 ms**

## *Semantic*

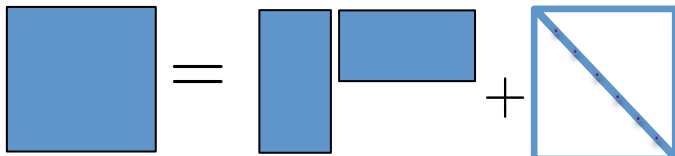


**t = 401 ms**

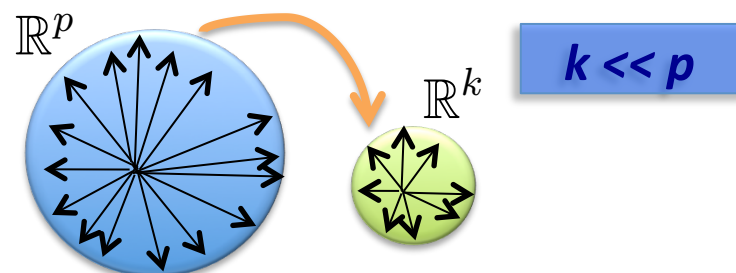


# Low-Dim Embedding Summary

Low Rank  $\Sigma = \Lambda\Lambda' + \Sigma_0$



The diagram illustrates the equation  $\Sigma = \Lambda\Lambda' + \Sigma_0$ . On the left, a solid blue square represents the matrix  $\Sigma$ . This is followed by an equals sign. To the right of the equals sign, there is a blue vertical rectangle representing  $\Lambda$ , followed by a blue horizontal rectangle representing  $\Lambda'$ , and then a plus sign. To the right of the plus sign is a square with a blue diagonal line representing  $\Sigma_0$ .

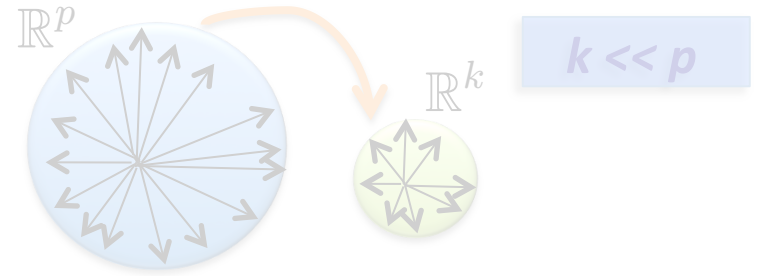
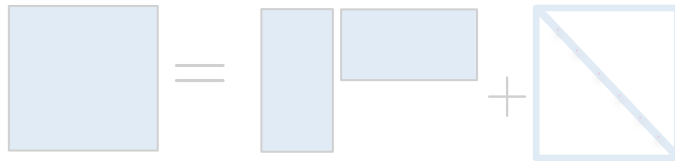


*Low-dimensional embedding*

- Latent factor models
  - Low-rank covariance approximation to high-dim i.i.d. Gaussian observations
- Dynamic latent factor model
  - Interpretation as state space model with low-dim state
  - Many approaches to modeling latent dynamics, including Gaussian processes
- Capturing changing correlations in high-dim setting
  - Factor structure within dynamic latent factor model
  - Gaussian process “dictionary” functions

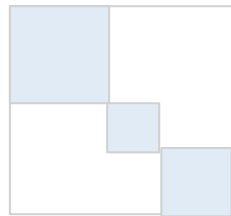
# Methods for Scaling to High Dimensions

Low Rank  $\Sigma = \Lambda\Lambda' + \Sigma_0$



Low-dimensional embedding

$\Sigma$  sparse



Independent groups of nodes

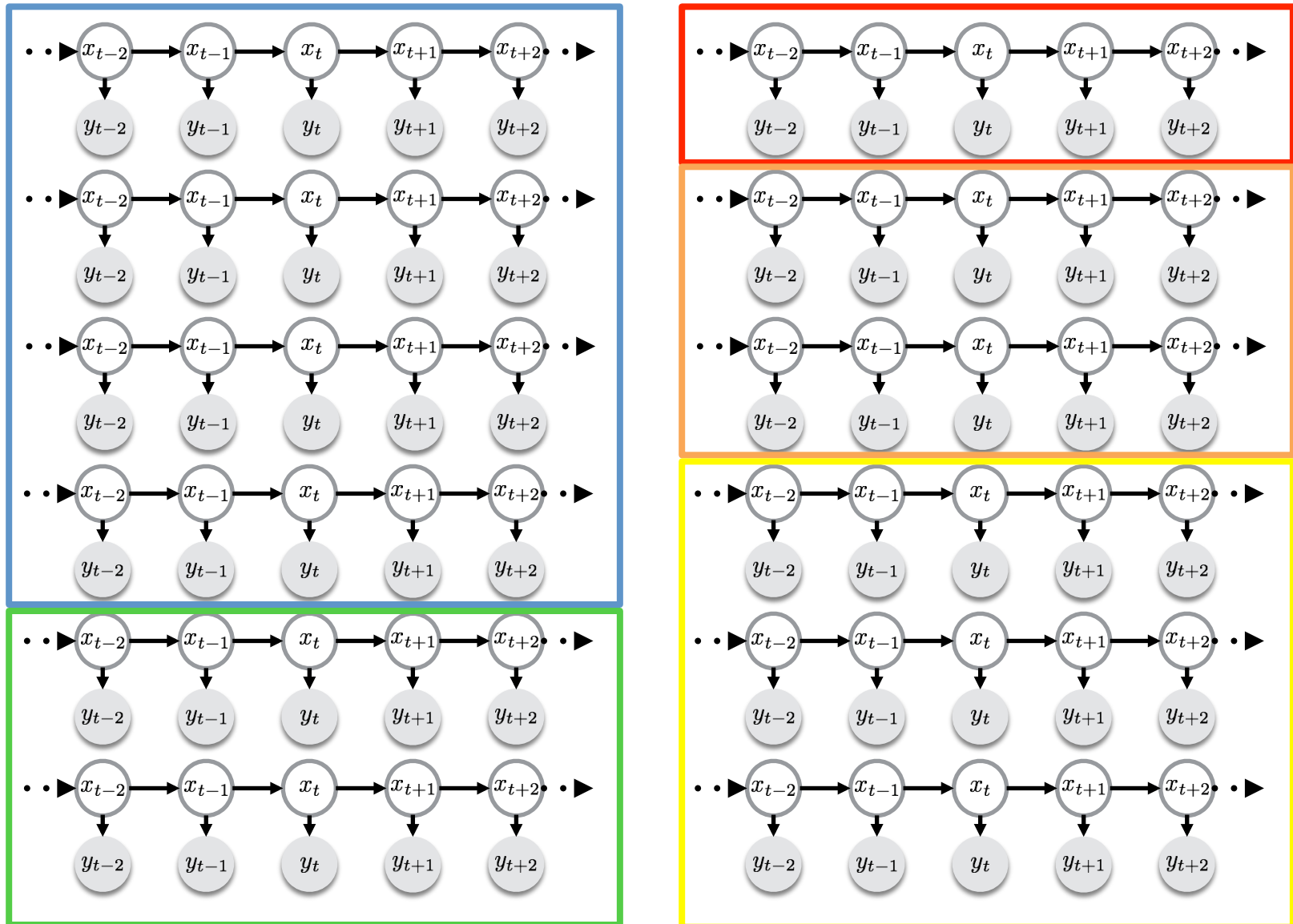
$\Sigma^{-1}$  sparse

Gaussian Graphical Model



Zeros = no edge in graph,  
Cond. ind. between nodes

# Clustering Time Series



# High-Resolution Housing Price Index

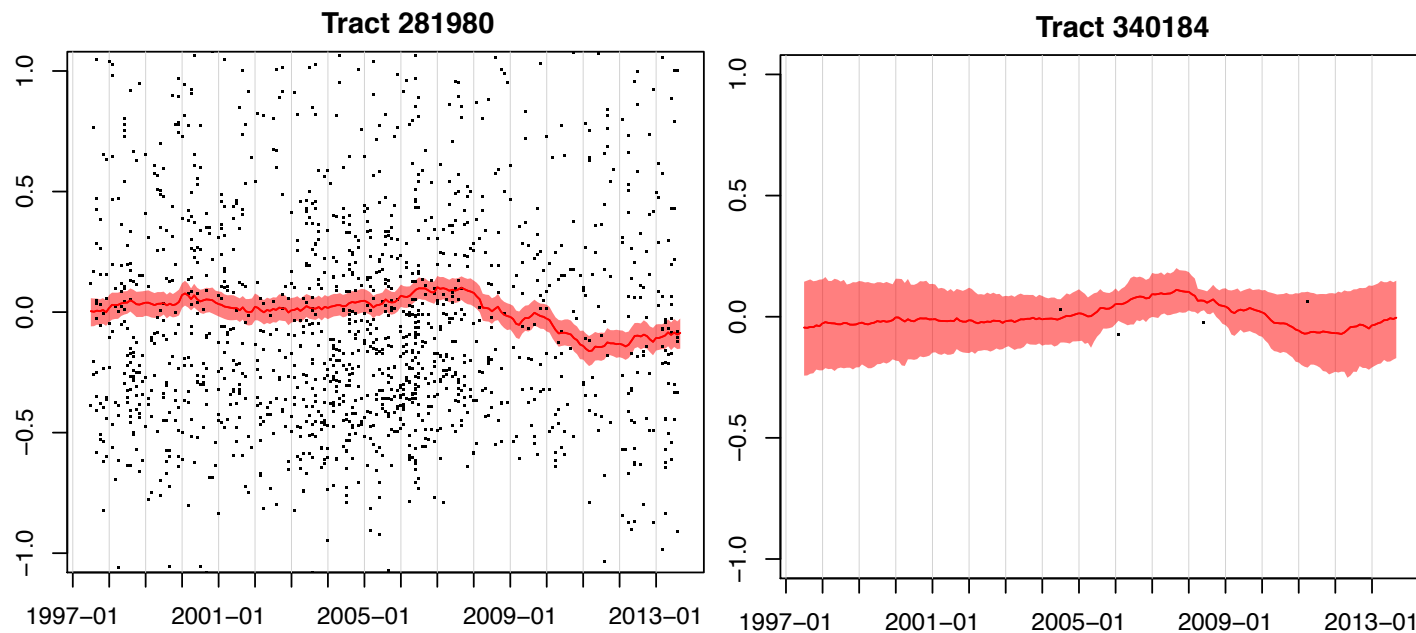
- **Goal:** Model neighborhood housing value over time based on observed house sales (with covariates)



# Challenge

**Issue:** Data are spatiotemporally sparse

Average monthly sales	< 1	< 3	< 5	< 7	< 9
Number of tracts	16	58	114	136	139
Percentage of tracts	11	41	81	97	99

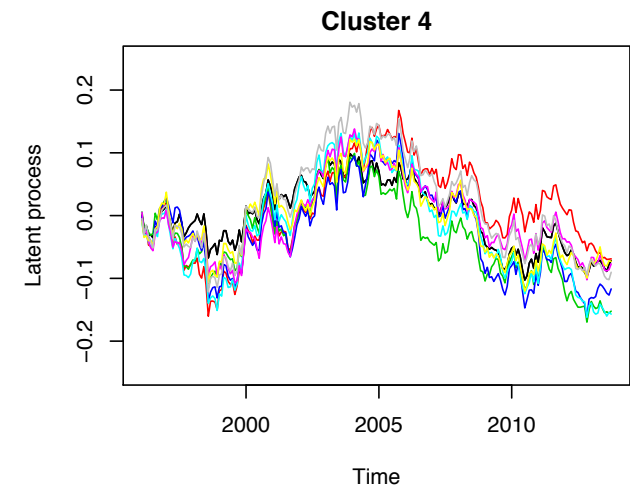
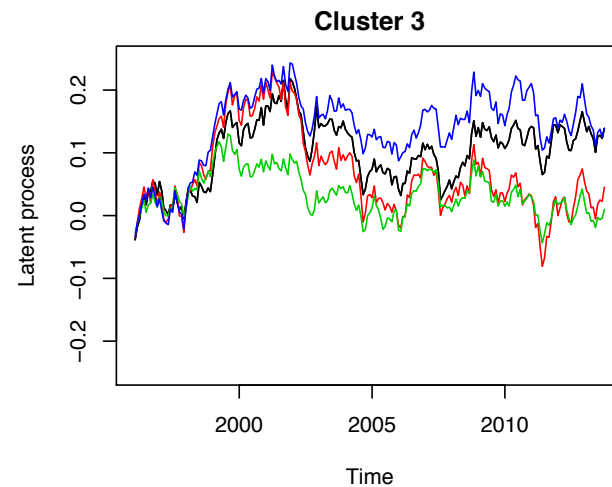
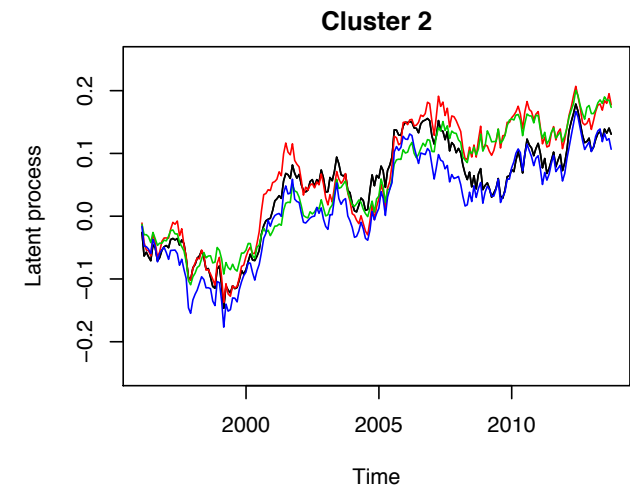
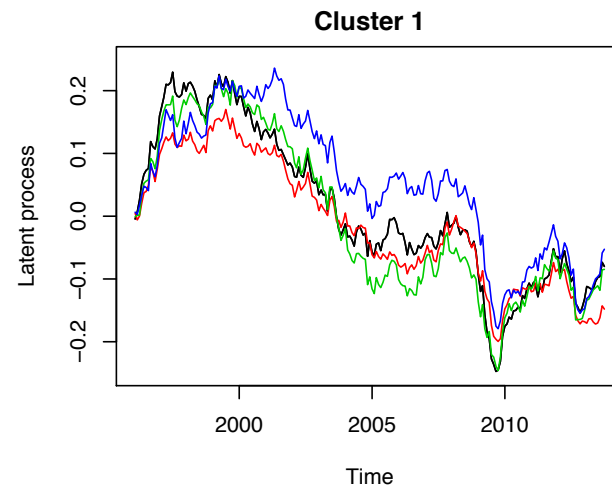


# Relate Time Series via Clustering

## Solution:

Discover groups of tracts with correlated dynamics

*Leverage observations jointly within group*



# State Space Model

tract  $i$

Hidden: global trend + seasonality

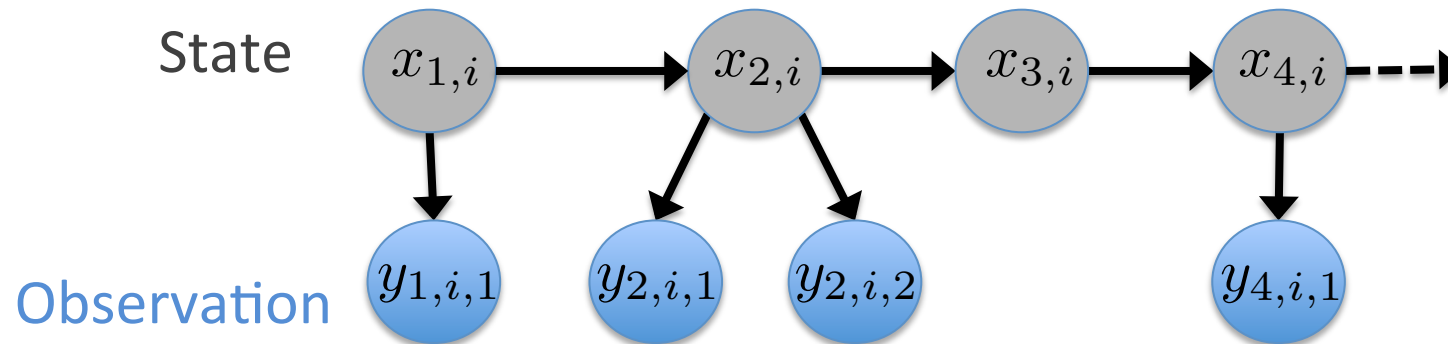
Latent price dynamics  $x_{t,i} = a_i x_{t-1,i} + \epsilon_{t,i} \quad \epsilon_{t,i} \sim \mathcal{N}(0, \sigma_i^2)$

Observed log(price)  $y_{t,i,l} = x_{t,i} + \sum_{h=1}^H \beta_{i,h} U_{l,r} + v_{t,i,l} \quad v_{t,i,l} \sim \mathcal{N}(0, R_i)$

$l^{\text{th}}$  sales

covariate effects

Discrete-time linear Gaussian state space model for census tract  $i$

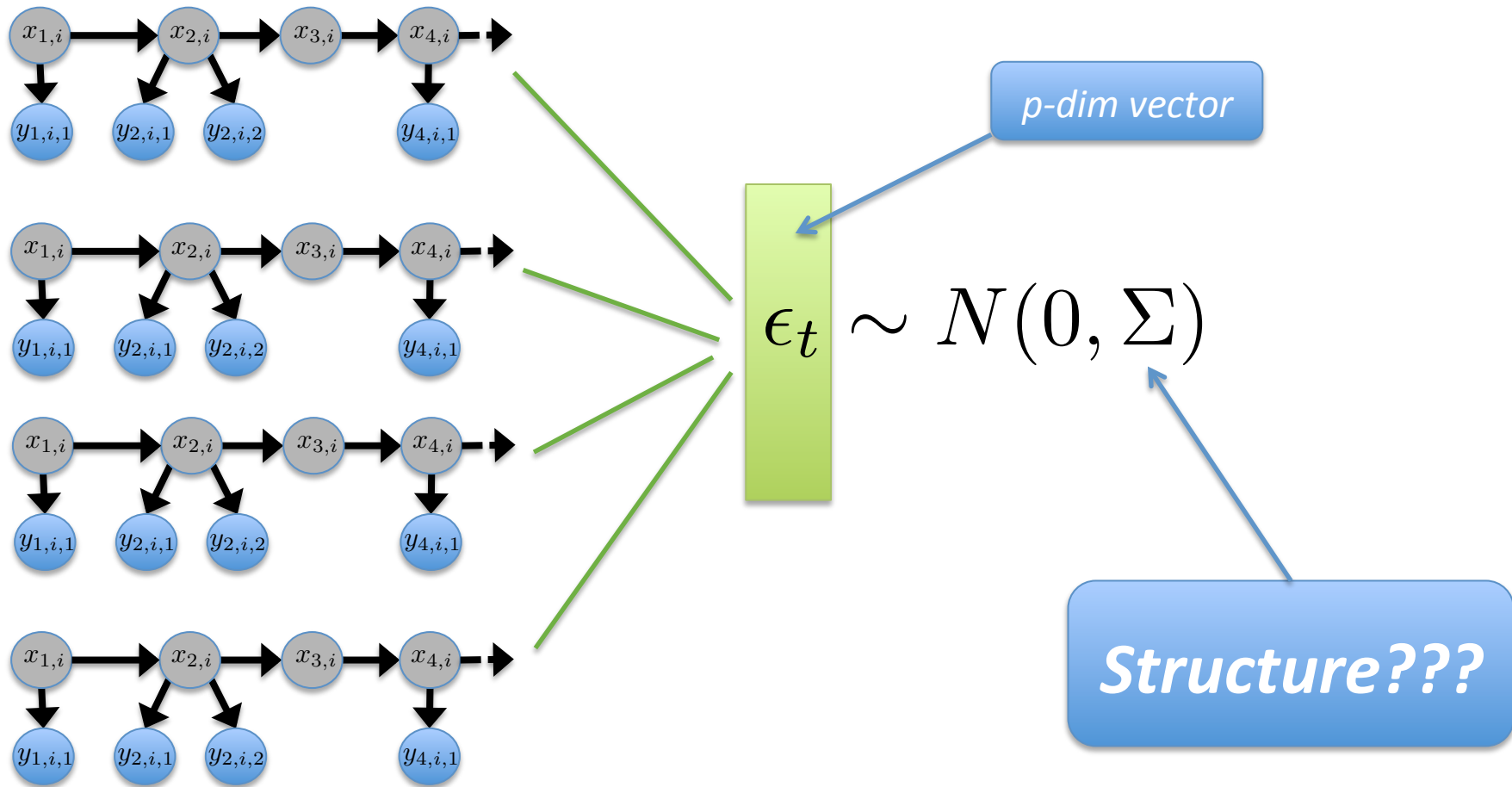


Ren, Fox, Bruce, arXiv 2015.



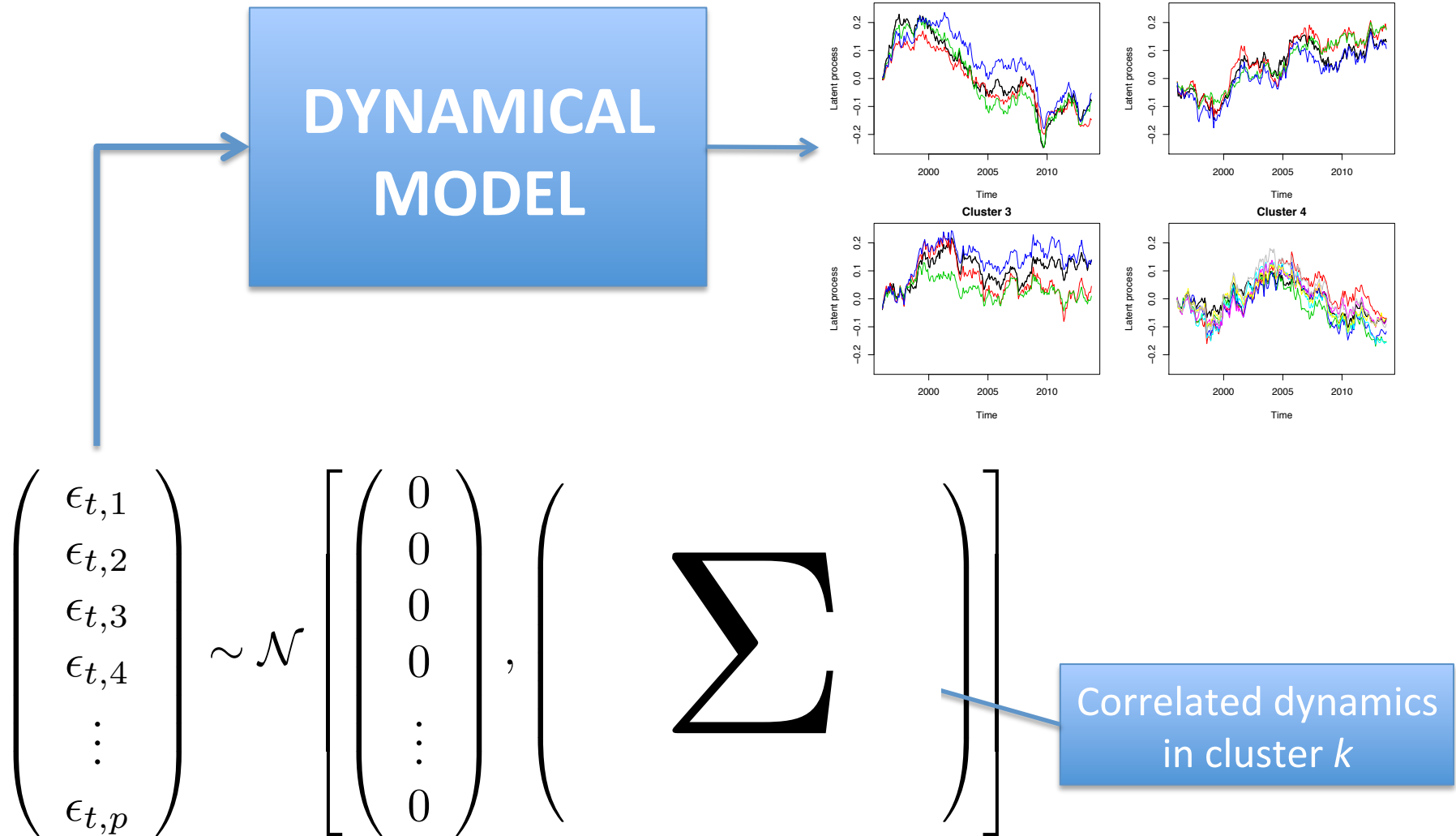
# Multiple Census Tract Model

Latent price dynamics:  $x_{t,i} = a_i x_{t-1,i} + \epsilon_{t,i}$   $\epsilon_{t,i} \sim \mathcal{N}(0, \sigma_i^2)$



Ren, Fox, Bruce, arXiv 2015.

