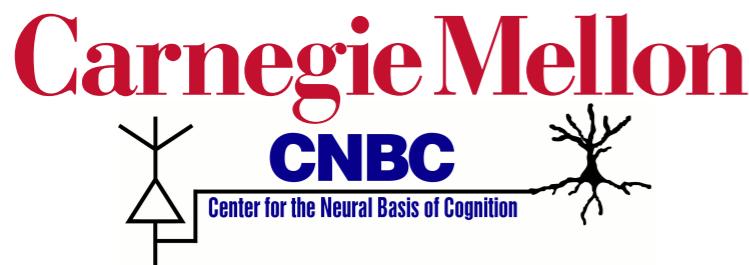


Hierarchical Latent Dictionaries for Models of Brain Activation

Alona Fyshe*, Emily Fox, David Dunson and Tom Mitchell



Wharton
UNIVERSITY OF PENNSYLVANIA

Duke
UNIVERSITY

afyshe@cs.cmu.edu

Multi-sensor time series

- Multiple noisy recordings of the same process
 - Many sensors (high dimension)
 - Changing sensor correlation

Multi-sensor time series

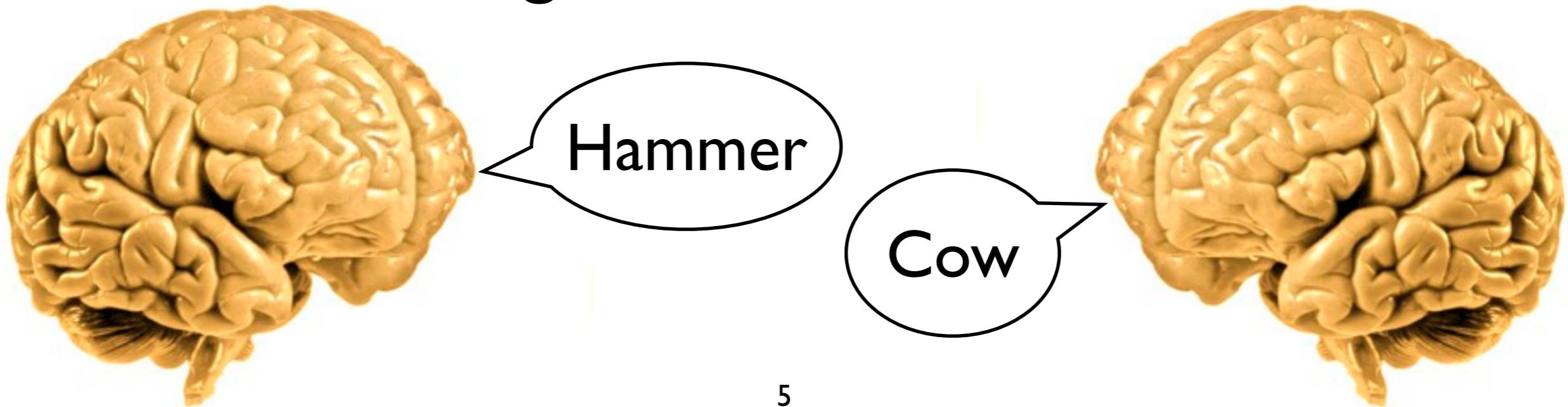
- Seismic data
- Weather sensors
- Physiological recordings
 - EKG, Ecog, MEG