# Cluster and Correlate Multiple Time Series



# Cluster and Correlate Multiple Time Series

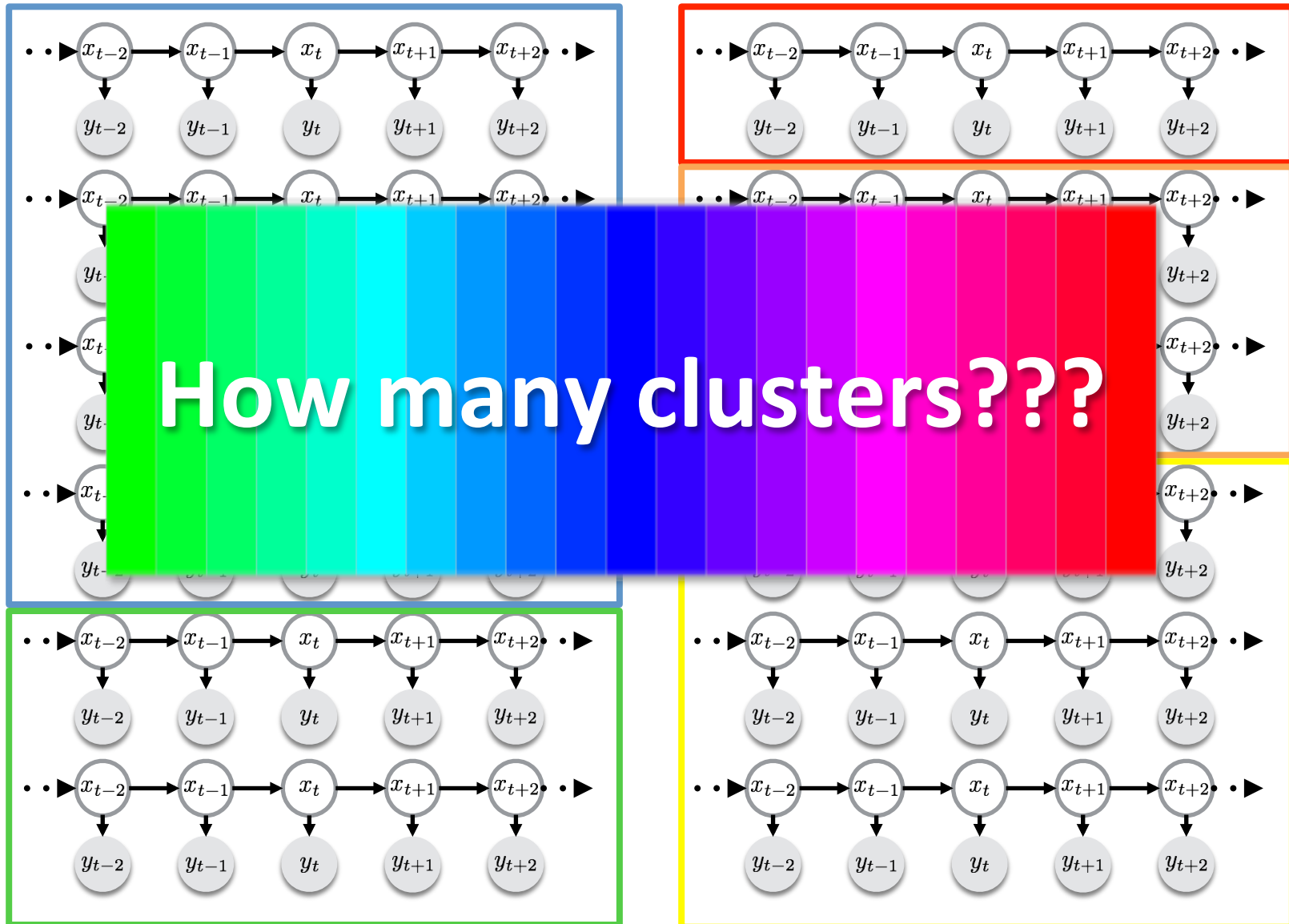
- **Challenge:**  
Unknown cluster structure = unknown # of blocks & size of each
- **Solution:** *Latent factor model* with Bayesian nonparametric prior on latent factor processes

$$\begin{pmatrix} \epsilon_{t,1} \\ \epsilon_{t,2} \\ \epsilon_{t,3} \\ \epsilon_{t,4} \\ \vdots \\ \epsilon_{t,p} \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma_1 & & & & 0 \\ & \Sigma_2 & & & \\ & & \ddots & & \\ 0 & & & & \Sigma_K \end{pmatrix} \right]$$

Correlated dynamics in cluster  $k$

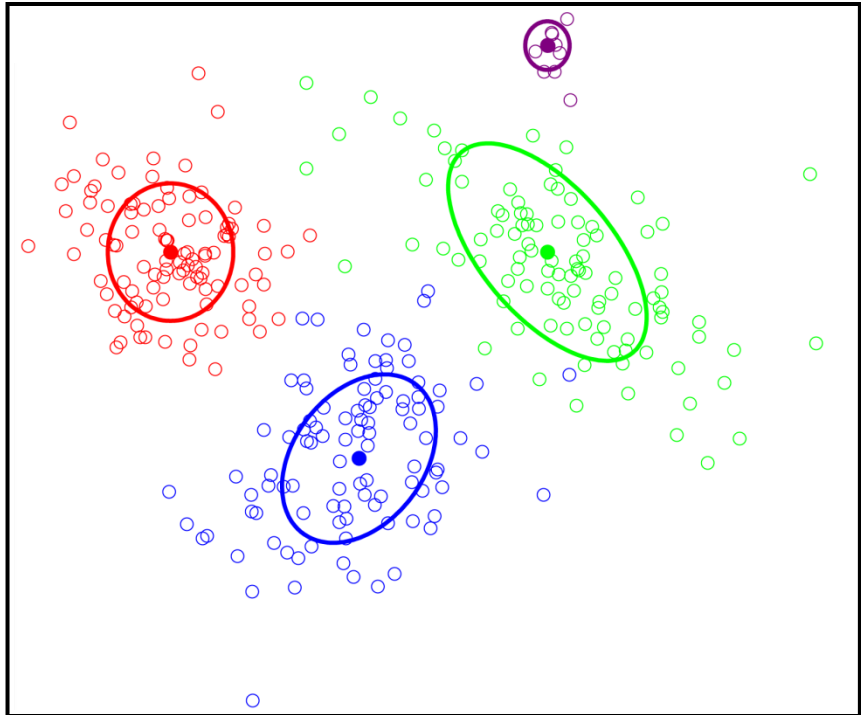


# Clustering Time Series



# Bayesian Nonparametric Clustering

- Bayesian nonparametric approach:
  - Allows infinite # clusters
  - Uses sparse subset
  - Model *complexity adapts* to observations



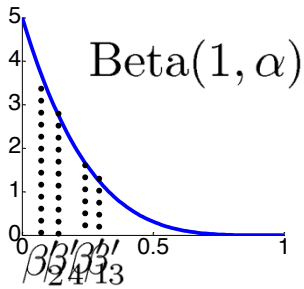
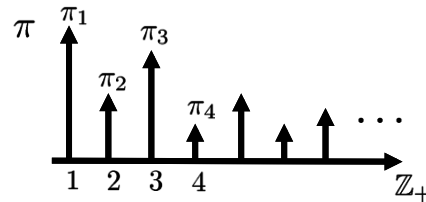
Mixture of Gaussians

$$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6 \quad \theta_7 \quad \dots$$

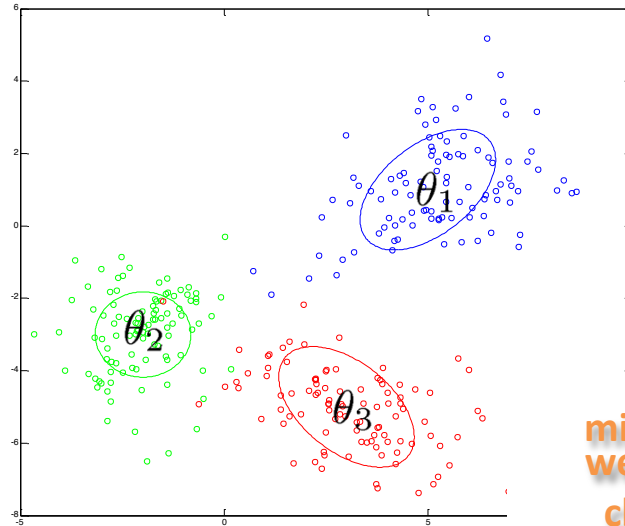
# Dirichlet Process Mixture Model

**REVIEW**

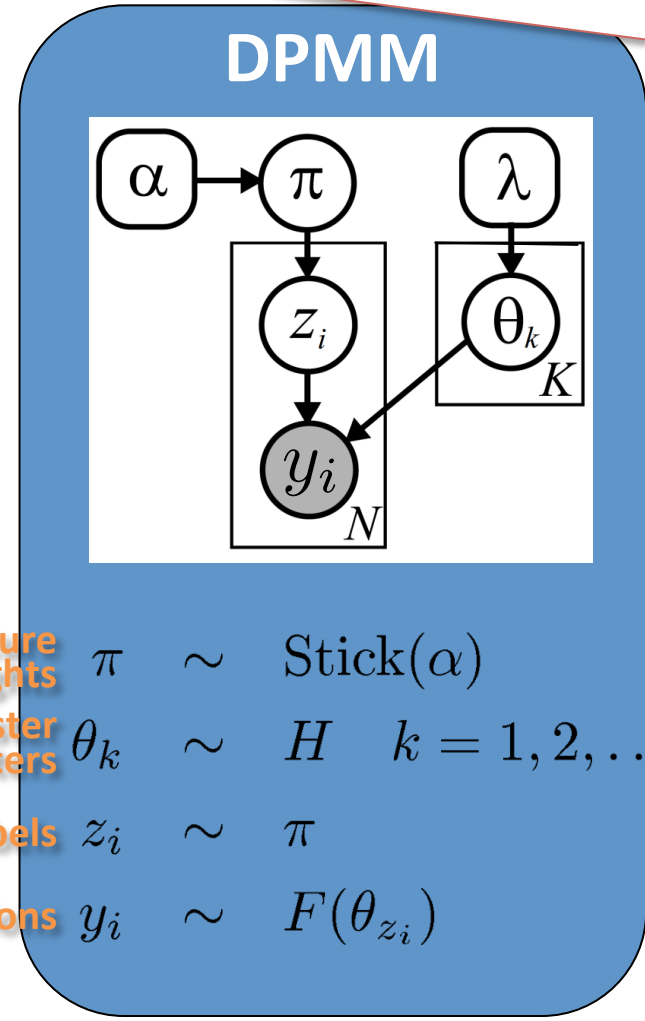
Stick-breaking construction for:  
DP( $\alpha H$ )



Stick of unit probability mass



- mixture weights**  $\pi \sim \text{Stick}(\alpha)$
- cluster parameters**  $\theta_k \sim H \quad k = 1, 2, \dots$
- cluster labels**  $z_i \sim \pi$
- observations**  $y_i \sim F(\theta_{z_i})$



# Chinese Restaurant Process (CRP)

**REVIEW**

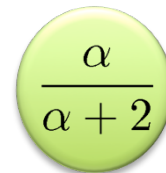
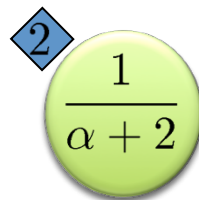
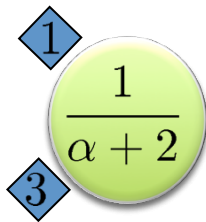
- Distribution on induced partitions described via the CRP
- Visualize clustering as a sequential process of customers sitting at tables in an (infinitely large) restaurant:

*customers*  $\longleftrightarrow$  *observed data to be clustered*

*tables*  $\longleftrightarrow$  *distinct clusters*

- The first customer sits at a table. Subsequent customers randomly select a table according to:

$$p(z_{N+1} = z \mid z_1, \dots, z_N, \alpha) = \frac{1}{\alpha + N} \left( \sum_{k=1}^K N_k \delta(z, k) + \alpha \delta(z, \bar{k}) \right)$$



• • •

Number of current assignments to parameter  $k$



# Cluster by Latent Factor Process

Latent price dynamics  $x_{t,i} = a_i x_{t-1,i} + \epsilon_{t,i} \quad \epsilon_{t,i} \sim \mathcal{N}(0, \sigma_i^2)$

Observed log(price)  $y_{t,i,l} = x_{t,i} + \sum_{h=1}^H \beta_{i,h} U_{l,r} + v_{t,i,l} \quad v_{t,i,l} \sim \mathcal{N}(0, R_i)$

Labels: tract  $i$ ,  $l^{\text{th}}$  sales, covariate effects

**Recall:** Desired structure attained by assuming that if tract  $i$  is from cluster  $k$ ,

$$\epsilon_{t,i} = \lambda_{ik} \eta_{t,k}^* + \tilde{\epsilon}_{t,i} \quad \tilde{\epsilon}_{t,i} \sim N(0, \sigma_0^2) \quad \eta_{t,k}^* \sim N(0, 1).$$

Labels: factor loadings, latent factor process for cluster  $k$

**Motivates:** Dirichlet process mixture model with

$$\theta_k = \eta_{1:T,k}^* \quad k = 1, 2, \dots$$

Gaussian i.i.d. version:  
[Palla et al., NIPS 2012]

# Alternative Clustering

Latent price dynamics  $x_{t,i} = a_i x_{t-1,i} + \epsilon_{t,i} \quad \epsilon_{t,i} \sim \mathcal{N}(0, \sigma_i^2)$

Observed log(price)  $y_{t,i,l} = x_{t,i} + \sum_{h=1}^H \beta_{i,h} U_{l,r} + v_{t,i,l} \quad v_{t,i,l} \sim \mathcal{N}(0, R_i)$

Labels: tract  $i$ ,  $l^{\text{th}}$  sales, covariate effects

Alternative: Dirichlet process mixture model with

$$\theta_k = \{x_{1:T,k}^*, \beta_{k,h}^*\} \quad k = 1, 2, \dots$$

Cluster-specific latent trend

Cluster-specific covariate model

[Niето-Barajas and Contreras-Cristán, 2014]

Assumes all census tracts in cluster have *same* latent value rather than just *correlated* latent value

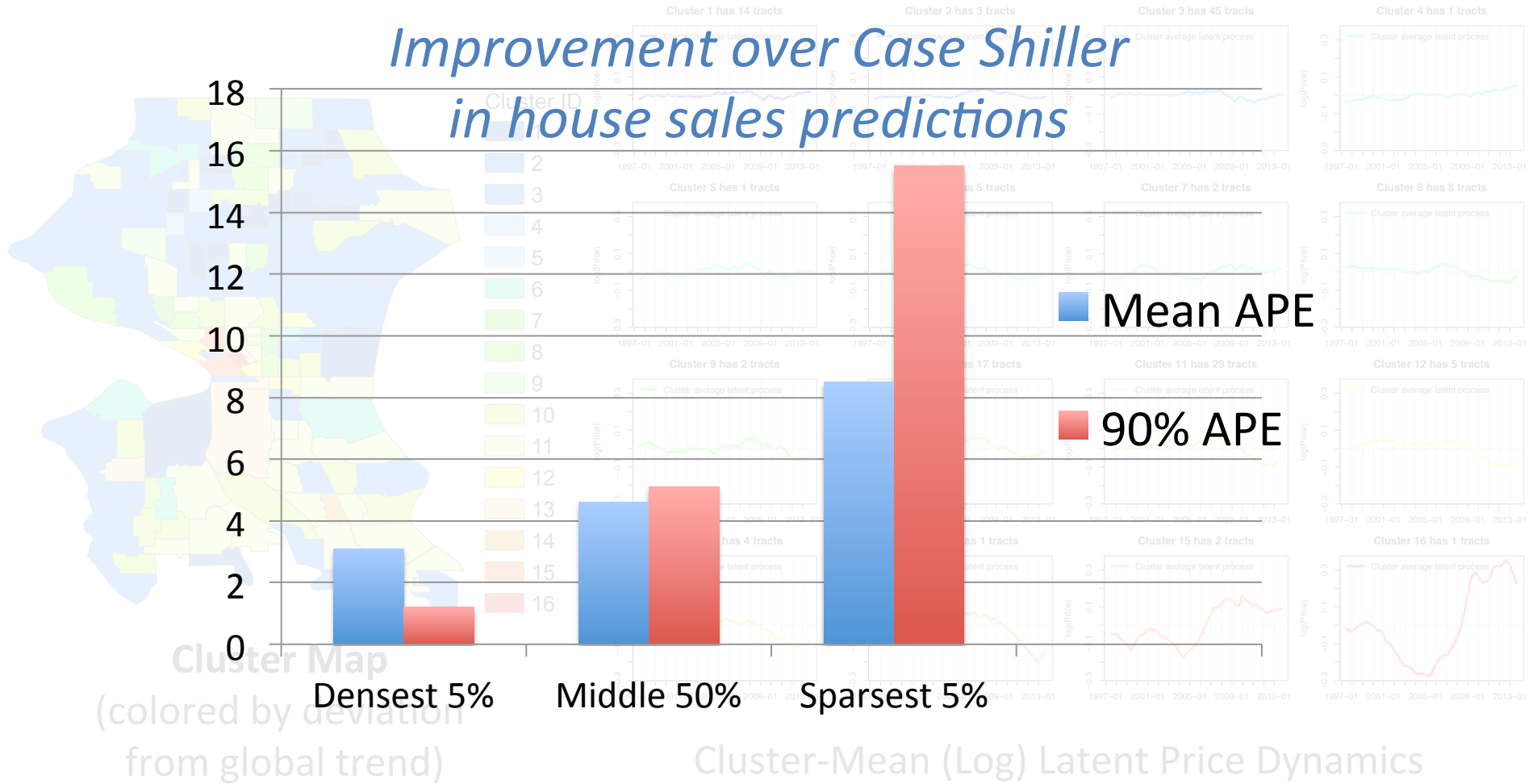
(also cluster parameter  $x_{T+1}$  depends on  $x_T$ , whereas  $\epsilon_{T+1}$  ind. of  $\epsilon_T$ )

# Housing Data Analysis

- Seattle City
  - 140 census tracts
  - 125k transactions during 17 years
- Computational details:
  - Parallel (collapsed) Dirichlet process MCMC sampler [Williamson et al., ICML 2013]
  - 10x speedup with 10 processors

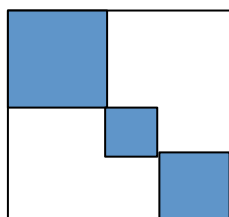
# Seattle City Analysis (17 years)

*Improvement over Case Shiller  
in house sales predictions*



# Clusters of Time Series Summary

$\Sigma$  *sparse*

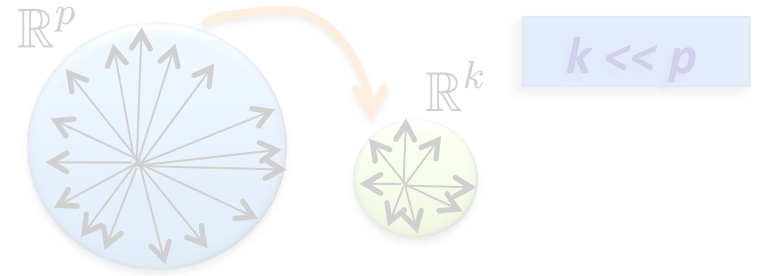
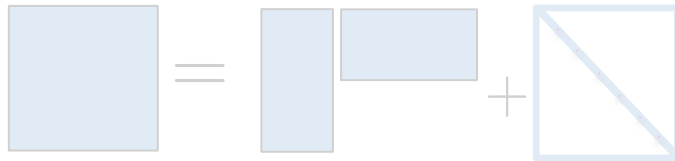


*Independent groups of nodes*

- 
- **Goal:** Cluster time series to share information
    - Individually not informative enough
    - Full joint model statistically and computationally infeasible
  - Cluster structure
    - Cluster on latent state process → clusters of *identical latent trends*
    - Assume latent factor model for AR innovations + cluster latent factor process → clusters of *correlated time series*
  - Bayesian nonparametric clustering
    - Dirichlet process prior allows unknown number of clusters

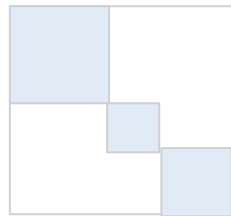
# Methods for Scaling to High Dimensions

Low Rank  $\Sigma = \Lambda\Lambda' + \Sigma_0$



Low-dimensional embedding

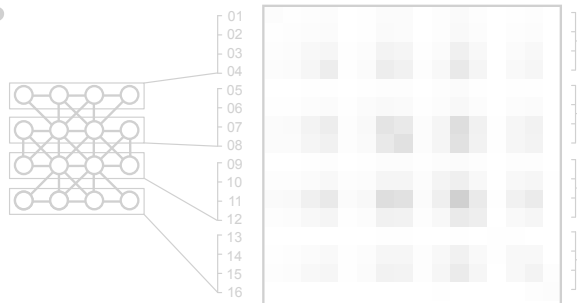
$\Sigma$  sparse



Independent groups of nodes

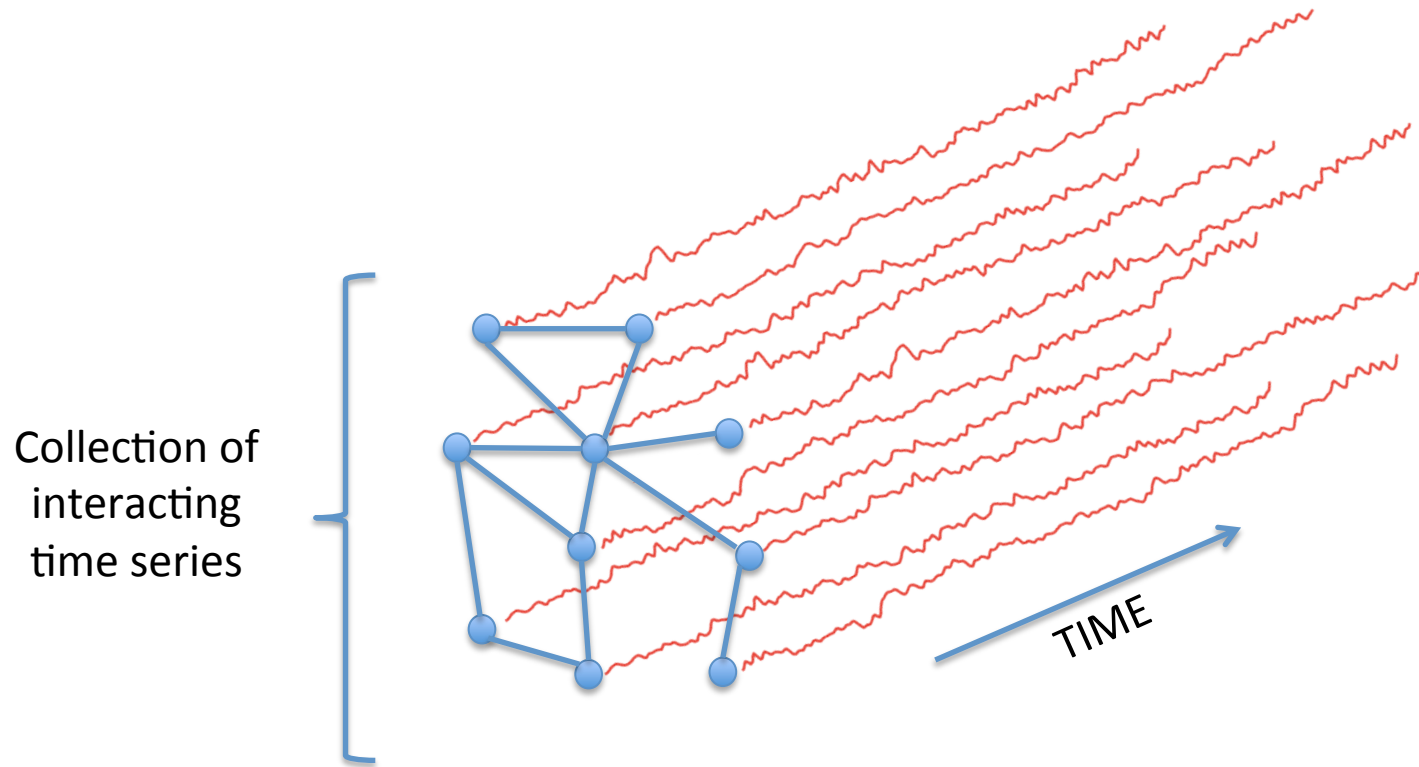
$\Sigma^{-1}$  sparse

Gaussian Graphical Model



Zeros = no edge in graph,  
Cond. ind. between nodes

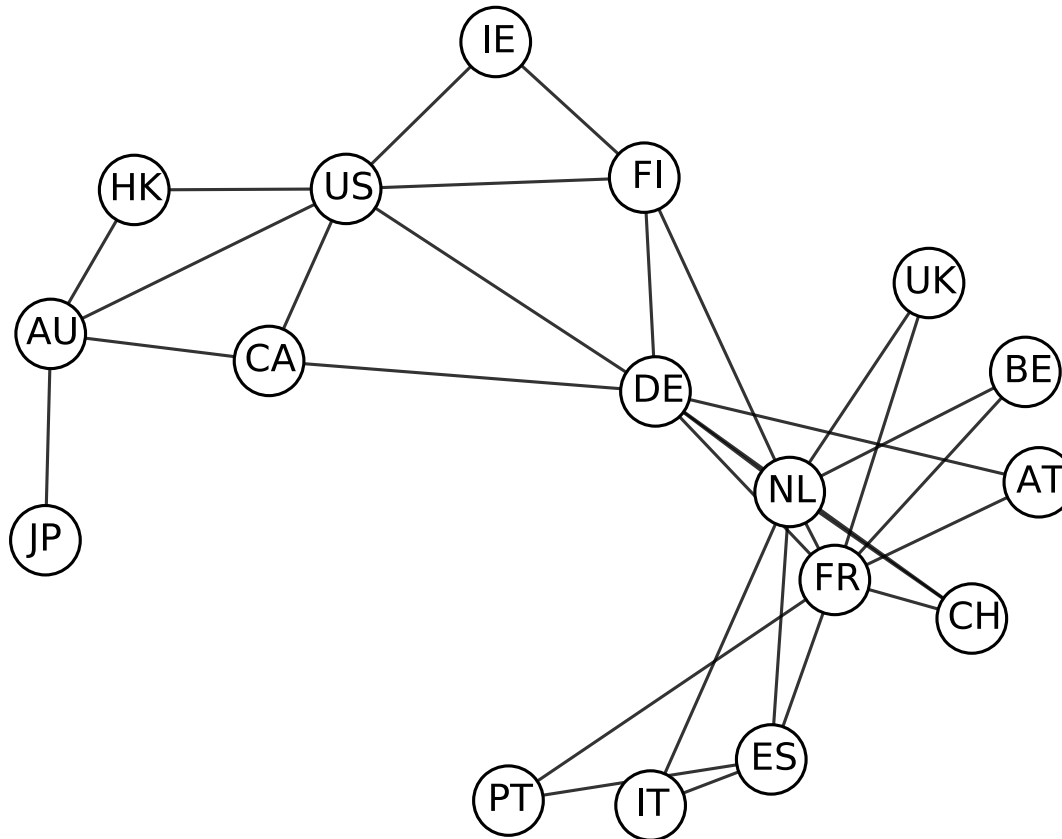
# Graphs of Time Series



# Conditional Independencies

**Collection of global stock indices**

**Data:** time series of daily returns



**Countries:**

Australia (AU)

Austria (AT)

Belgium (BE)

Canada (CA)

Finland (FI)

France (FR)

Germany (DE)

Hong Kong (HK)

Ireland (IE)

Italy (IT)

Japan (JP)

Netherlands (NL)

Portugal (PT)

Spain (ES)

Switzerland (CH)

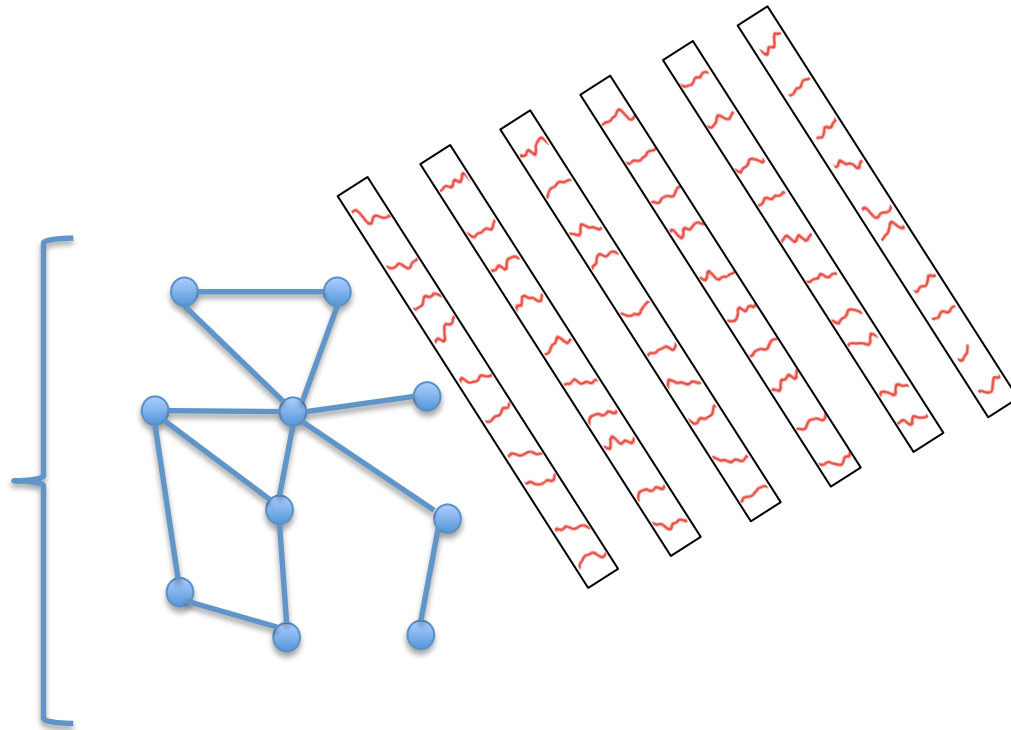
United Kingdom (UK)

United States (US)



# Graphs of i.i.d. Data

Collection of  
interacting  
random  
variables

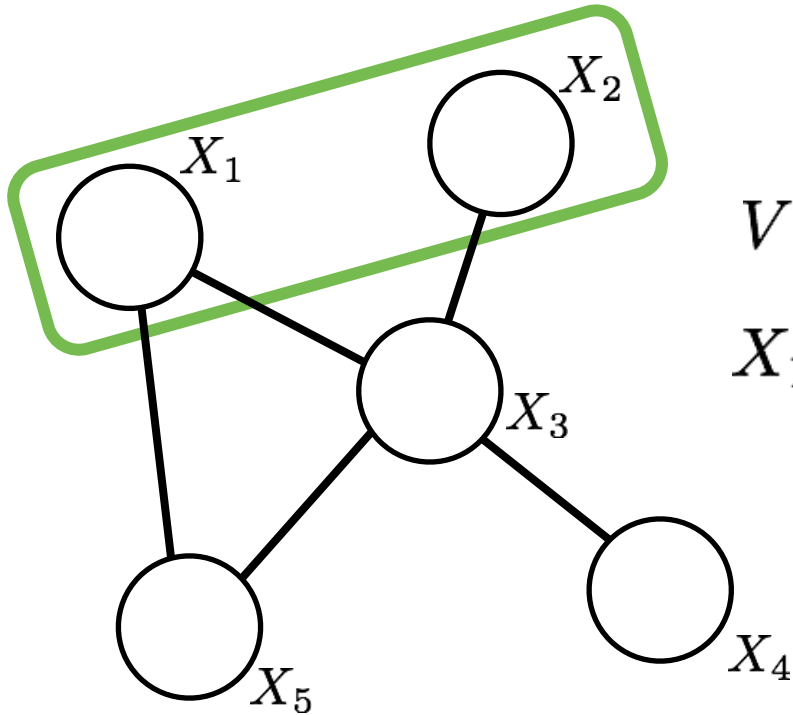


# Graphical Models for Random Variables

- Graph  $G=(V,E)$  encodes conditional independence statements

nodes edges

no edge  $(i,j) \Rightarrow X_i, X_j$  cond. ind. given rest



$$V = \{1, 2, 3, 4, 5\}$$

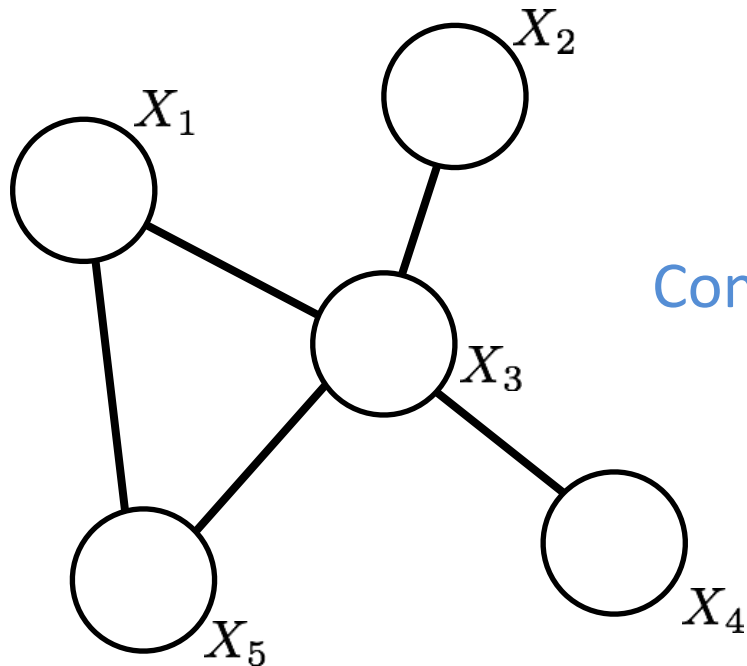
$$X_1 \perp\!\!\!\perp X_2 | (X_3, X_5, X_4)$$

# Gaussian Graphical Models

- Assume Gaussian random vector  $X \sim \mathcal{N}(0, \Sigma)$

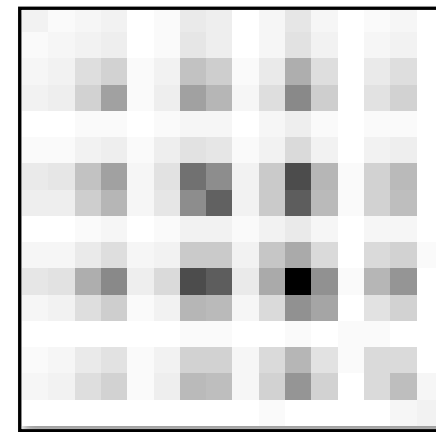
no edge  $(i,j) \Rightarrow X_i, X_j$  cond. ind. given rest

$$\begin{array}{c} \Updownarrow \\ (\Sigma^{-1})_{ij} = 0 \end{array}$$



Cond. ind. encoded in *precision* matrix

$$\Sigma^{-1} =$$



# Information Form Gaussian

- Motivations for considering “information form” of multivariate normal
  - Easier to read off conditional densities
  - Has log-linear form in terms of “information parameters”

$$x \sim N(\mu, \Sigma) \rightarrow \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\begin{array}{|c|} \hline \Omega = \Sigma^{-1} \\ \eta = \Sigma^{-1} \mu \\ \hline \end{array}$$

$$x \sim N^{-1}(\eta, \Omega) \rightarrow \propto e^{\eta^T x - \frac{1}{2} x^T \Omega x}$$

# Info. Gaussian Conditional Densities

- Assume a model with

$$x \sim N^{-1}(\eta, \Omega)$$

and divide the dimensions into two sets  $A, \bar{A}$

- Then,

$$\begin{bmatrix} x_A \\ x_{\bar{A}} \end{bmatrix} \sim N^{-1} \left( \begin{bmatrix} \eta_A \\ \eta_{\bar{A}} \end{bmatrix}, \begin{bmatrix} \Omega_{AA} & \Omega_{A\bar{A}} \\ \Omega_{\bar{A}A} & \Omega_{\bar{A}\bar{A}} \end{bmatrix} \right)$$


$$p(x_A | x_{\bar{A}}) = N^{-1}(\eta_A - \Omega_{A\bar{A}}x_{\bar{A}}, \Omega_{AA})$$

# Info. Gaussian Conditional Densities

- Let  $A = \{s, t\}$  and  $\bar{A}$  everything else

$$p(x_A | x_{\bar{A}}) = N^{-1}(\eta_A - \Omega_{A\bar{A}}x_{\bar{A}}, \Omega_{AA})$$

$\begin{bmatrix} \Omega_{ss} & \Omega_{st} \\ \Omega_{ts} & \Omega_{tt} \end{bmatrix}$   
inverse cov. of  
 $p(x_s, x_t | x_{\setminus st})$



- What if  $\Omega_{st} = 0$ ?

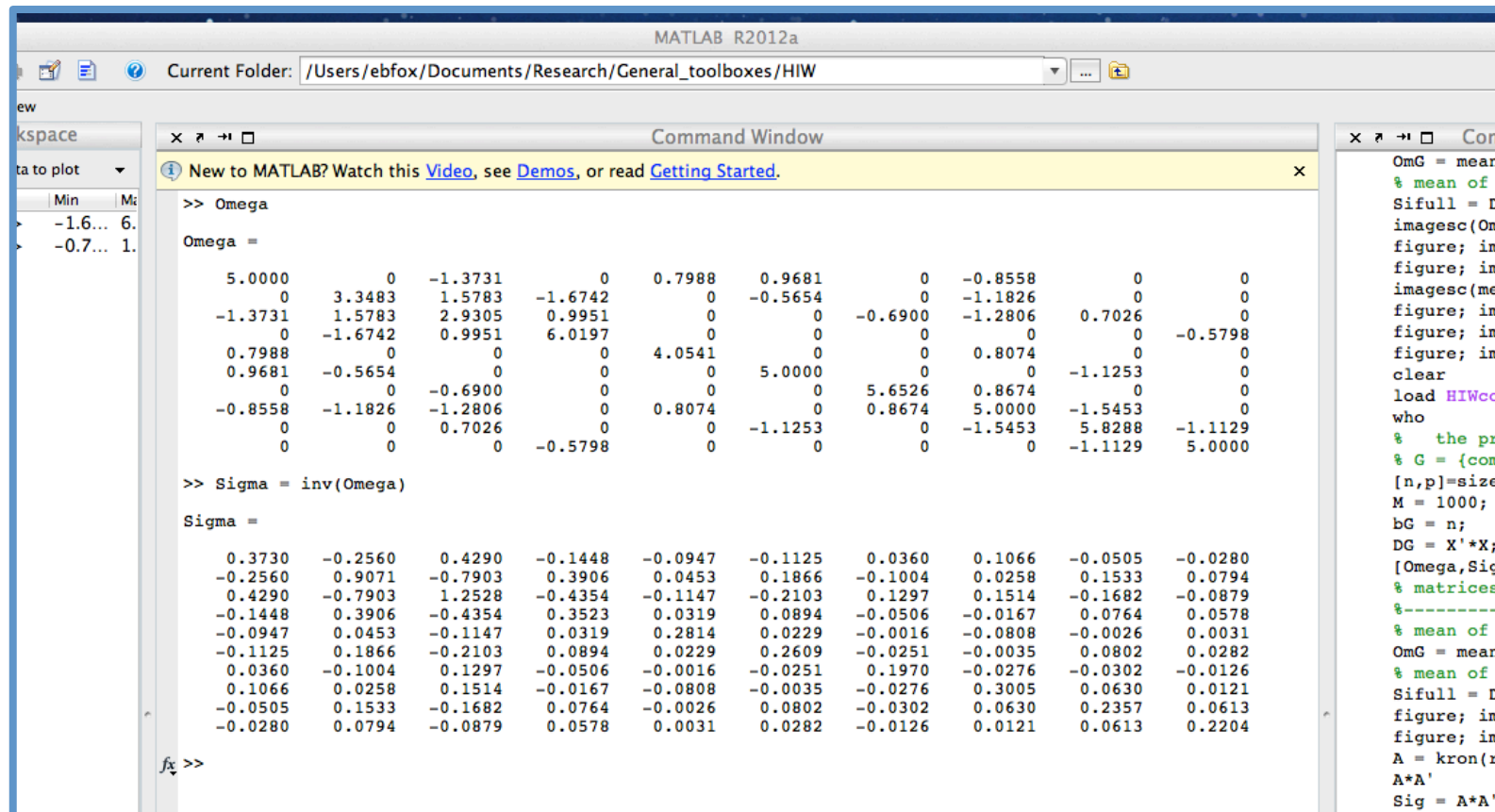
$$\text{COV}(x_s, x_t | x_{\setminus st}) = \Omega_{AA}^{-1} = \begin{bmatrix} \Omega_{ss}^{-1} & 0 \\ 0 & \Omega_{tt}^{-1} \end{bmatrix}$$

$$\Leftrightarrow x_s \perp\!\!\!\perp x_t | x_{\setminus st}$$

- Precision matrix encodes conditional independencies*

# Sparse Precision vs. Covariance

- For a sparse precision matrix, the covariance need not be



The image shows a MATLAB R2012a Command Window with the following content:

```
MATLAB R2012a
Current Folder: /Users/ebfox/Documents/Research/General_toolboxes/HIW

>> Omega
Omega =
    5.0000    0    -1.3731    0    0.7988    0.9681    0    -0.8558    0    0
    0    3.3483    1.5783   -1.6742    0   -0.5654    0    -1.1826    0    0
   -1.3731    1.5783    2.9305    0.9951    0    0   -0.6900   -1.2806    0.7026    0
    0   -1.6742    0.9951    6.0197    0    0    0    0    0   -0.5798
    0.7988    0    0    0    4.0541    0    0    0.8074    0    0
    0.9681   -0.5654    0    0    0    5.0000    0    0    -1.1253    0
    0    0   -0.6900    0    0    0    5.6526    0.8674    0    0
   -0.8558   -1.1826   -1.2806    0    0.8074    0    0.8674    5.0000   -1.5453    0
    0    0    0.7026    0    0    -1.1253    0   -1.5453    5.8288   -1.1129
    0    0    0   -0.5798    0    0    0    0   -1.1129    5.0000

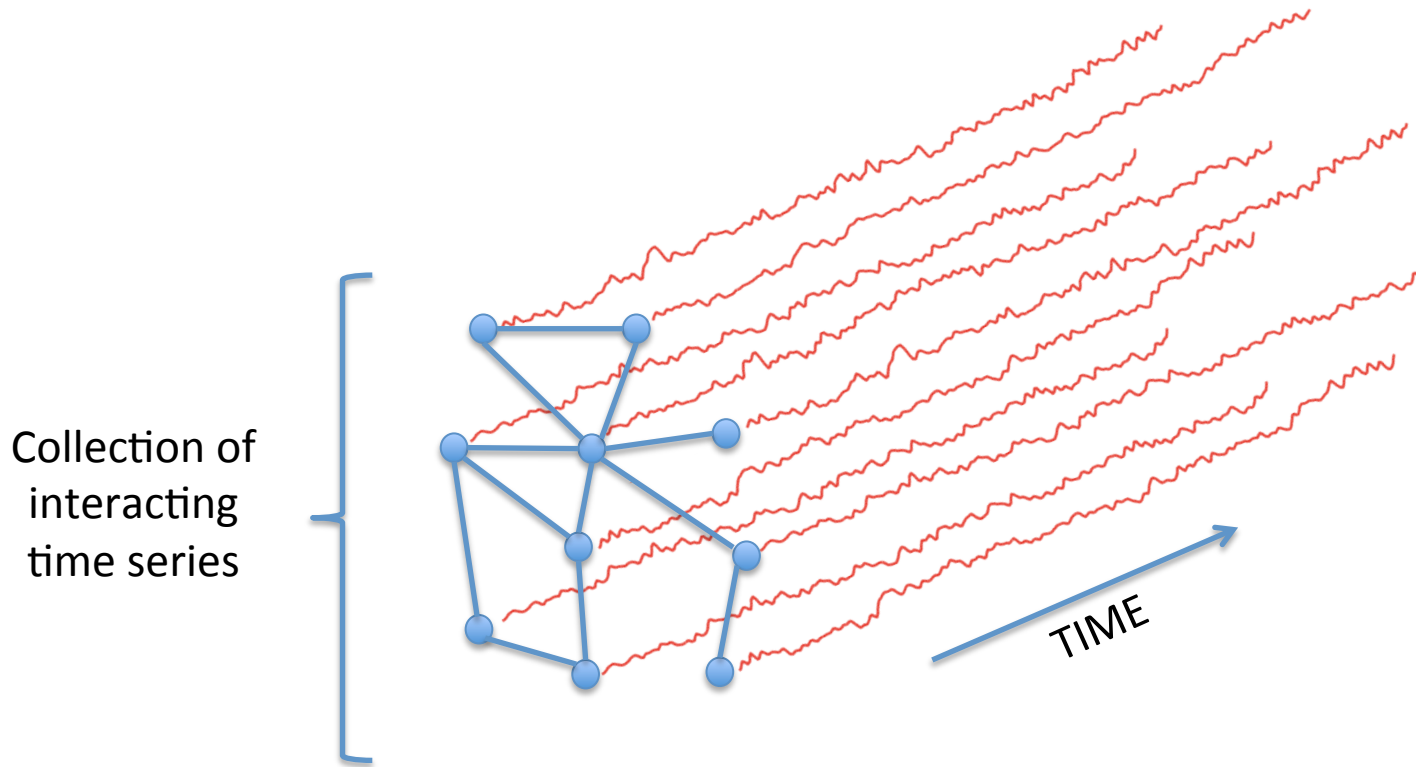
>> Sigma = inv(Omega)
Sigma =
    0.3730   -0.2560    0.4290   -0.1448   -0.0947   -0.1125    0.0360    0.1066   -0.0505   -0.0280
   -0.2560    0.9071   -0.7903    0.3906    0.0453   -0.1866   -0.1004    0.0258    0.1533    0.0794
    0.4290   -0.7903    1.2528   -0.4354   -0.1147   -0.2103    0.1297    0.1514   -0.1682   -0.0879
   -0.1448    0.3906   -0.4354    0.3523    0.0319    0.0894   -0.0506   -0.0167    0.0764    0.0578
   -0.0947    0.0453   -0.1147    0.0319    0.2814    0.0229   -0.0016   -0.0808   -0.0026    0.0031
   -0.1125    0.1866   -0.2103    0.0894    0.0229    0.2609   -0.0251   -0.0035    0.0802    0.0282
    0.0360   -0.1004    0.1297   -0.0506   -0.0016   -0.0251    0.1970   -0.0276   -0.0302   -0.0126
    0.1066    0.0258    0.1514   -0.0167   -0.0808   -0.0035   -0.0276    0.3005    0.0630    0.0121
   -0.0505    0.1533   -0.1682    0.0764   -0.0026    0.0802   -0.0302    0.0630    0.2357    0.0613
   -0.0280    0.0794   -0.0879    0.0578    0.0031    0.0282   -0.0126    0.0121    0.0613    0.2204

fx >>
```

On the right side of the Command Window, a portion of a script is visible:

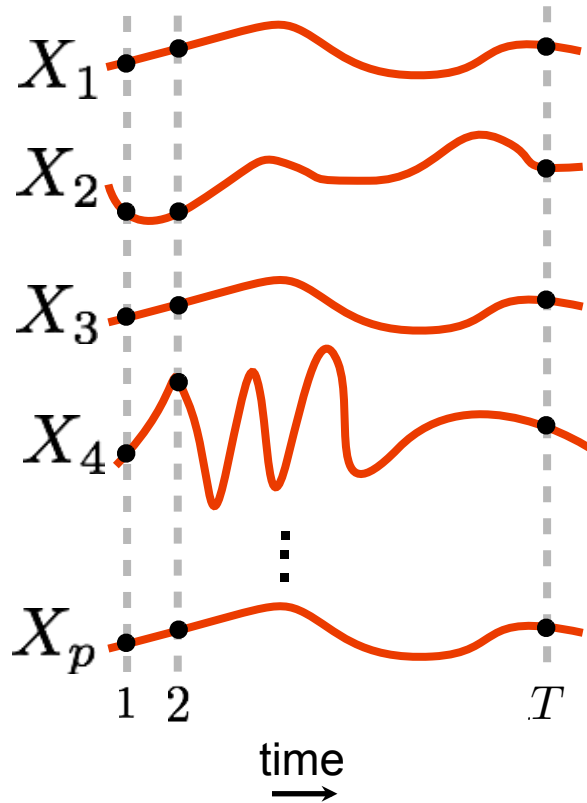
```
OmG = mean
% mean of
Sifull = D
imagesc(Om
figure; im
figure; im
imagesc(me
figure; im
figure; im
figure; im
clear
load HIWco
who
% the pr
% G = {com
[n,p]=size
M = 1000;
bG = n;
DG = X'*X;
[Omega,Sig
% matrices
%-----
% mean of
OmG = mean
% mean of
Sifull = D
figure; im
figure; im
A = kron(r
A*A'
Sig = A*A'
```

# Defining Graphs of Time Series





# Random Variables $\rightarrow$ Stochastic Processes



**Goal:** Represent and infer conditional independence relations between *time series*

Assume *stationarity*:

$$E(X(t)) = \mu$$

$$\text{Cov}(X(t), X(t+h)) = \Gamma(h)$$

For simplicity, zero mean

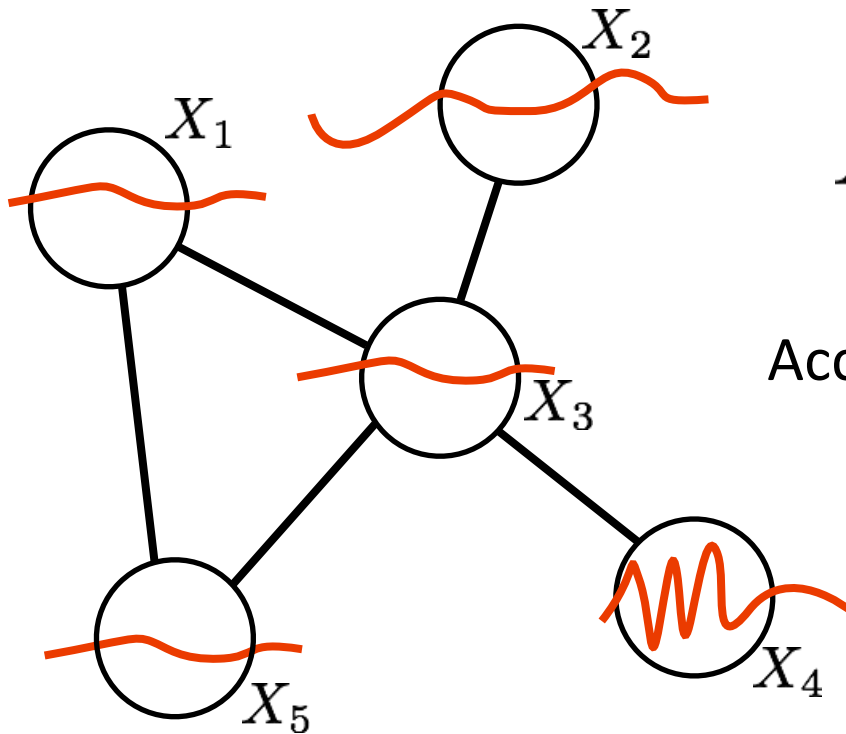
# Graphical Models for Time Series

no edge  $(i,j) \Rightarrow$  *time series*  $X_i, X_j$  cond. ind.  
given *entire histories* of other series