Understanding Cognition

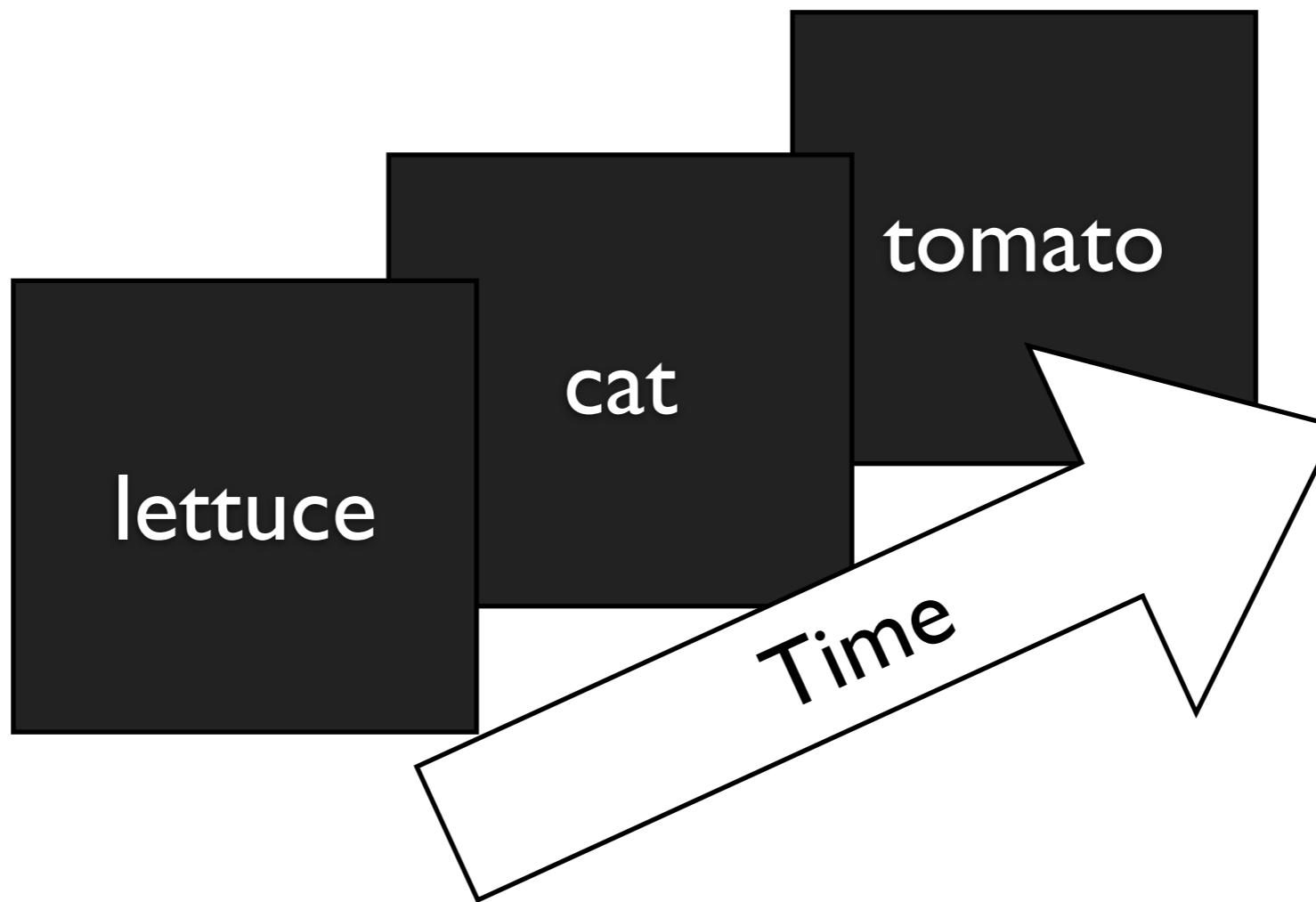
- How does the brain code concepts?
 - e.g. animals, food...
- How does brain activity change over time?

Understanding Cognition

- Show words to subjects
- Record & model brain activity
- Distinguish between word categories
 - e.g. animals vs tools?