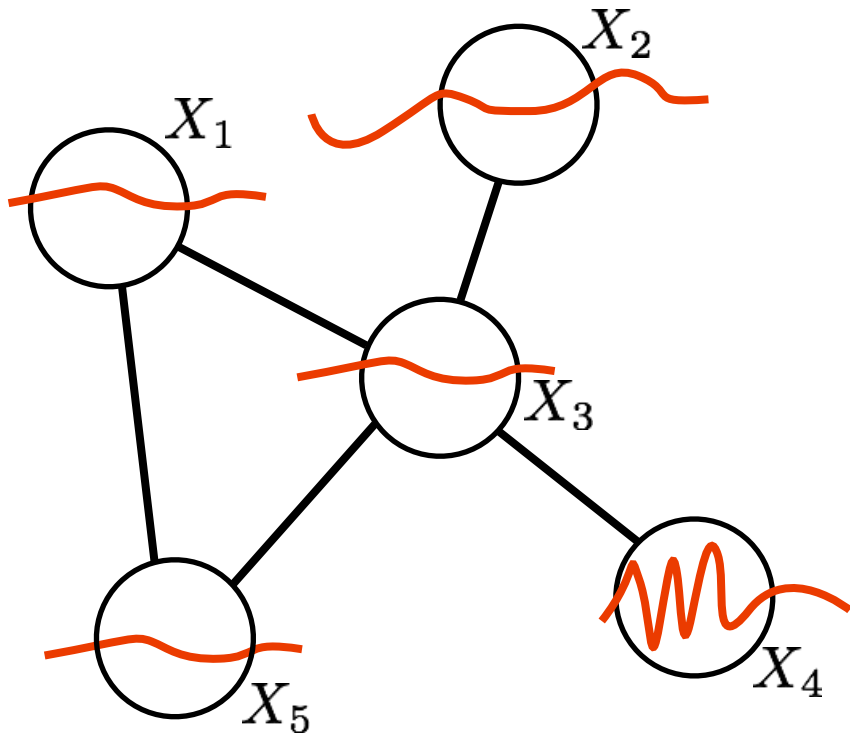


$$X_1(t) \perp\!\!\!\perp X_2(t+h) | X_{V \setminus \{1,2\}}$$

Accounts for interactions at *any lag*



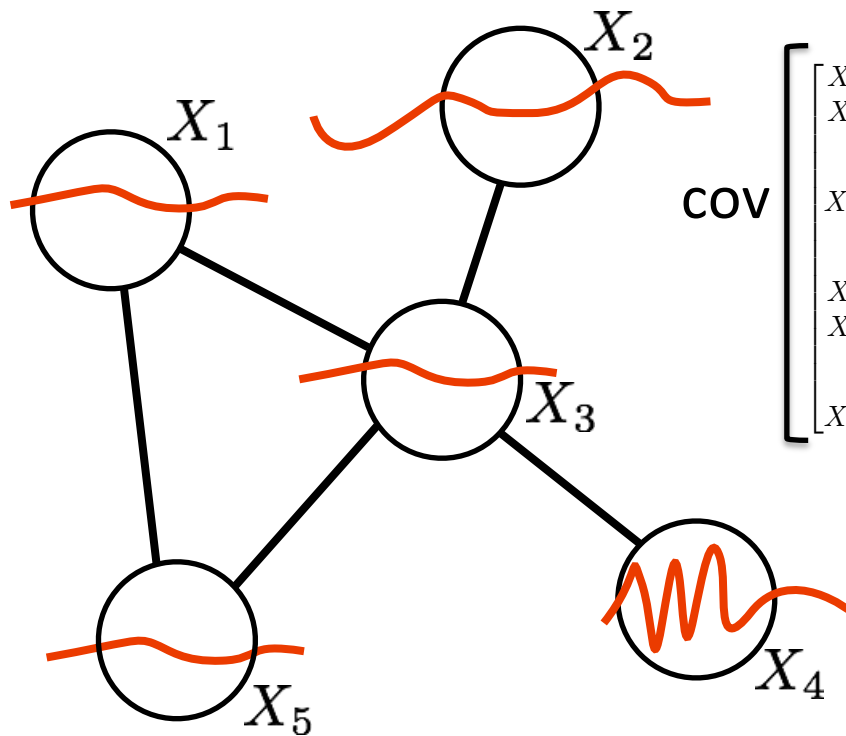
# Naïve Approach to Structure Learning



$$\tilde{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = \begin{bmatrix} X_1(1) \\ X_1(2) \\ \vdots \\ X_1(T) \\ \vdots \\ X_p(1) \\ X_p(2) \\ \vdots \\ X_p(T) \end{bmatrix}$$

# Naïve Approach to Structure Learning

$$\tilde{X} \sim N(0, \Sigma)$$



COV

$$\begin{bmatrix} X_1(1) \\ X_1(2) \\ \vdots \\ X_1(T) \\ \vdots \\ X_p(1) \\ X_p(2) \\ \vdots \\ X_p(T) \end{bmatrix}$$

$$= \begin{bmatrix} \Gamma(0) & \Gamma(1) & \Gamma(2) & \cdots & \Gamma(T-1) \\ \Gamma(1)' & \Gamma(0) & \Gamma(1) & \cdots & \Gamma(T-2) \\ \vdots & & \ddots & & \vdots \\ \Gamma(T-1)' & \Gamma(T-2)' & \cdots & \Gamma(1)' & \Gamma(0) \end{bmatrix}$$

*$Tp \times Tp$  matrix*

*Lagged covariance matrix*

$$\Gamma(h) = \text{Cov}(X(t), X(t+h))$$

# Previous Approaches

- **Songsiri et al. 2011**
  - Assume parametric time series model
  - Determine conditions on parameters leading to conditional independencies
  - Optimize penalized likelihood

**VAR(q) process:** 
$$X(t) = \sum_{i=1}^q A_i X(t-i) + \epsilon(t) \quad \epsilon(t) \sim N(0, \Sigma_\epsilon)$$

Define: 
$$B_k = \Sigma_\epsilon^{-\frac{1}{2}} A_k \quad Y^k = \sum_{\ell=0}^{p-k} B_\ell^T B_{k+\ell}$$

**Main Result:**

$$X_a \perp X_b \mid X_{V \setminus \{a,b\}} \iff Y_{ab}^k = Y_{ba}^k = 0 \quad \forall k$$

**Objective:** Penalized likelihood with a group penalty to enforce group sparsity

**Tool:** Optimize convex relaxation

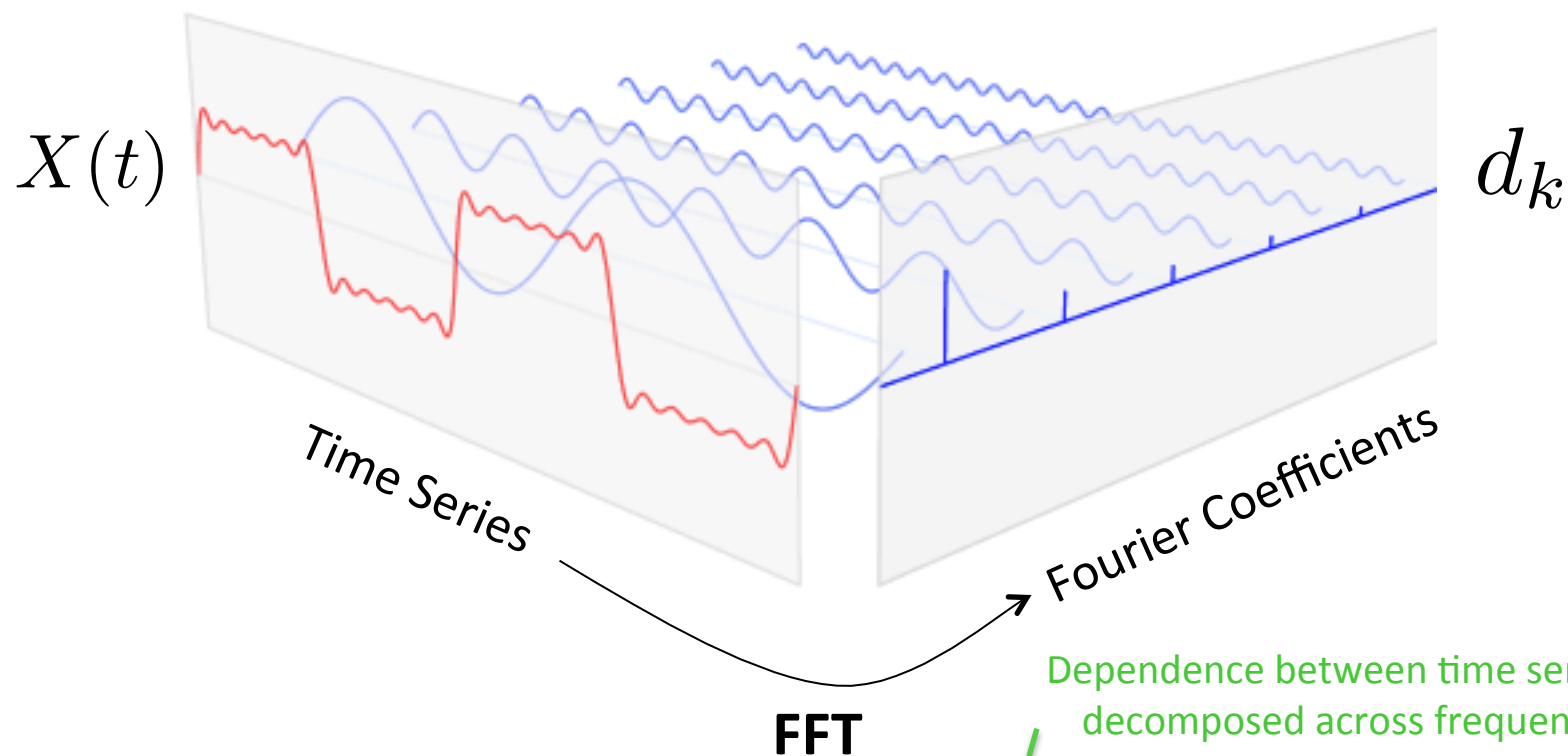
# Previous Approaches

- **Songsiri et al. 2011**
  - Assume parametric time series model
  - Determine conditions on parameters leading to conditional independencies
  - Optimize penalized likelihood
- **Dahlhaus 2000, Matsuda 2006, Wolstenholme and Walden 2015, Bach and Jordan 2004, Jung et al. 2014**
  - Transform to frequency domain
  - Determine conditions on spectral parameters leading to conditional independencies
  - Hypothesis test to see if conditions are satisfied

OR

  - Optimize a Whittle-approximated (penalized) likelihood

# Model in the Frequency Domain



Dependence between time series,  
decomposed across frequencies

*Spectral density matrix*

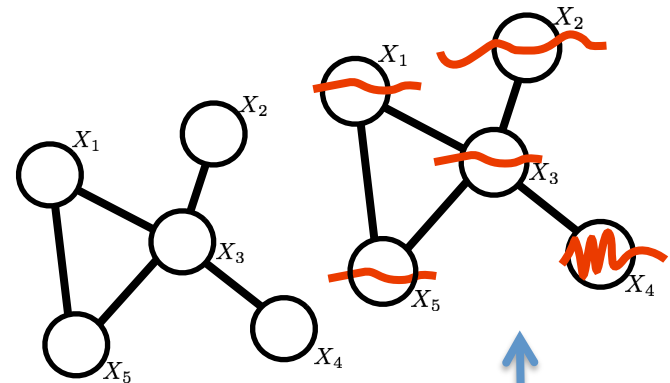
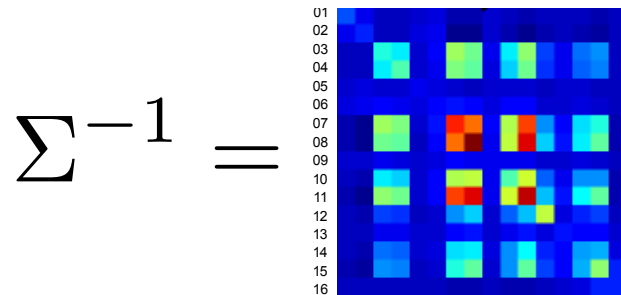
*Lagged covariance matrix*

$$\Gamma(h) = \text{Cov}(X(t), X(t+h)) \longrightarrow S(\lambda) = \sum_{h=-\infty}^{\infty} \Gamma(h) e^{-i\lambda h}$$

What conditions on  $S$  leads to conditional independence?

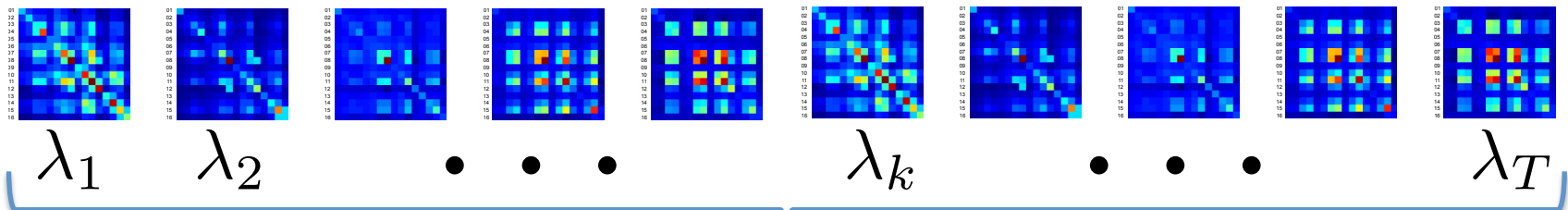
# Encoding Time Series Structure

For **Gaussian** i.i.d. random variables,



(Dahlhaus, 2000) For **Gaussian** stationary time series,

$S^{-1}(\lambda)$  :



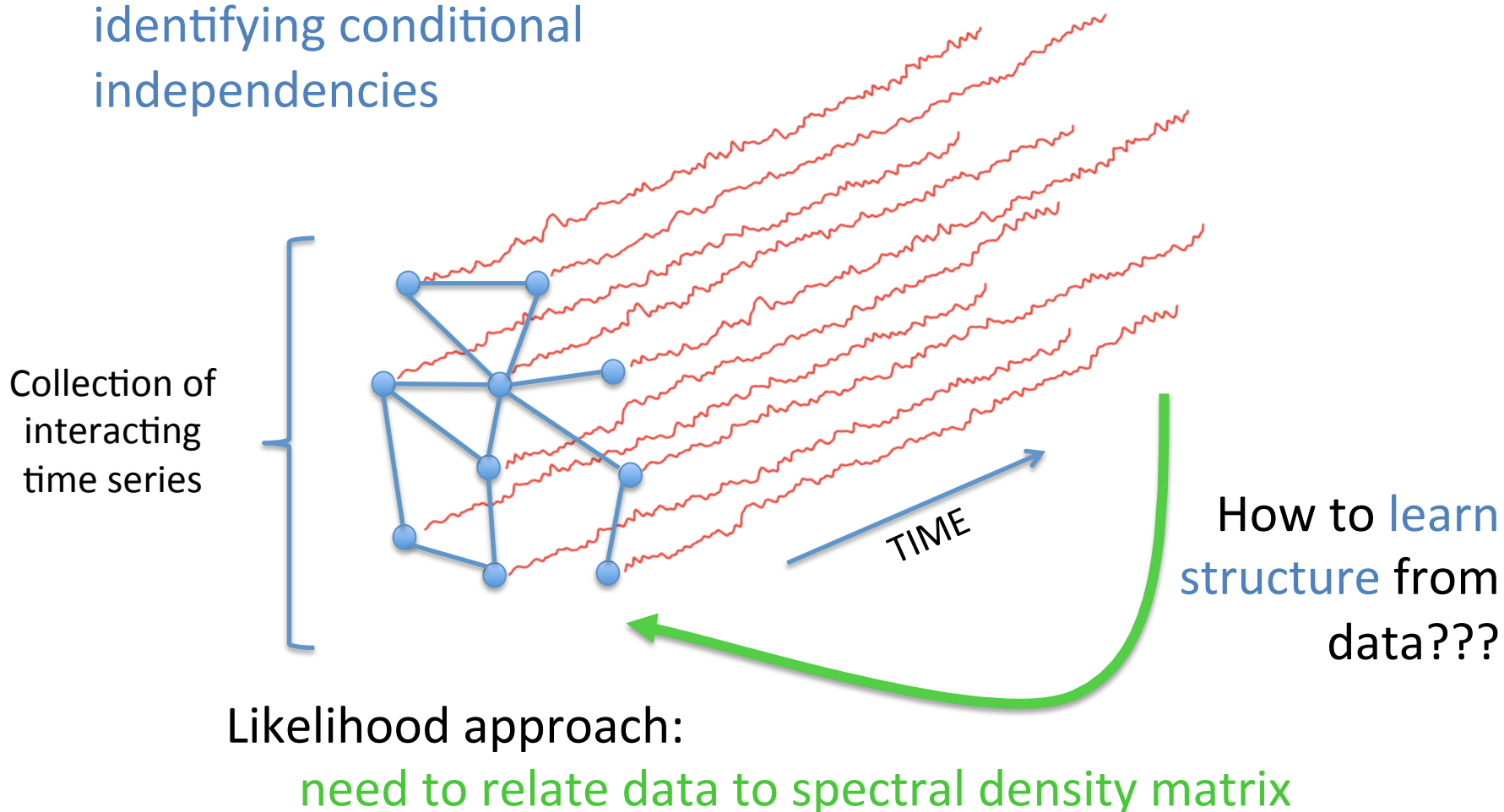
*complex* inverse **spectral density** matrices

$$X_i \perp X_j \mid X_{V \setminus \{i,j\}} \iff S(\lambda)_{ij}^{-1} = 0 \quad \forall \lambda \in [0, 2\pi]$$



# Learning Graphs of Time Series

Dalhaus gives conditions for identifying conditional independencies



# Whittle Approximation (no graph)

$$d_k = \frac{1}{T} \sum_{t=0}^{T-1} X(t) e^{-i\lambda_k t} \quad \lambda_k = \frac{2\pi k}{T}$$

- *Fourier coeff. asymptotically independent*

$$d_k \sim \mathcal{N}_c(0, S_k) \quad k = 0, \dots, T-1$$

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = \tilde{X} \sim N(0, \Sigma)$$

- Instead of likelihood depending on  $Tp \times Tp$  covariance, decomposes over  $p \times p$  spectral density matrices:

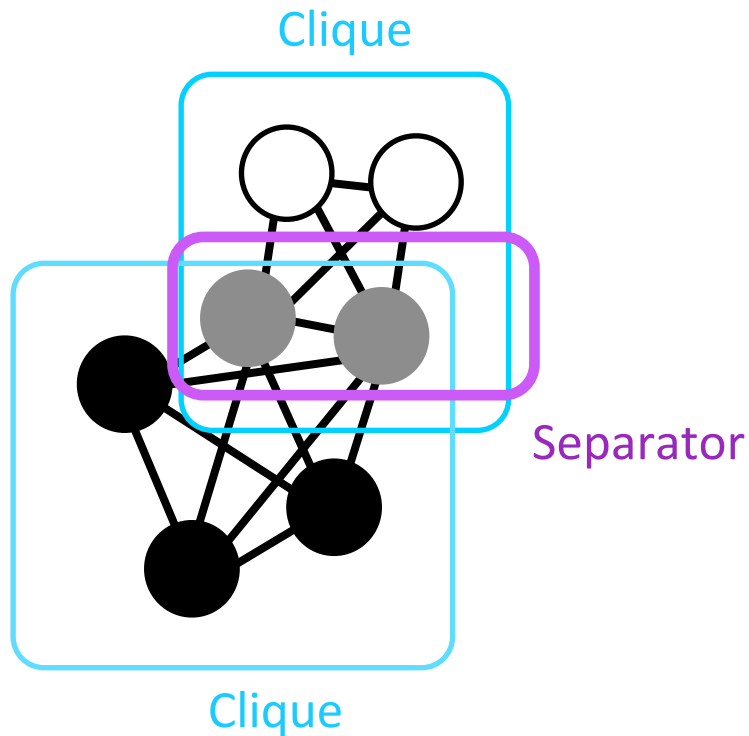
$$p(X_{1:p} | S_{0:T-1}) \approx \prod_{k=0}^{T-1} \frac{1}{\pi^p |S_k|} e^{-d_k^* S_k^{-1} d_k}$$

\* indicates the conjugate transpose

Close to Gaussian likelihood with iid data...

# Decomposable Graphs

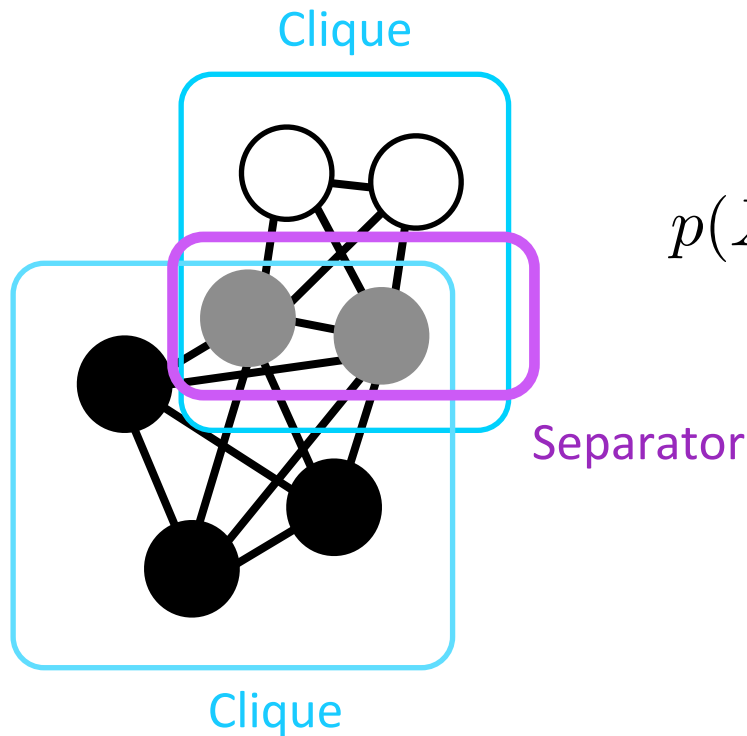
- If graph is **decomposable**, joint distribution decomposes over **cliques**  $\mathcal{C}$  and **separators**  $\mathcal{S}$



$$p(x) = \frac{\prod_{C \in \mathcal{C}} p(x_C)}{\prod_{S \in \mathcal{S}} p(x_S)}$$

# Whittle on Decomposable Graphs

$$p(X_{1:p} | S_{0:T-1}) \approx \prod_{k=0}^{T-1} \frac{1}{\pi^p |S_k|} e^{-d_k^* S_k^{-1} d_k}$$



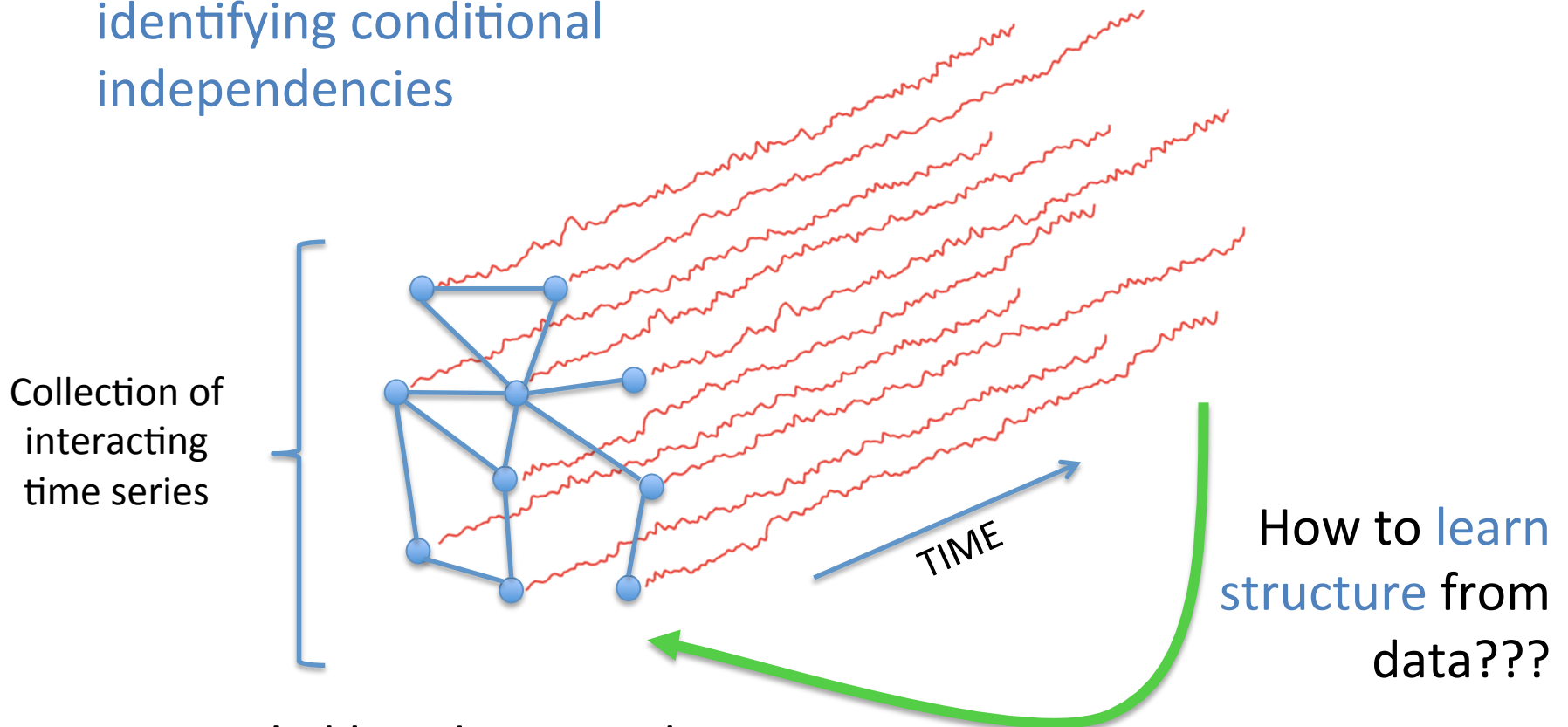
$$p(X_{1:p} | \mathbf{G}, S_{0:(T-1)}) \approx$$

$$\prod_{k=0}^{T-1} \frac{\prod_{C \in \mathcal{C}} \frac{1}{\pi^{p|C|} |S_{kC}|^p} e^{-\text{tr} P_{kC} S_{kC}^{-1}}}{\prod_{S \in \mathcal{S}} \frac{1}{\pi^{p|S|} |S_{kS}|^p} e^{-\text{tr} P_{kS} S_{kS}^{-1}}}$$

Periodogram at frequency  $\lambda_k$   
 $d_k d_k^*$

# Learning Graphs of Time Series

Dalhaus gives conditions for identifying conditional independencies



Likelihood approach:

need to relate data to spectral density matrix

**Bayesian approach: specify priors**

# Bayesian Approach to Structure Learning

## Gaussian Graphical Model

Graph prior  $\mathcal{G} \sim p(\mathcal{G})$

Prior on *graph-constrained* covariance  
(inverses with specified 0's)  $\Sigma \sim p(\Sigma | \mathcal{G})$

Normal obs.  $X \sim \mathcal{N}(0, \Sigma)$

## Time Series Graphs

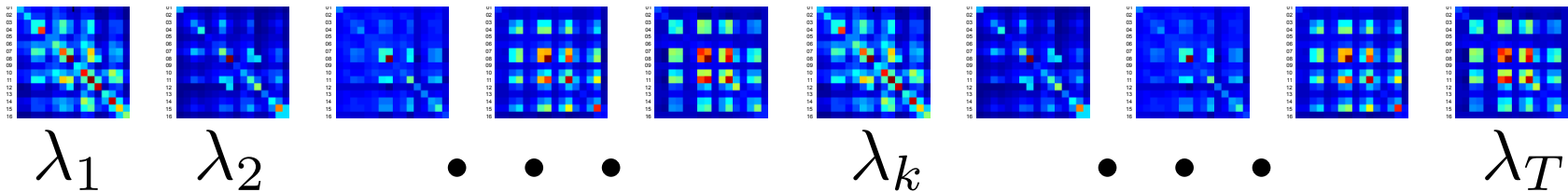
Graph prior  $\mathcal{G} \sim p(\mathcal{G})$

Prior on *graph-constrained* spectral density matrices  
(inverses with specified 0's)  $S_k \sim p(S_k | \mathcal{G})$

Complex normal Fourier coeff.  $d_k \sim \mathcal{N}_c(0, S_k)$

Fourier coefficients

Whittle approx.



*Graph-constrained, complex* spectral density matrices

# Defining a Conjugate Prior

## Gaussian Graphical Model

Graph prior  $\mathcal{G} \sim p(\mathcal{G})$

Prior on *graph-constrained* covariance  
(inverses with specified 0's)  
**HIW distribution**

Normal obs.  $X \sim \mathcal{N}(0, \Sigma)$

## Time Series Graphs

$\mathcal{G} \sim p(\mathcal{G})$

**complex HIW distribution**

$d_k \sim \mathcal{N}_c(0, S_k)$

Graph prior

Prior on *graph-constrained* spectral density matrices  
(inverses with specified 0's)

Complex normal Fourier coeff.

Conjugate prior  $\rightarrow$  marginal likelihood  $p(X | \mathcal{G})$   
graph score

*Use existing graph search algorithms*

# Marginal Likelihood

**REVIEW**

- **Hyper complex inverse Wishart prior** on  $\Sigma$


$$p(\Sigma | \delta, W, G) = h(W, \delta, G) \mathbf{1}_{\Sigma \in M^+(G)} |\Sigma|^{-(\delta+2p)} e^{-\text{tr} W \Sigma^{-1}}$$

- Complex normal observation  $d$

$$p(d | G, \Sigma) = \frac{1}{\pi^p} |\Sigma|^{-1} e^{-\text{tr} P \Sigma^{-1}}$$

- Marginal likelihood

$$\begin{aligned}
 p(d | G) &= \int_{\Sigma} p(d | \Sigma) p(\Sigma) d\Sigma \\
 &= \int_{\Sigma} \frac{1}{\pi^p} h(W, \delta, G) \mathbf{1}_{\Sigma \in M^+(G)} |\Sigma|^{-(\delta+1+2p)} e^{-\text{tr}(W+P)\Sigma^{-1}} d\Sigma \\
 &= \frac{1}{\pi^p} \frac{h(W, \delta, G)}{h(W+P, \delta+1, G)}
 \end{aligned}$$

$\frac{1}{h(W+P, \delta+1, G)} p(\Sigma | \delta+1, W+P, G)$   




# Marginal Likelihood

REVIEW

- Generically:

$$\begin{aligned} \Sigma \mid G &\sim HIW_c(\delta, W, G) \\ d \mid \Sigma &\sim N_c(0, \Sigma) \end{aligned} \quad \longrightarrow \quad p(d \mid G) = \frac{1}{\pi^p} \frac{h(W, \delta, G)}{h(W + P, \delta + 1, G)}$$

- For time series graph:

$$S_k \mid G \sim HIW_c(\delta_k, W_k, G)$$

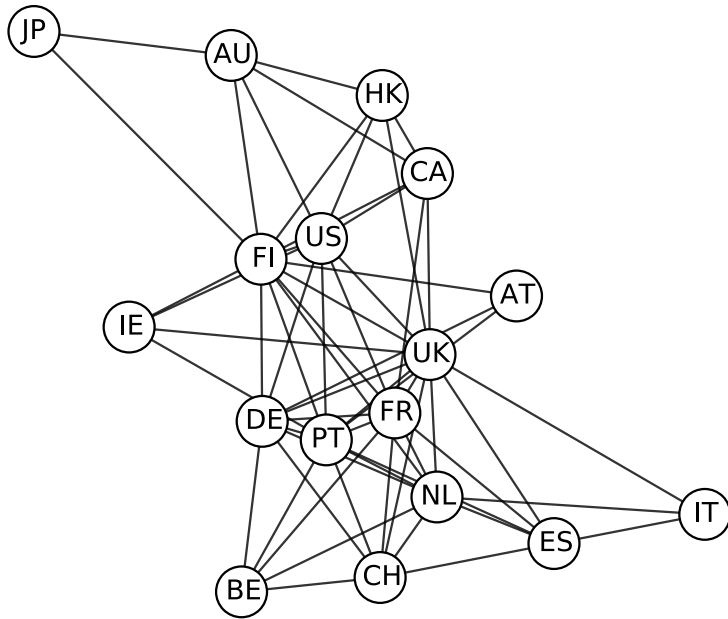
$$d_k \mid S_k \sim N_c(0, S_k)$$

$$\longrightarrow p(X_{1:p} \mid G) \approx \prod_{k=1}^T \frac{1}{\pi^p} \frac{h(W_k, \delta_k, G)}{h(W_k + P_k, \delta_k + 1, G)}$$

- For decomposable graphs

- Prior decomposes over cliques  $C$  and separators  $S$
- $h(W, \delta, G)$  decomposes over cliques  $C$  and separators  $S$
- **Marginal likelihood** decomposes over cliques  $C$  and separators  $S$

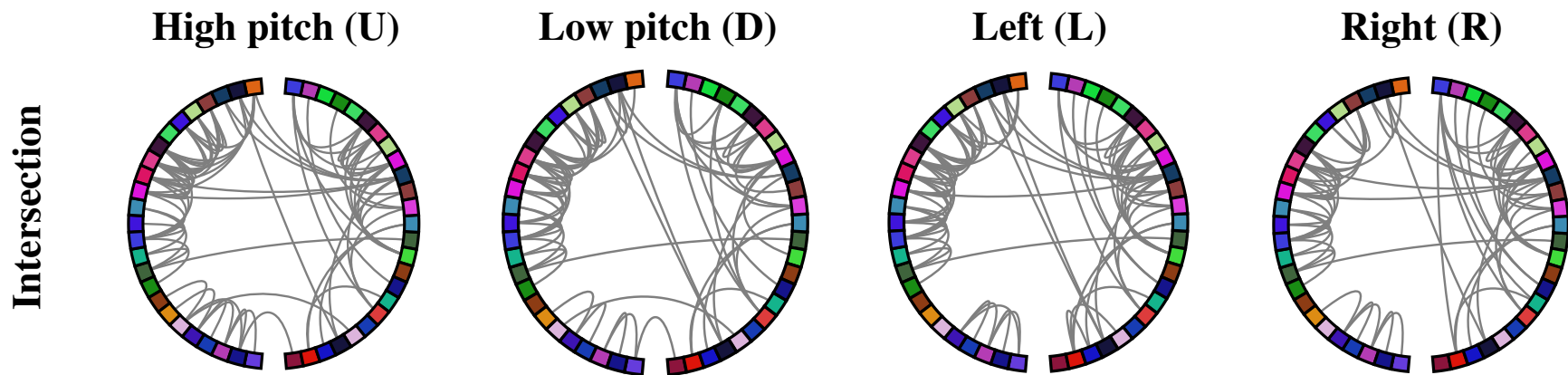
# Global Stock Indices



i.i.d. graph

# MEG Auditory Attention Task

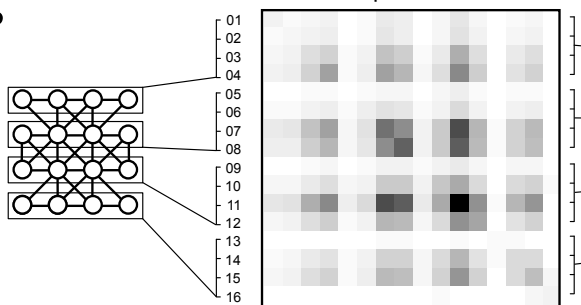
- Two tasks: (1) Focus attention, (2) Switch attention
- Four setups: high pitch, low pitch, left sound, right sound



# Graphs of Time Series Summary

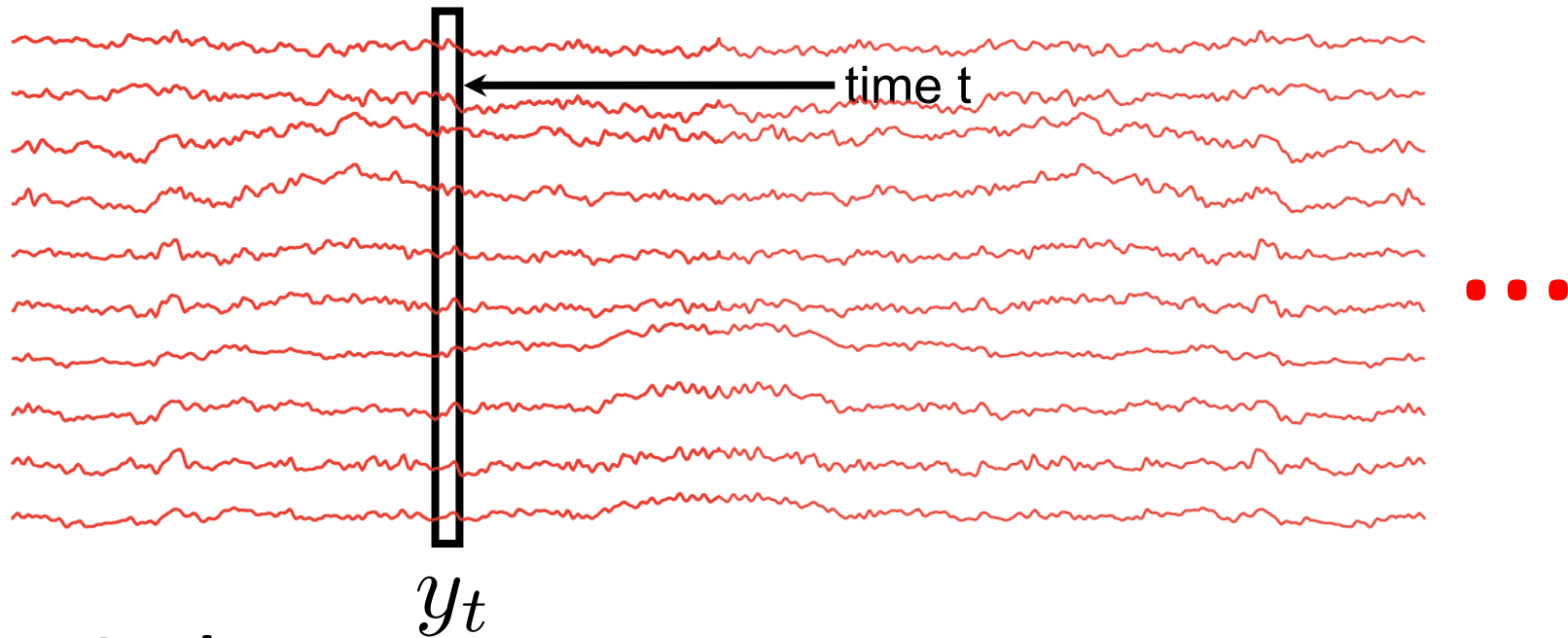
$\Sigma^{-1}$  sparse

Gaussian  
Graphical  
Model



Zeros = no edge in graph,  
Cond. ind. between nodes

- **Goal:** Infer conditional independencies between time series
- Efficient representation via **spectral density matrix**
  - Conditional independencies encoded by zeros in *inverse spectral density matrices*
- Whittle likelihood approximation defines tractable likelihood of data (Fourier coefficients) given spectral density matrices
- Defined **hyper complex inverse Wishart** prior
  - **Conjugate** prior on **graph-constrained** spectral density matrices
  - Enables closed-form **marginal likelihood** of data given graph



## Goals:

- Evolution – Dynamics across time
- Relational structure – Dependencies between series

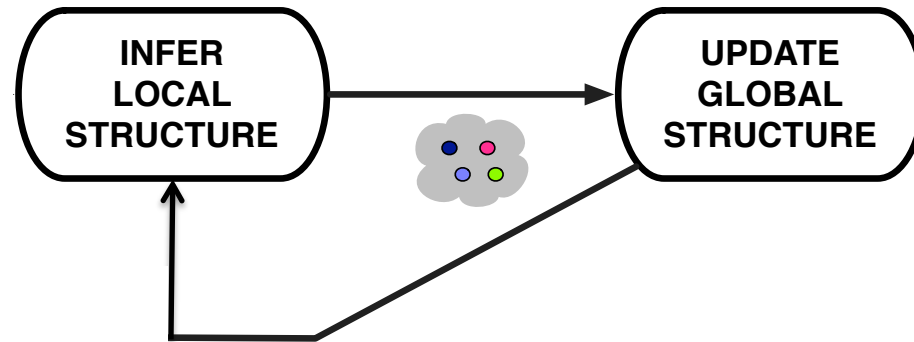
## Modeling challenges:

- Large  $p$  – Many dimensions/series
- Irregular grid of observations
- Missing values
- Heterogeneous data sources
- ...

## Computational challenges:

- Large  $n$  – Long time series
- Streaming data –  
Continuum of observations

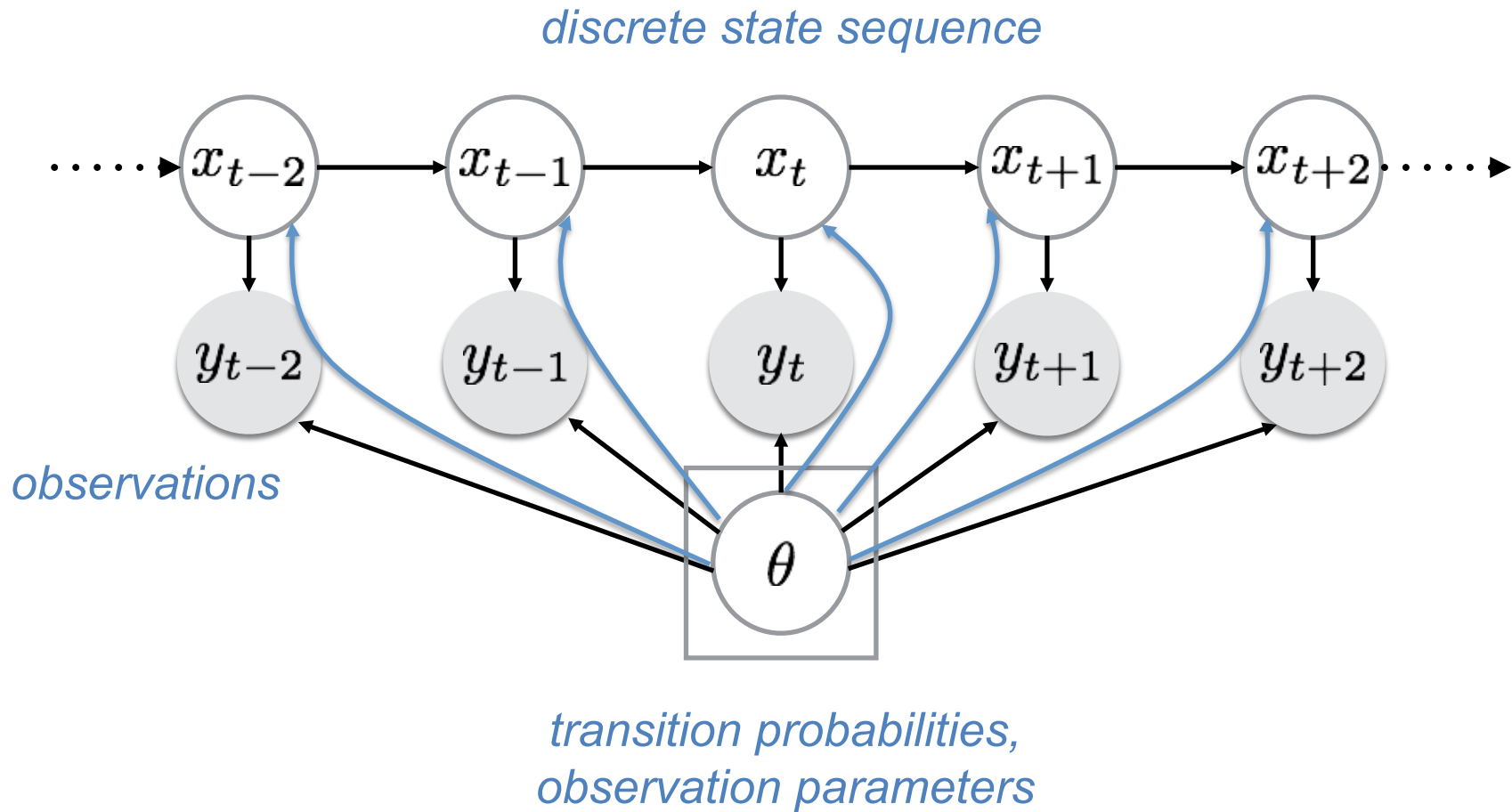
# Minibatch-Based Algorithms



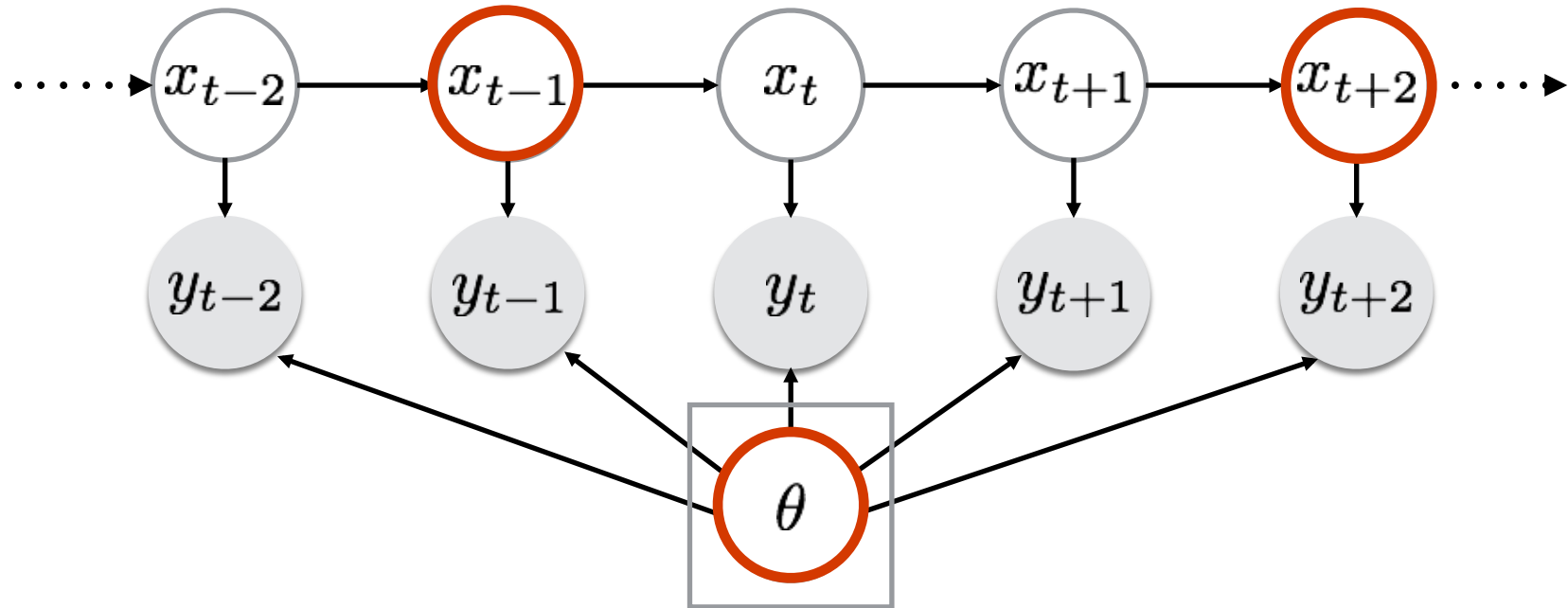
- Many ML/stat algorithms (e.g., gradient descent, Gibbs sampling,...) iterate between
  - operations involving **all data**
  - **updating parameters**
- Costly for large data / infeasible for streaming data
- Common approach for scalability:
  - **subsample data** → noisy operation
  - **noisy update** of parameters

**Not appropriate for dependent data**

# Hidden Markov Models (HMMs)



# Minibatches for HMMs

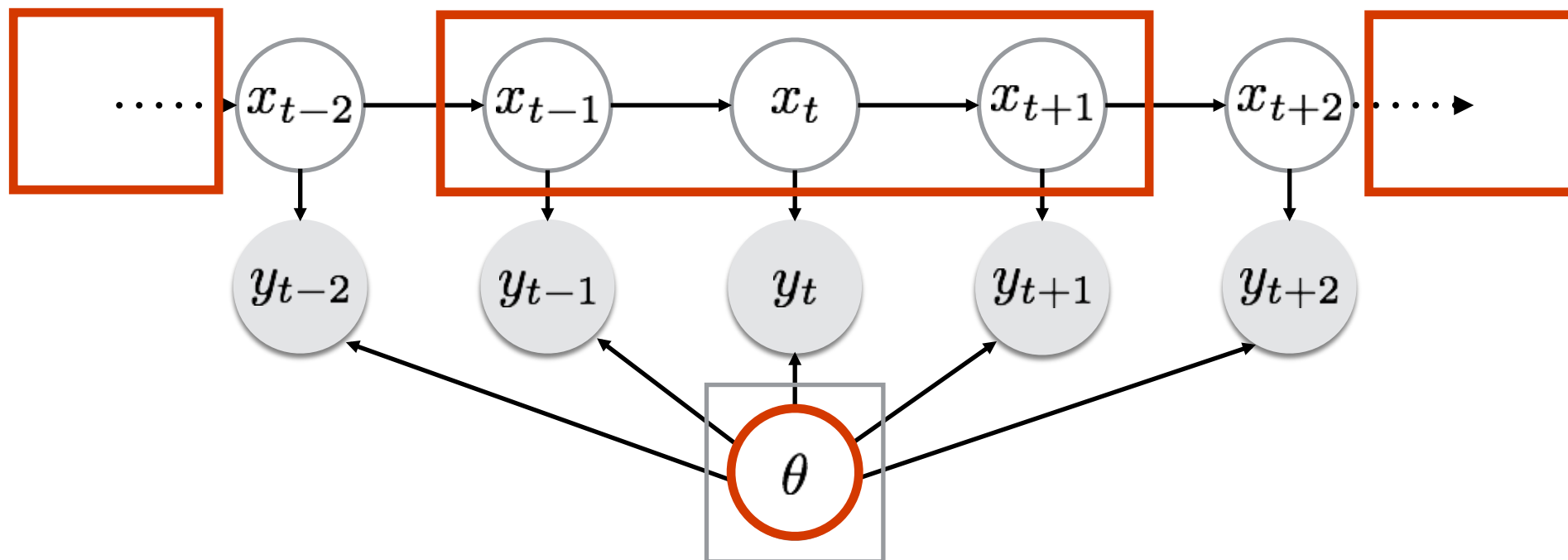


- Why not just subsample observations independently?
- Cannot learn transition structure

$$p(\mathbf{y}, \mathbf{x}, \theta) = p(\theta)\pi(x_1) \prod_{t=2}^T p(x_t | x_{t-1}, \theta_A) p(y_t | x_t, \theta_\phi)$$

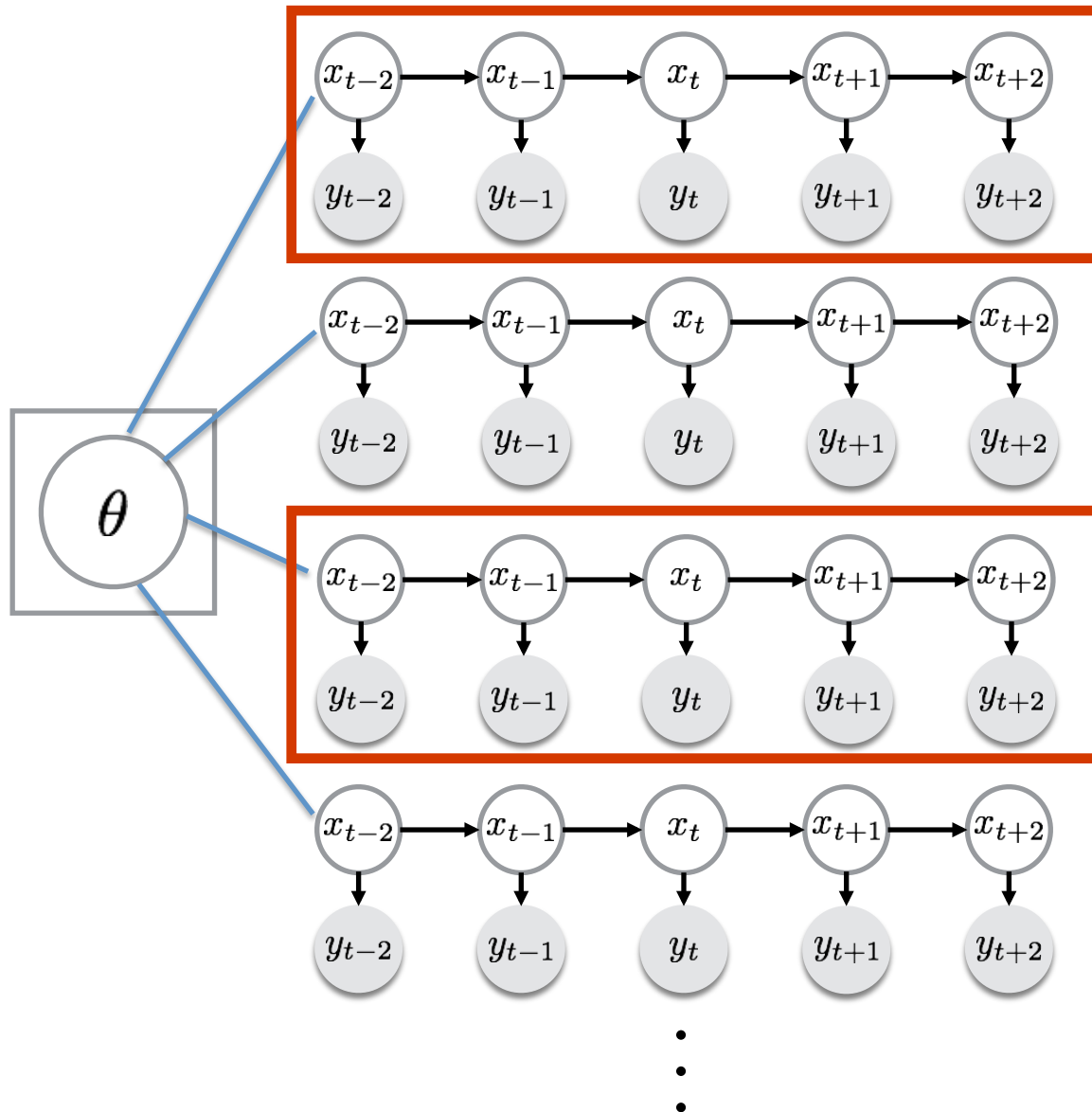


# Minibatches for HMMs



- How about sampling *subchain*?  $x^S = (x_{t-L}, \dots, x_t, \dots, x_{t+L})$
- Do we just *sever dependencies* between subchains and analyze *separately*?

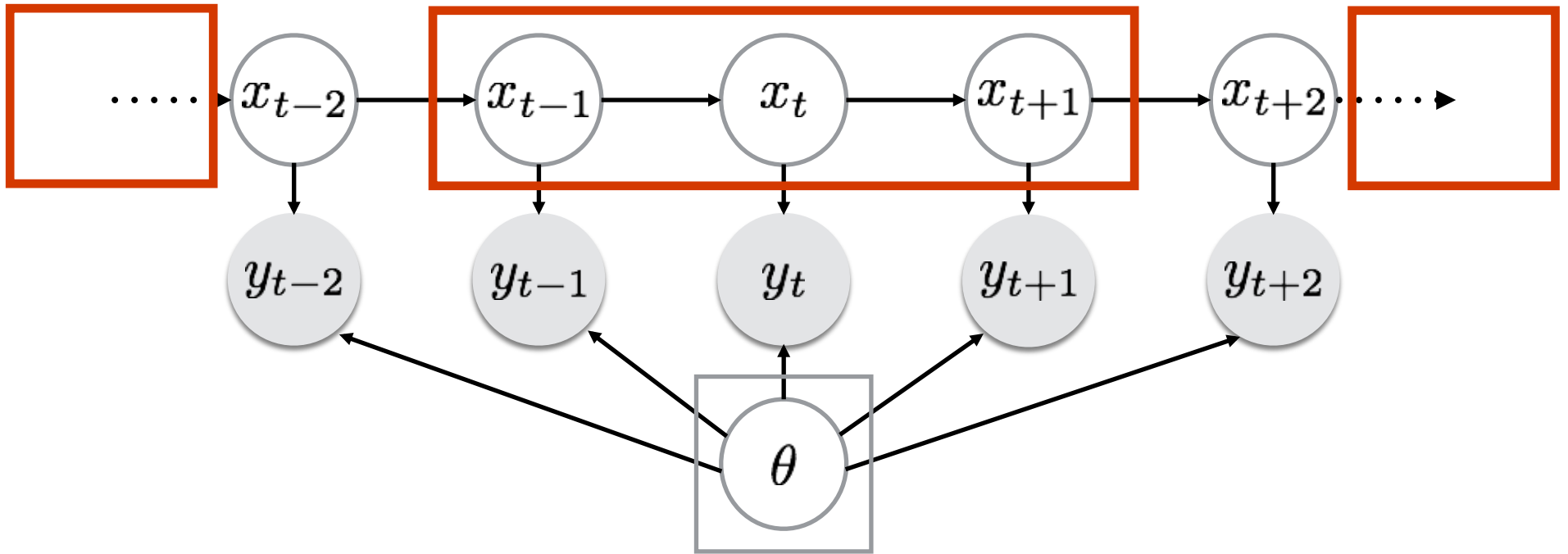
# Large Collections of Short Chains



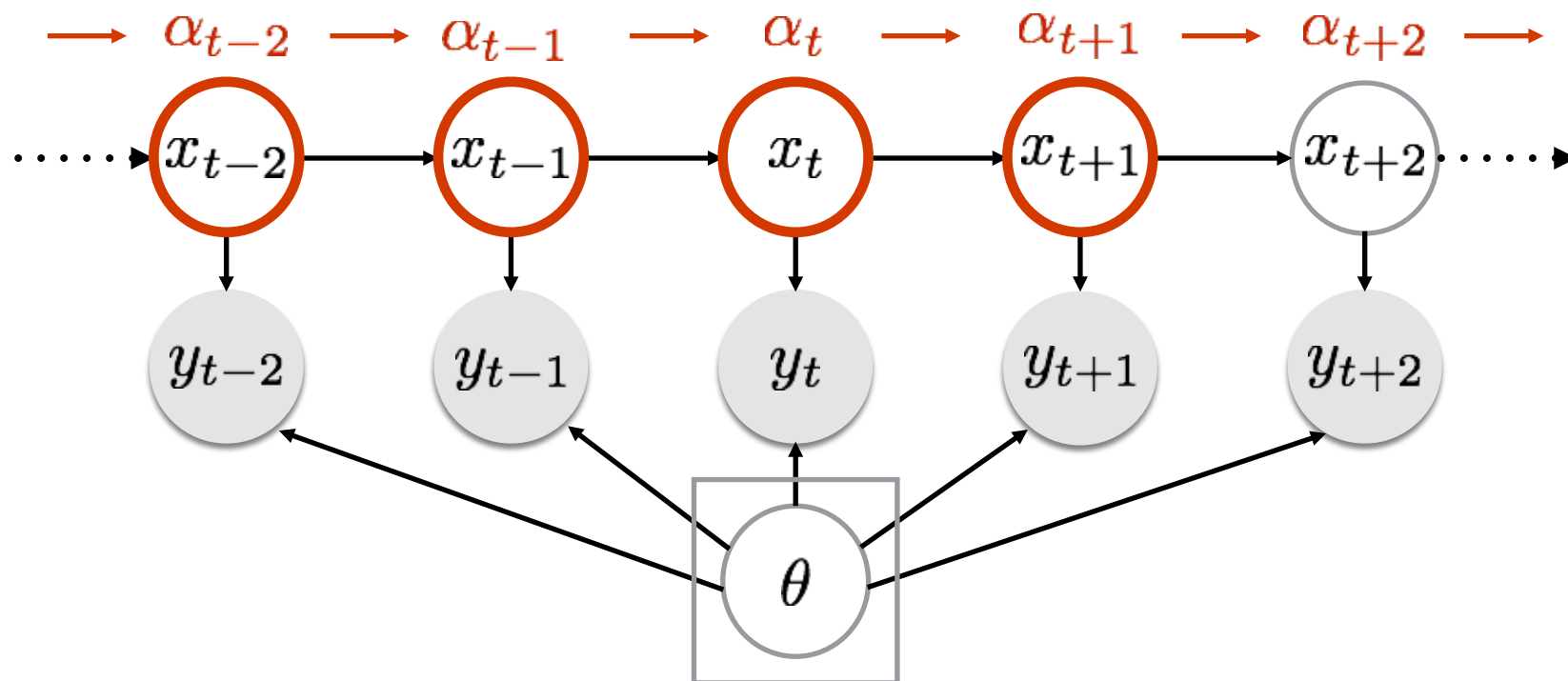
Johnson and Willsky,  
ICML 2014

Hughes et al.,  
preprint

# One Long Chain



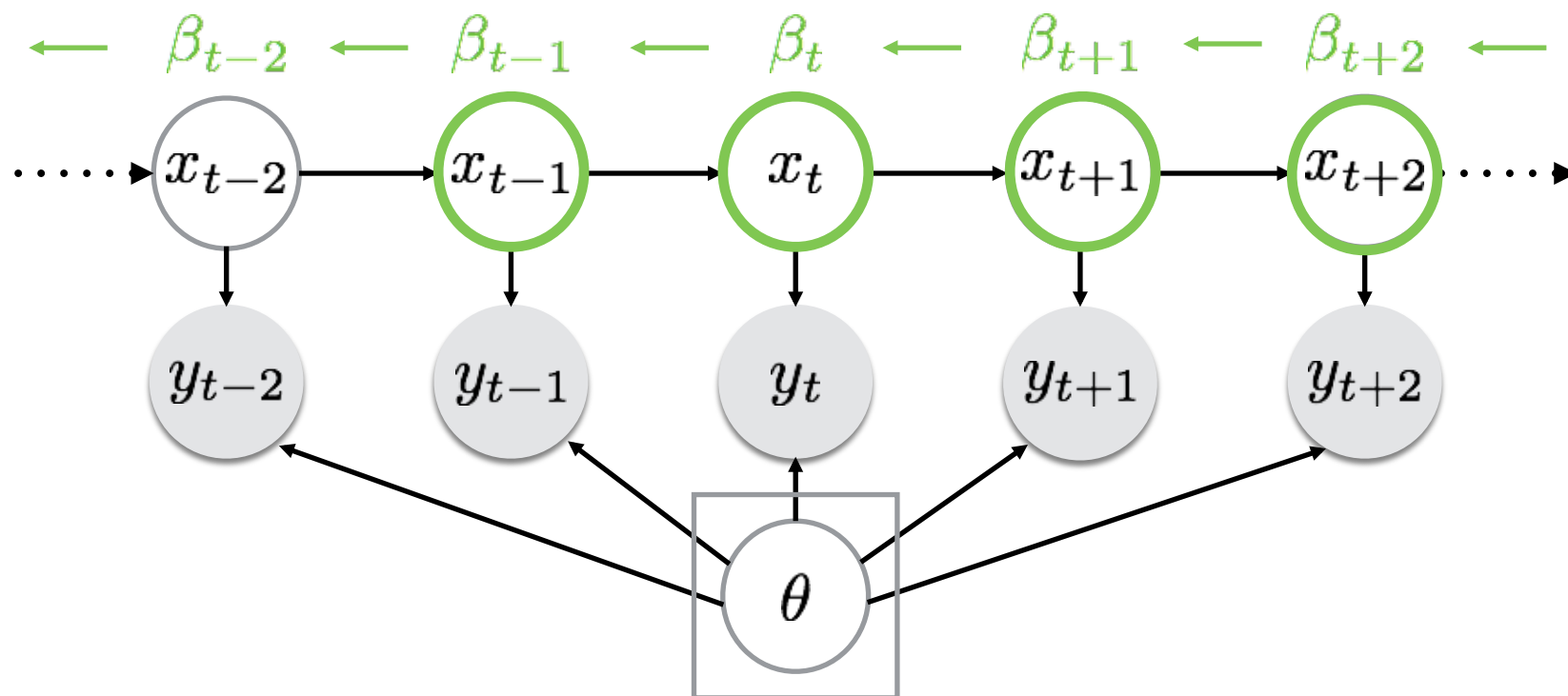
# Batch Learning for HMMs



- Use current  $\theta$  to form local state beliefs:
  - Propagate info forwards to form  $\alpha_t = p(y_1, \dots, y_t, x_t)$

$$\alpha_{t+1,k} = p(y_{t+1} \mid x_{t+1} = k) \sum_{j=1}^K \alpha_{t,j} p(x_{t+1} = k \mid x_t = j)$$

# Batch Learning for HMMs

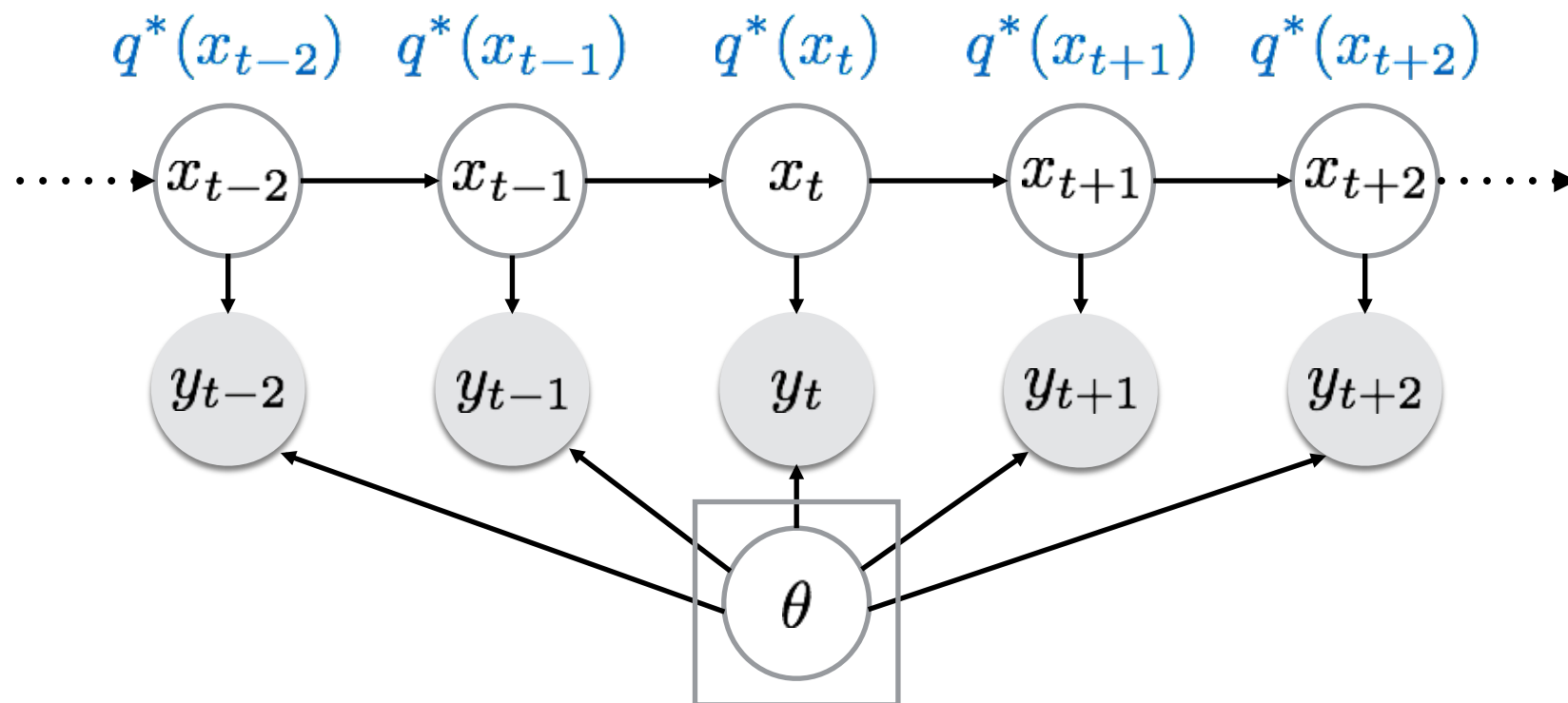


- Use current  $\theta$  to form local state beliefs:

– Propagate info backwards  $\beta_t = p(y_{t+1}, \dots, y_T | x_t)$

$$\beta_{t,k} = \sum_{j=1}^K p(y_{t+1} | x_{t+1} = j) p(x_{t+1} = j | x_t = k) \beta_{t+1,k}$$

# Batch Learning for HMMs

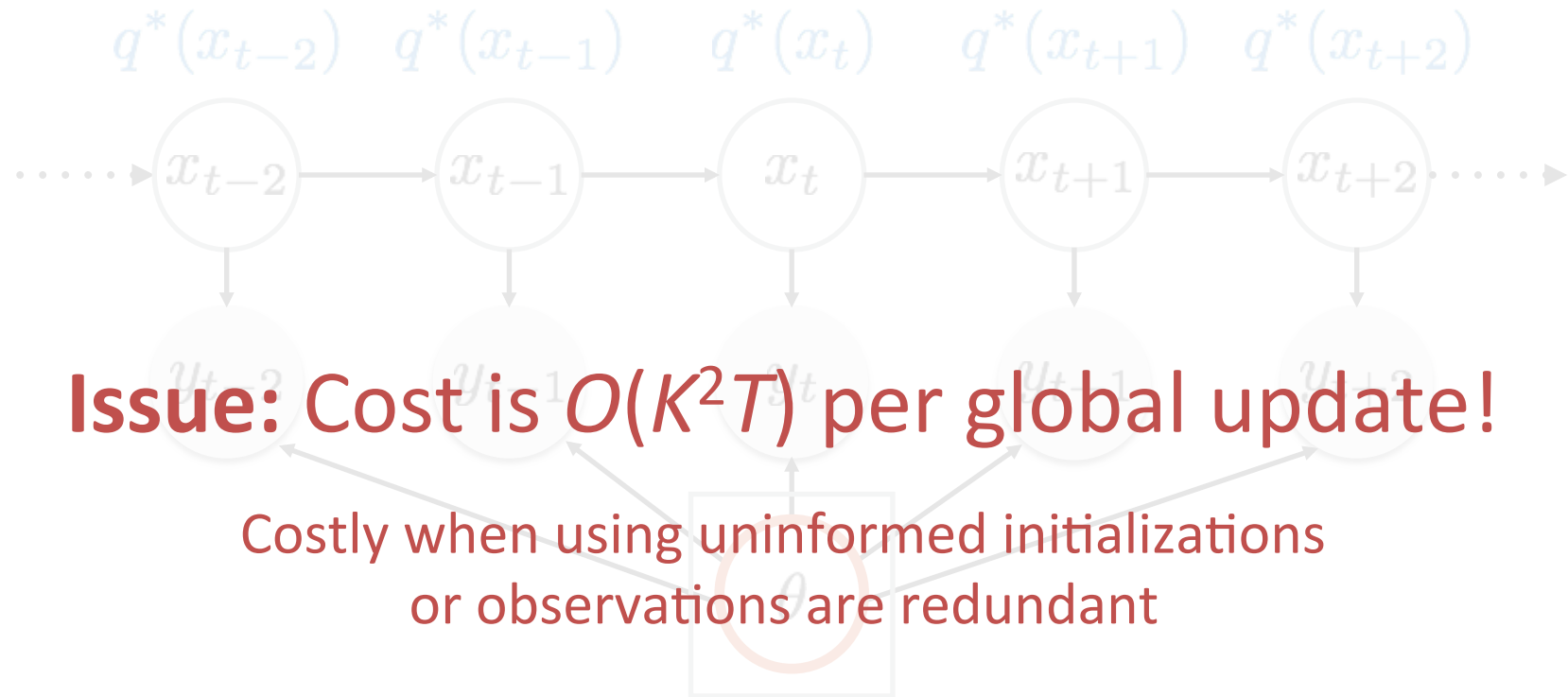


- Combine to form *smoothed* local state belief:

$$p(x_t \mid y_1, \dots, y_T) \propto \alpha_t \beta_t$$

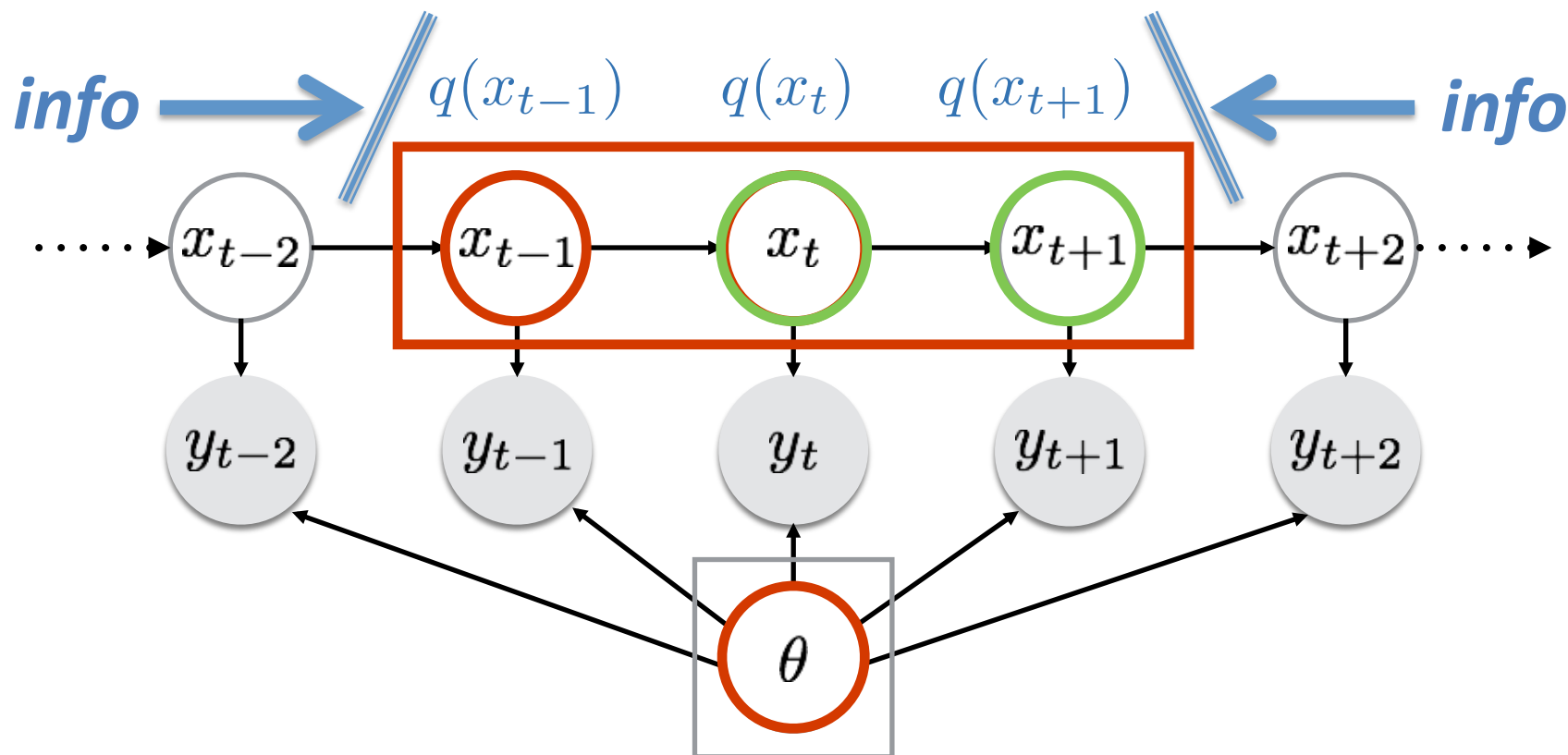
$q^*(x_t) \propto \alpha_t \beta_t$

# Batch Learning for HMMs



- Given local beliefs, update global parameter  **$T=250$  million**

# Minibatch Inference for HMMs



- Form **local** beliefs  $q(x_t) \propto \tilde{\alpha}_t \tilde{\beta}_t \rightarrow$  perform **global** update

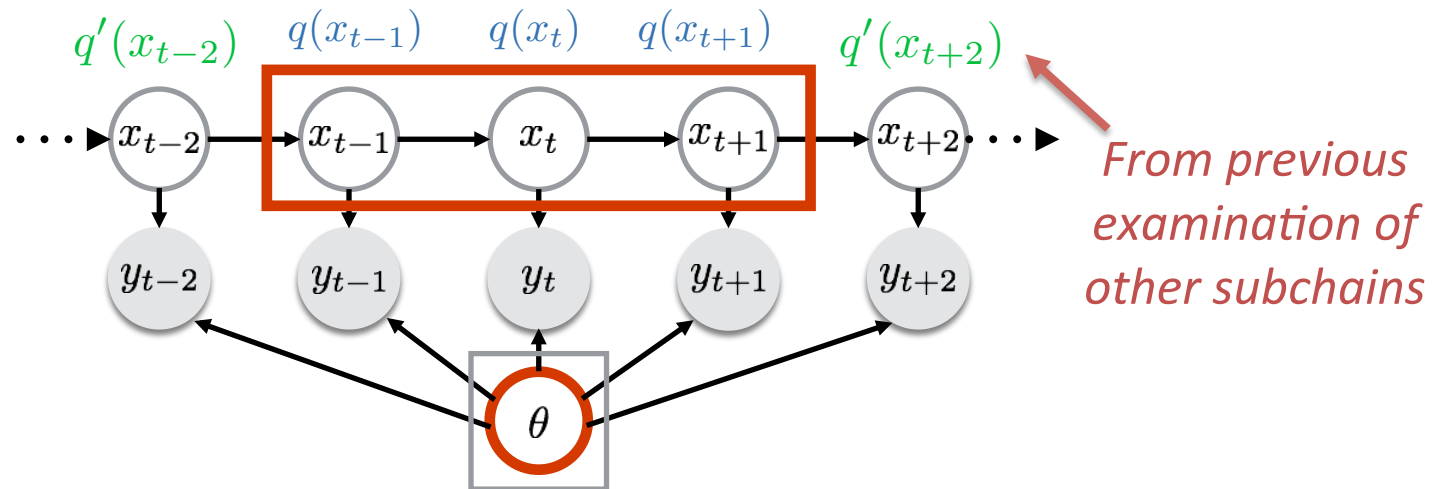
Local forward message

Local backward message



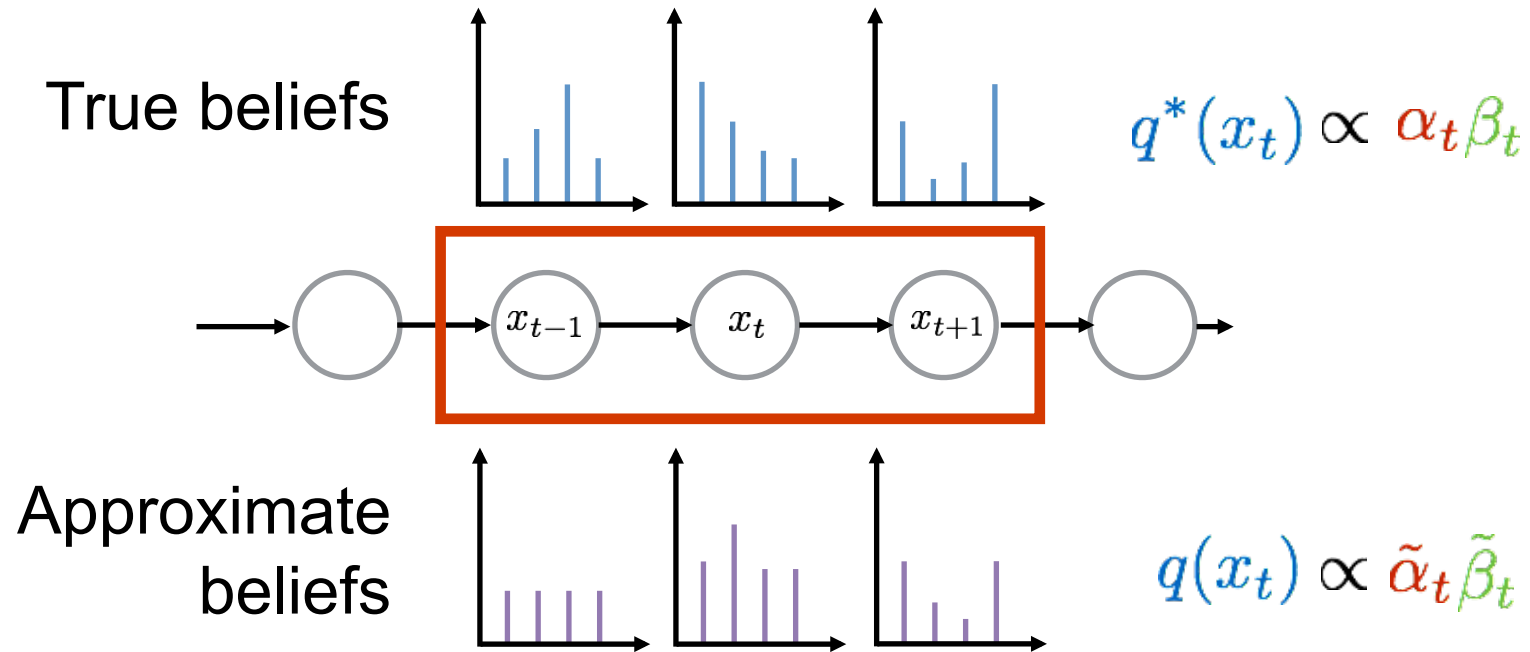
# Storage Limitations

- Can local message passing harness previous beliefs on nodes outside the subchain? **NO!**



- T=250 M obs x K=25 latent states  $\longrightarrow$  25 GB storage
- Need constant space algorithm  
 $\rightarrow$  *can't remember past beliefs*

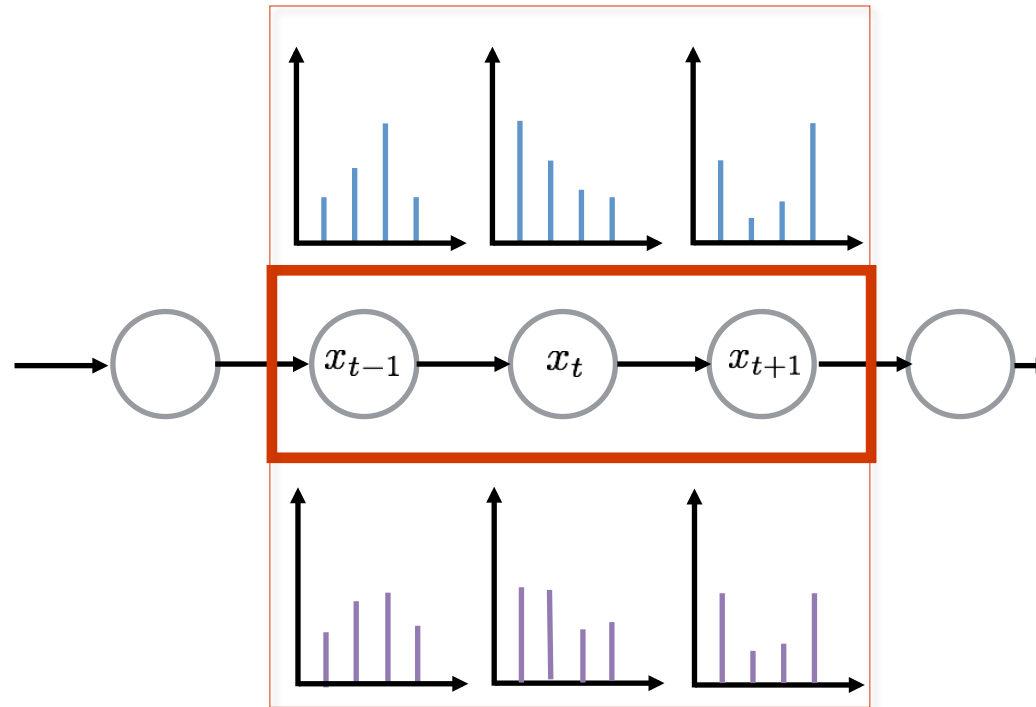
# Harnessing Memory Decay



Do we expect  $x_t$  to influence  $x_{t+1,000,000}$ ?

**Leverage memory decay**

# Buffering Subchains



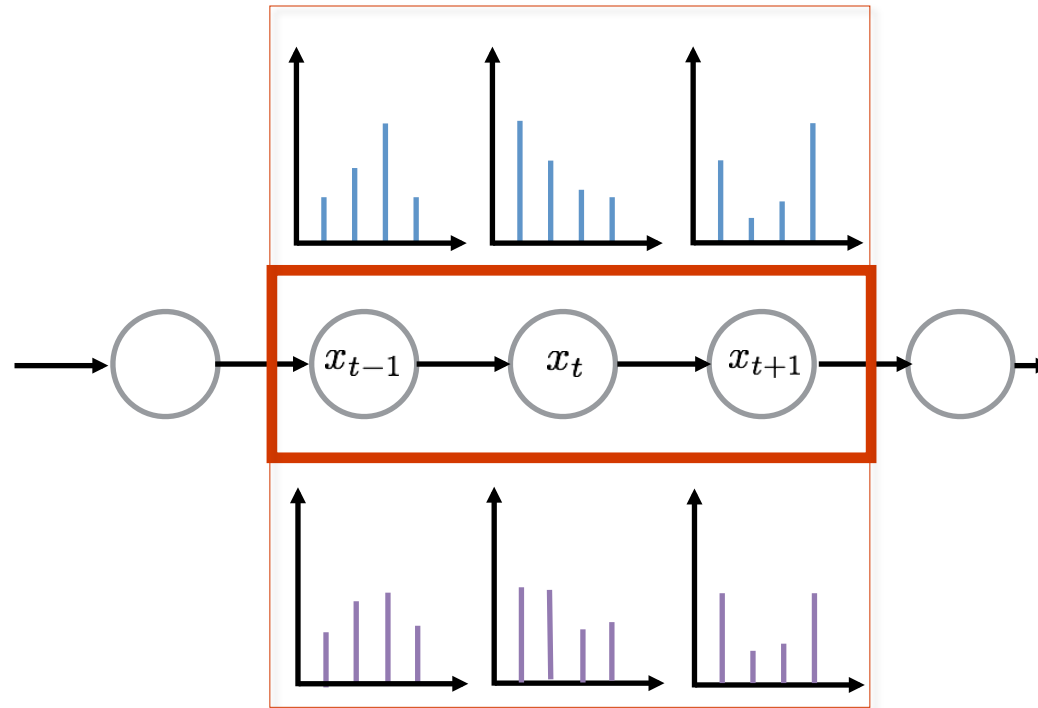
$$q^*(x_t) \propto \alpha_t \beta_t$$

$$q(x_t) \propto \tilde{\alpha}_t \tilde{\beta}_t$$

Check that subchain marginals are approximated well:

$$\max_{i \in S} \|q(x_i) - q^*(x_i)\| < \epsilon$$

# Buffering Subchains



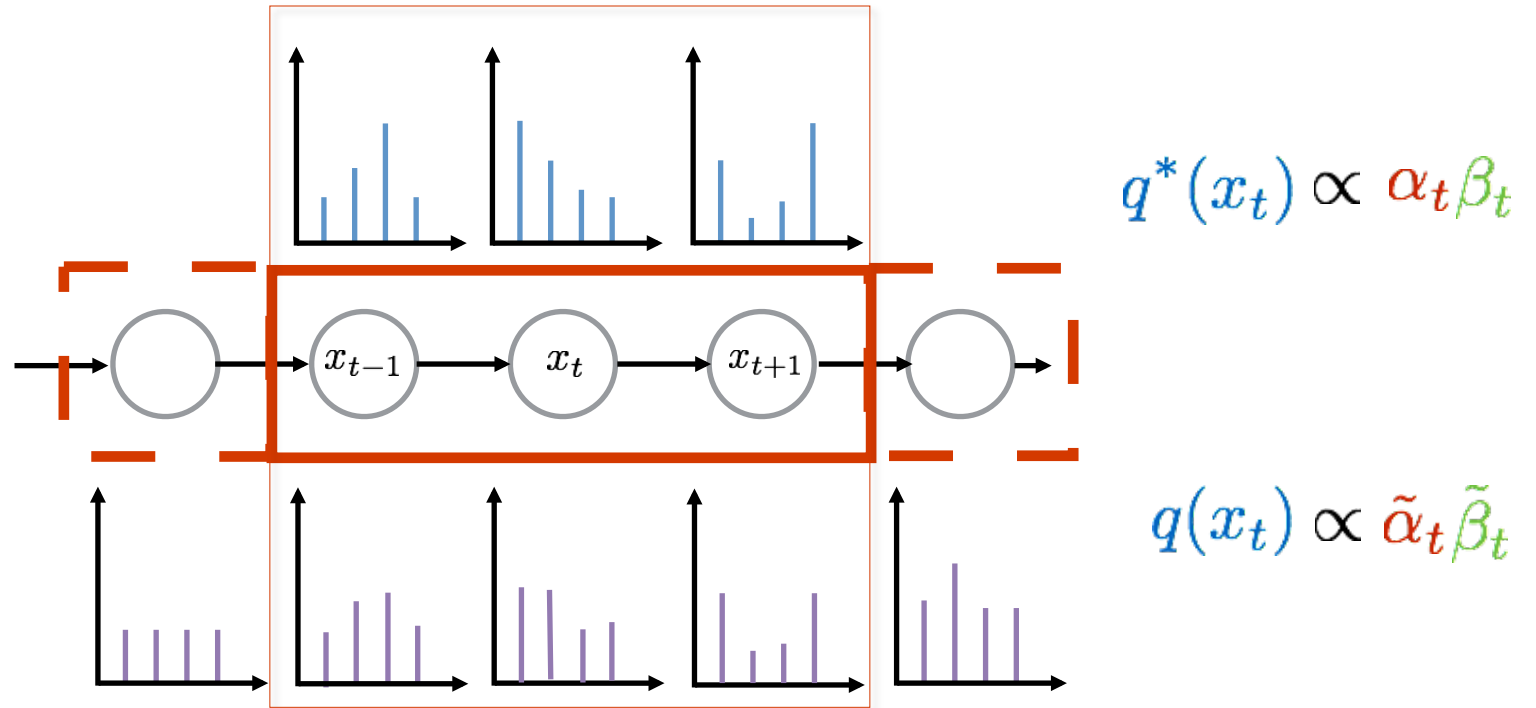
$$q^*(x_t) \propto \alpha_t \beta_t$$

$$q(x_t) \propto \tilde{\alpha}_t \tilde{\beta}_t$$

?

$$\max_{i \in S} \|q(x_i) - q^*(x_i)\| < \epsilon$$

# Buffering Subchains



?

$$\max_{i \in S} \|q(x_i) - q^*(x_i)\| < \epsilon$$

# Buffering Subchains

– Only need limited buffer

$$q^*(x_t) \propto \alpha_t \beta_t$$

– Complexity is now  $O(K^2 L_{buffer}^{x_{t+1}})$  per iteration

Large savings for  $L + \text{buffer} \ll T$

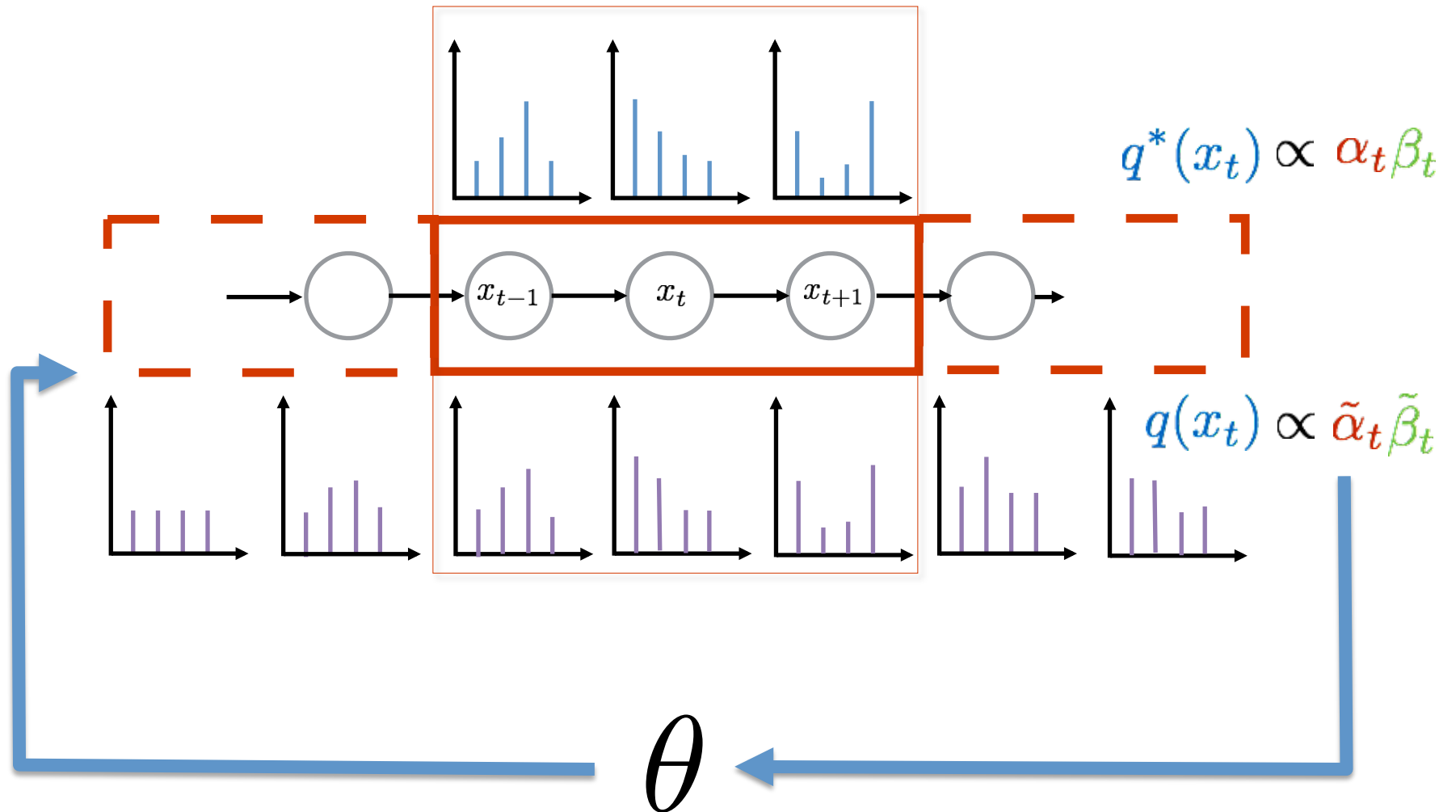
$$q(x_t) \propto \tilde{\alpha}_t \tilde{\beta}_t$$

– Similar idea as Splash BP (parallelizing BP)

[Gonzalez, et. al. 2009]

*But, uncertain parameter setting here*

# Buffering for Learning

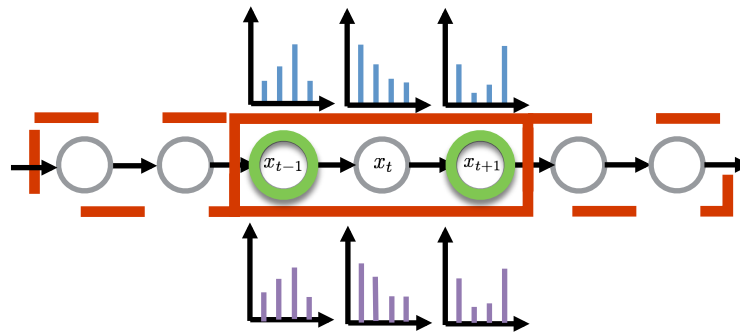


# Buffering in Practice

- We do not actually know the true marginals
- Monitor changes in approximate subchain beliefs:

$$\max_{i \in S} \left| \left| q(x_i)^{\text{new}} - q(x_i)^{\text{old}} \right| \right| < \epsilon$$

- Chain structuring implies that only endpoints must be checked



- During buffer expansions, forward-backward passes can reuse computations of previous buffer



# Variational Bayes (VB)

- Approximate posterior with variational distribution

$$\begin{array}{c} \text{parameters} \\ \downarrow \\ p(x, \theta | y) = \frac{p(y|x, \theta)p(x, \theta)}{p(y)} \approx q(x, \theta) \\ \uparrow \quad \uparrow \\ \text{latent variables} \quad \text{observations} \end{array}$$

- Minimize  $\text{KL}(q||p) \leftrightarrow$  maximize “ELBO”:

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(y, x, \theta)] - \mathbb{E}_q[\log q(x, \theta)] \leq \log p(y)$$

- Common to make mean-field assumption:

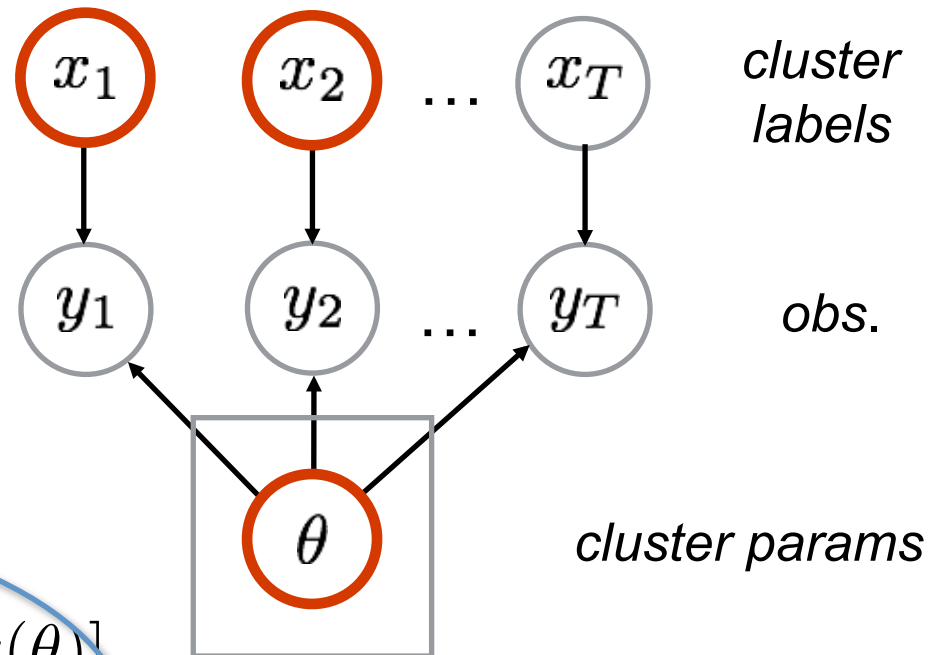
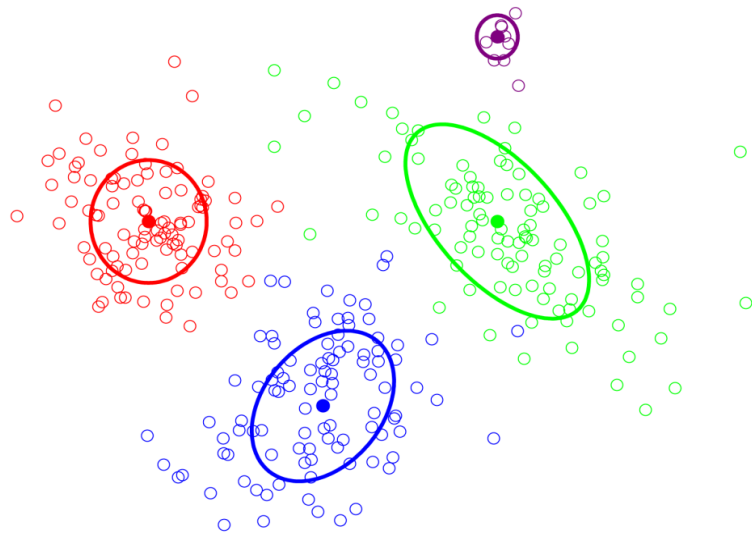
$$q(x, \theta) = q(x)q(\theta)$$

# Variational Methods Cartoon

- Cartoon of goal:
- Variational distribution parameterized by *variational free parameters*
- *Objective: optimize over free parameters to find “closest” distribution in variational family*

# VB Example: Mixture of Gaussians

Maximize ELBO with *coordinate-ascent*  $\frac{\partial \mathcal{L}}{\partial q(\mathbf{x})} = 0 \longleftrightarrow \frac{\partial \mathcal{L}}{\partial q(\theta)} = 0$



$$\mathcal{L} = E_{q(\theta)} [\ln p(\theta)] - E_{q(\theta)} [\ln q(\theta)] + \sum_{i=1}^T E_{q(x_i)} [\ln p(y_i, x_i | \theta)] - E_{q(x_i)} [\ln q(x_i)]$$

# Stochastic Variational Inference (SVI)

- Batch VB global step requires touching *all of the data*

$$\mathcal{L} = E_{q(\theta)} [\ln p(\theta)] - E_{q(\theta)} [\ln q(\theta)] \\ + \sum_{i=1}^T E_{q(x_i)} [\ln p(y_i, x_i | \theta)] - E_{q(x_i)} [\ln q(x_i)]$$

- SVI uses stochastic gradient descent (SGD) for global update [Hoffman, et. al. 2013]

- Sample observation:  $x^S \sim \text{Unif}(x_1, \dots, x_T)$
- Follow noisy, unbiased estimate of natural gradient of  $\mathcal{L}$ .

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + \rho_t \tilde{\nabla}_{\mathbf{w}} \mathcal{L}^S \quad \mathbb{E}_S[\tilde{\nabla}_{\mathbf{w}} \mathcal{L}^S] = \tilde{\nabla}_{\mathbf{w}} \mathcal{L}$$

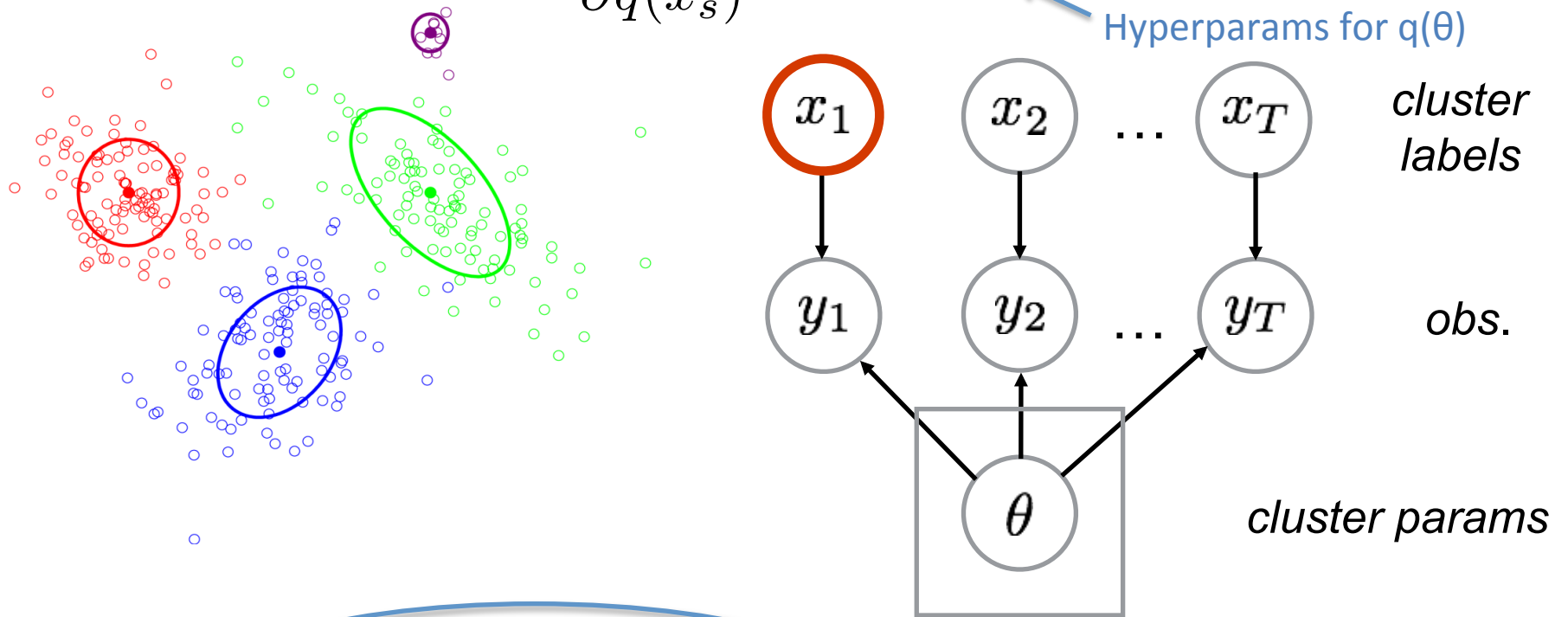
 Variational parameters defining  $q(\theta)$

# SVI Example: Mixture of Gaussians

Maximize ELBO with *stochastic gradient descent*

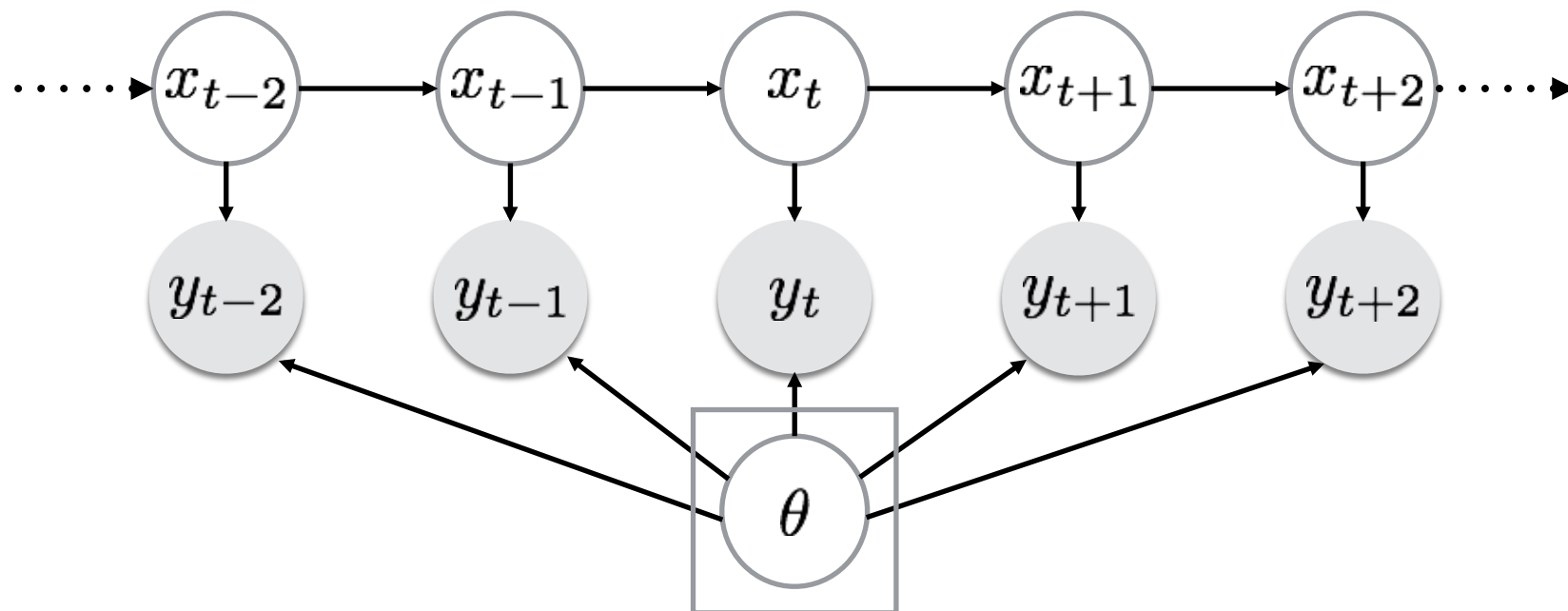
$$\frac{\partial \mathcal{L}^s}{\partial q(x_s)} = 0 \iff \mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + \rho_t \tilde{\nabla}_{\mathbf{w}} \mathcal{L}^s$$

Hyperparams for  $q(\theta)$



$$\mathcal{L}^s = E_{q(\theta)} [\ln p(\theta)] - E_{q(\theta)} [\ln q(\theta)] + T \cdot (E_{q(x_s)} [\ln p(y_s, x_s | \theta)] - E_{q(x_s)} [\ln q(x_s)])$$

# Structured Mean Field Approximation



- Use structured mean-field approximation:

$$p(x_1, x_2, \dots, x_T, \theta \mid y_1, y_2, \dots, y_T) \approx q(x_1, x_2, \dots, x_T)q(\theta)$$

# SVI for HMMs

(Approx) coordinate ascent step:

Function of  $q(\theta)$

$$\alpha_{t+1,k} = \tilde{p}(y_{t+1} | x_{t+1} = k) \sum_{j=1}^K \alpha_{t,j} \tilde{A}_{j,k}$$

$$\beta_{t,k} = \sum_{j=1}^K \tilde{p}(y_{t+1} | x_{t+1} = j) \tilde{A}_{k,j} \beta_{t+1,k}$$

$q^*(x_t) \propto \alpha_t \beta_t$

$q(x_t) \propto \tilde{\alpha}_t \tilde{\beta}_t$

$$q(\theta)$$

Stochastic natural gradient step:

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + \rho_t \tilde{\nabla}_{\mathbf{w}} \mathcal{L}^S$$

Function of  $q(\mathbf{x})$

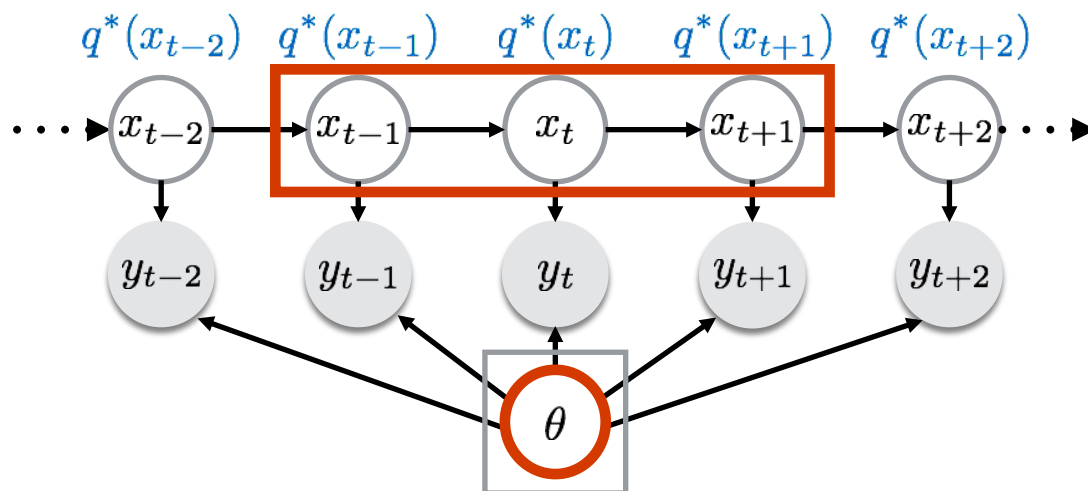
# Differences from i.i.d. Case

- Minibatches are *correlated*
  - Data in one is not independent of data in another
- Minibatch marginals  $\neq$  batch marginals
  - Impact of latent chain
  - Mitigated by buffering



# Correlated Minibatches

- Pretend we have exact local distribution  $q^*(x^S)$



As if we had run batch forward-backward

- Typical arguments for convergence to local mode rely on *unbiased* + *independent* noisy gradients [c.f., Bottou 1998, Hoffman 2013]
  - Our SGs are *dependent* since subchains are correlated
- Using [Polyak and Tsytkin 1973], unbiasedness suffices for convergence of  $\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + \rho_t \tilde{\nabla}_{\mathbf{w}} \mathcal{L}^S$

# Global Update – Unbiasedness

- In mixture model case with **uniform** sampling of **observation  $\mathbf{s}$** , unbiasedness was preserved via:

$$\mathcal{L}^s = E_{q(\theta)} [\ln p(\theta)] - E_{q(\theta)} [\ln q(\theta)] \\ + T \cdot \left( E_{q(x_s)} [\ln p(y_s, x_s | \theta)] - E_{q(x_s)} [\ln q(x_s)] \right)$$

- In HMM case, our ELBO data term is

$$\ln p(\mathbf{y}, \mathbf{x} | \theta) = \ln \pi(x_1) + \sum_{t=2}^T \ln A_{x_{t-1}, x_t} + \sum_{i=1}^T \ln p(y_t | x_t)$$

- Does not decompose over individual  $x_t$
- Need to scale transition and emission terms separately
- Straightforward for **uniform** sampling of **subchains  $\mathbf{S}$**  of length  $L$ , assuming chain is observed at *stationarity*

# Effect of Approximated Marginals

## SVI-HMM iterates:

buffer minibatches to **approx  $q(x)$**   $\leftrightarrow$  **update  $q(\theta)$**   
*coordinate gradient step* *stochastic*  
*(natural) gradient step*

For  $\epsilon$  sufficiently small (sufficiently long buffer)

- Approximate marginals “close enough” to true marginals
- Noisy gradient in same half-plane as true gradient

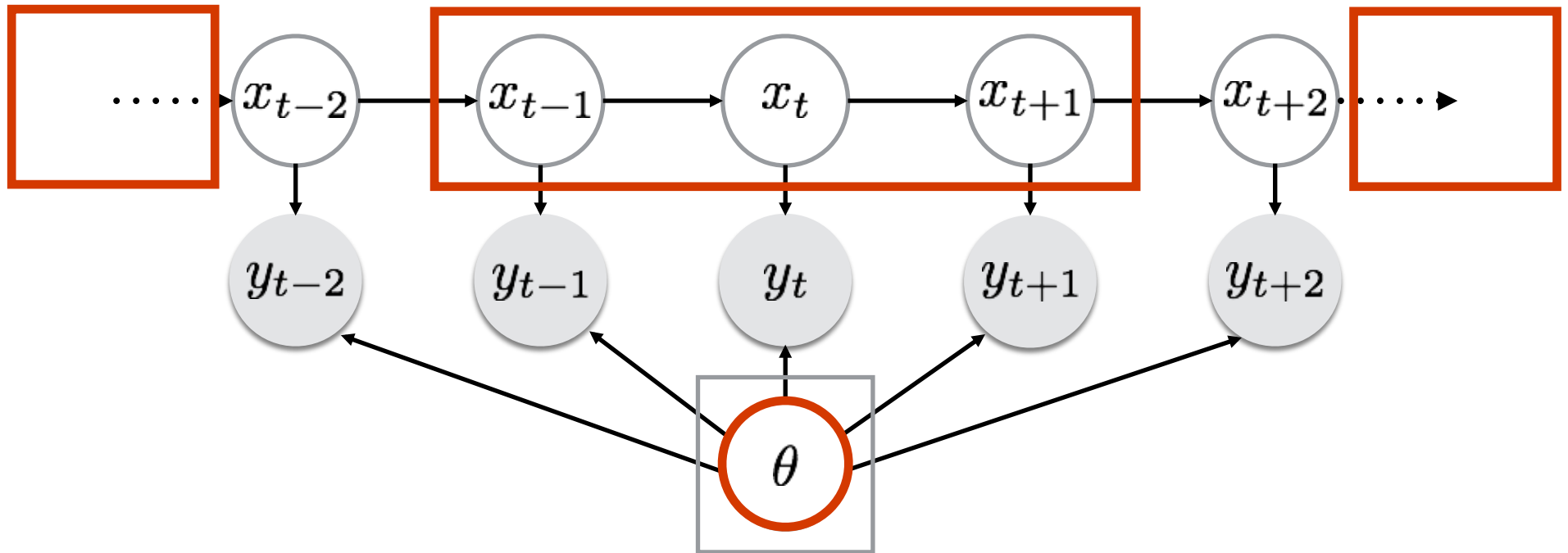


iterative algorithm converges to local mode of ELBO

# Experiments

- Synthetic data:
  - **Diagonally Dominant:** Long memory chain with large self-transitions
  - **Reversed Cycles:** Two overlapping cycles with opposite directions
- **Human chromatin application**

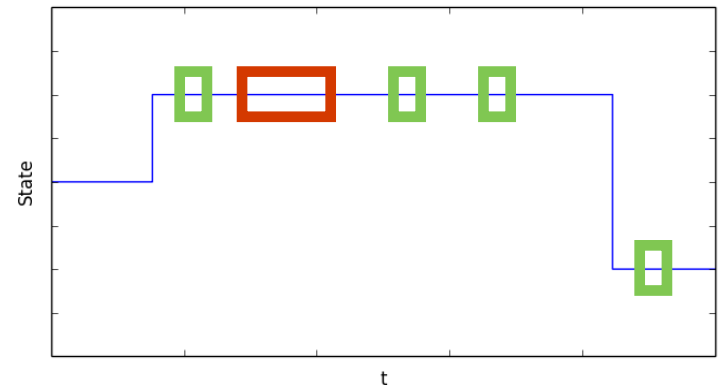
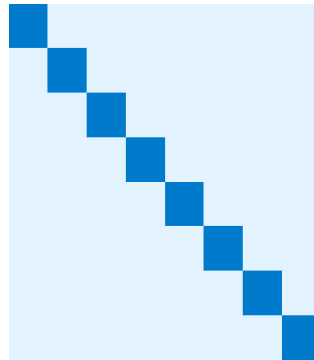
# Minibatch of Subchains



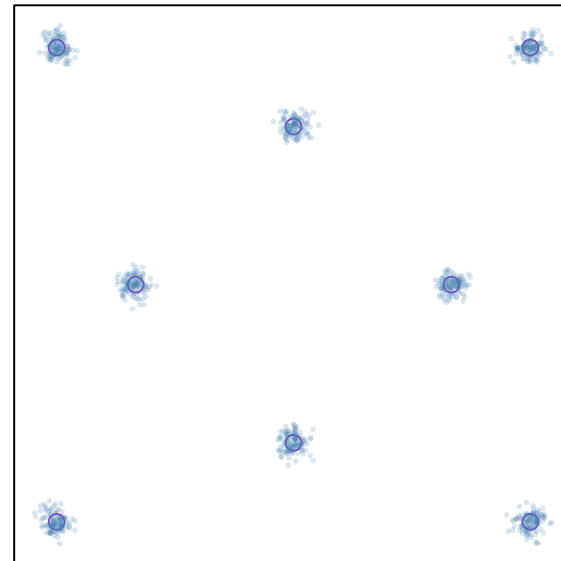
Minibatch consists of  $M$  subchains each of length  $L$

# Diagonally Dominant

- 8 latent states
- 2d Gaussian emissions

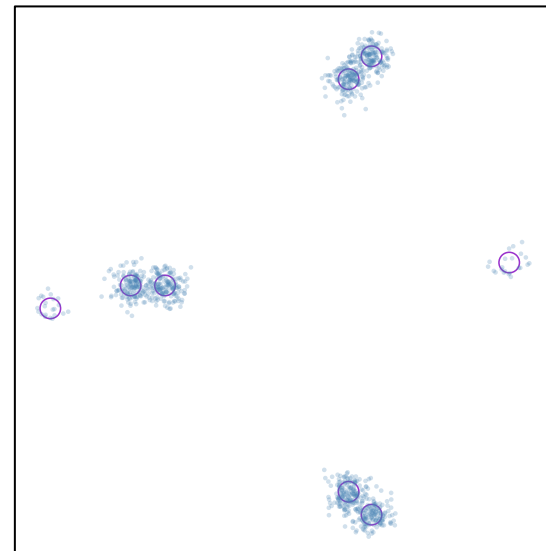
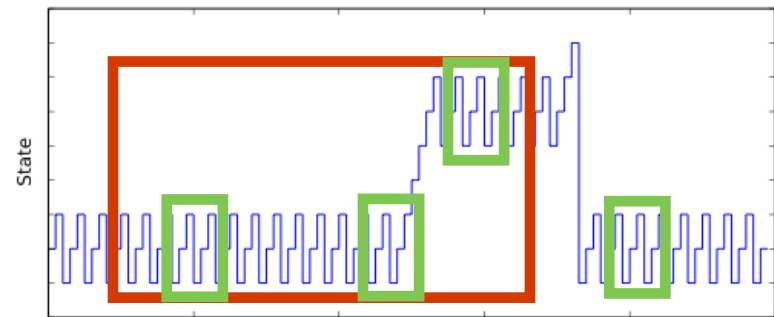
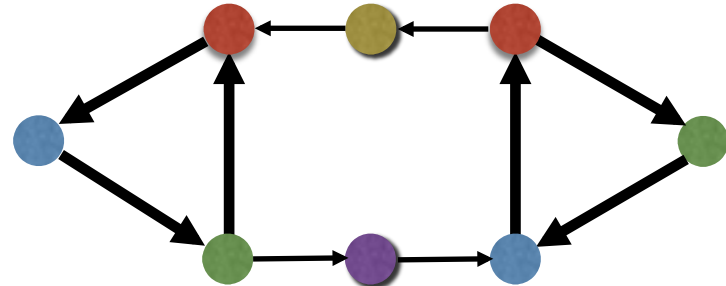


- *High auto-correlation*  
→ few long subchains converge slowly (small  $M$ , large  $L$ )
- *Emissions identifiable*  
→ many small subchains perform better (large  $M$ , small  $L$ )

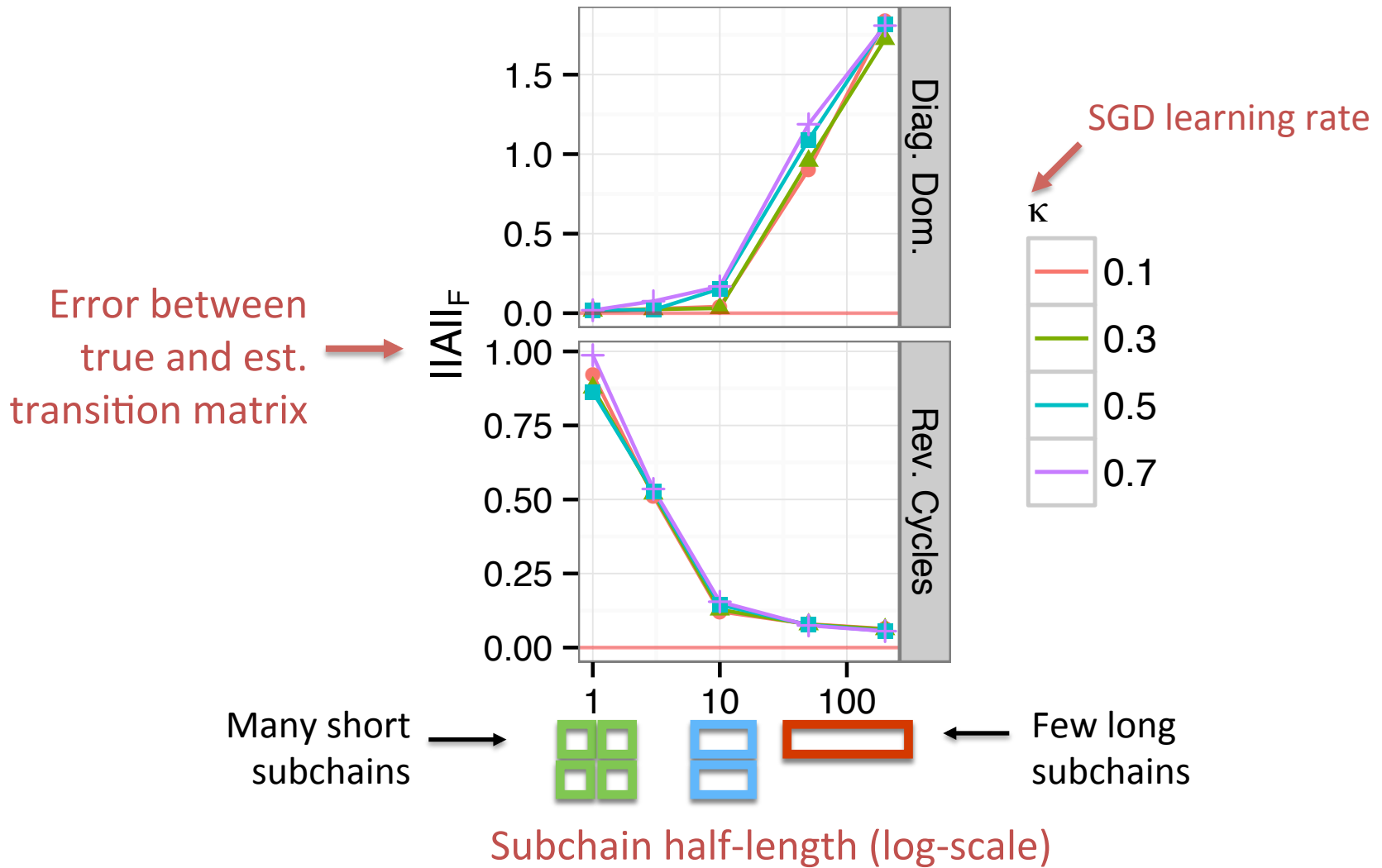


# Reversed Cycles

- 8 latent states
- 2d Gaussian emissions
- *Emission distributions overlap*
- *Direction* of cycles important to identify states
  - Singleton observations insufficient
  - Without buffering, need  $L > 3$  to learn effectively
- Longer subchains more likely to capture structure

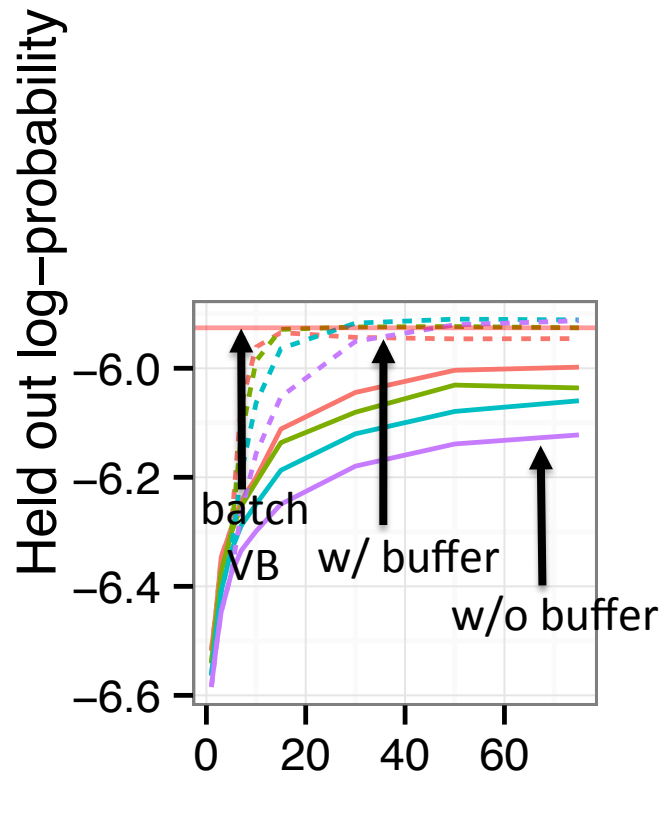
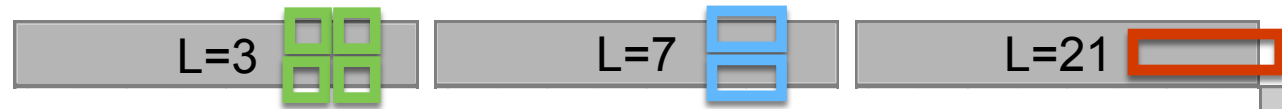


# Transition Matrix Recovery





# Subchain Buffering

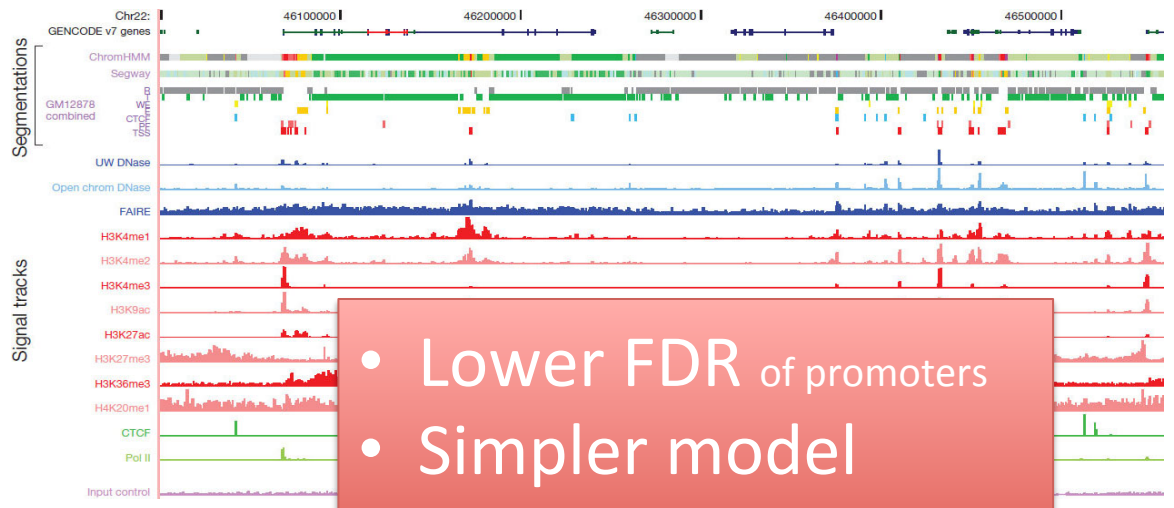


Diag. Dom.

Rev. Cycles

# Human Chromatin Segmentation

- Chromosome data from ENCODE project
- 12 dimensional observations
- **Goal:** segment sequences
- **T = 250 million**



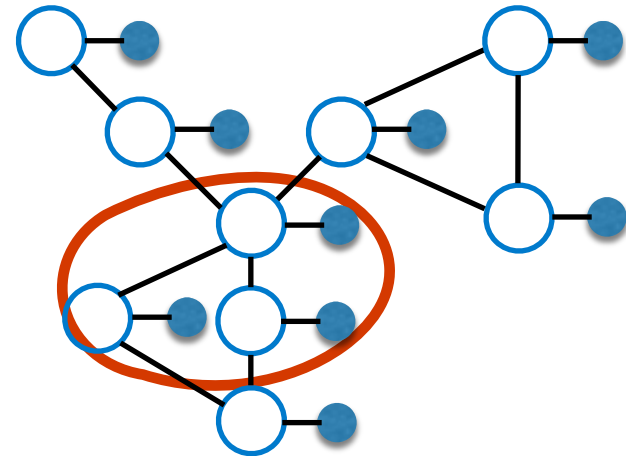
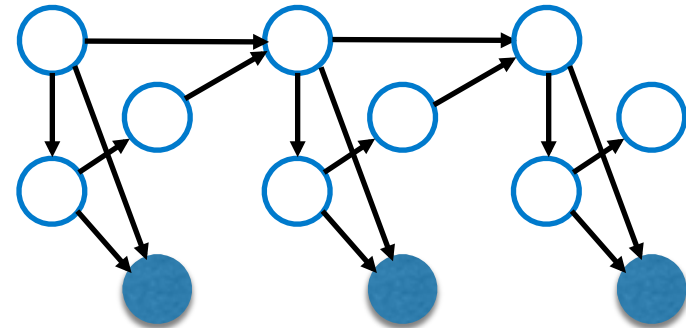
• Lower FDR of promoters  
• Simpler model  
• Uses all of the data

- [Hoffman et. al. 2012] used **dynamic Bayesian network**
    - Broke sequence into pieces to perform inference via EM
    - Severs long-range dependencies
- ← Runtime = **days**

- Adaptive subsampling on **HMM** (*simpler model*)
- Runtime = under **1 hr** →

# BNP and Other Extensions

- Presented finite HMM case, but ideas could generalize to:
  - Nonparametric HMMs
    - Truncation plus split-merge to change the number of states  
[Bryant & Sudderth, 2012]
  - DBN and MRF models
- Applications to:
  - Large spatial fields
  - Spatio-temporal data, etc.



# Overall Summary

- Scalable Bayesian dynamic modeling:
  - *Low-dimensional embeddings* with application to MEG word classification
  - *Clusters* for forming high-resolution housing value index
  - *Graphs of time series* with application to stocks + functional connectivity
- Scalable Bayesian computations in dynamic models
  - Harness *memory decay* to use subset-based methods in HMMs