Understanding Cognition



Challenging Data

- High noise & high dimension
- Need to predict with single (noisy) trials

Data Collection

- MEG recordings
- 2 subjects
- 20 words

Data Collection

- 4 word categories



Animals



Food



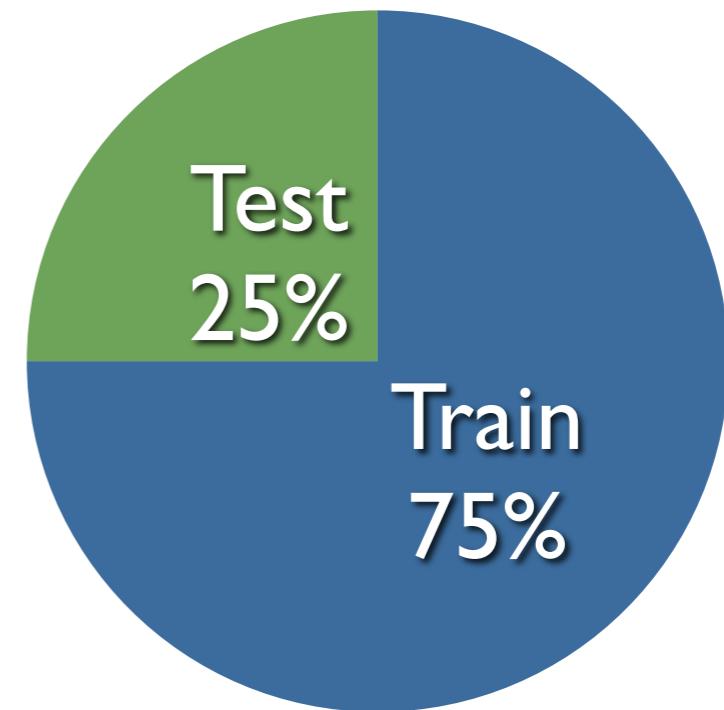
Tools



Buildings

Data Collection

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

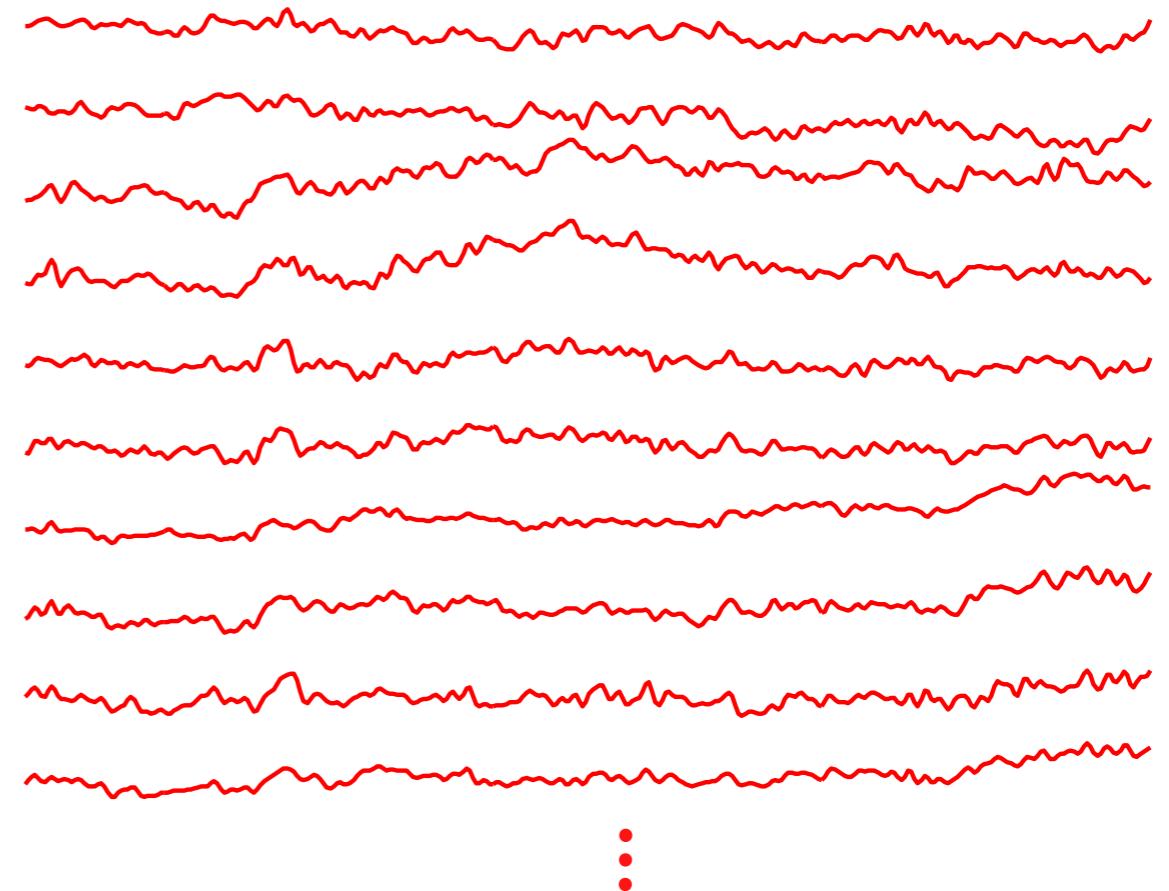


Brain Imaging

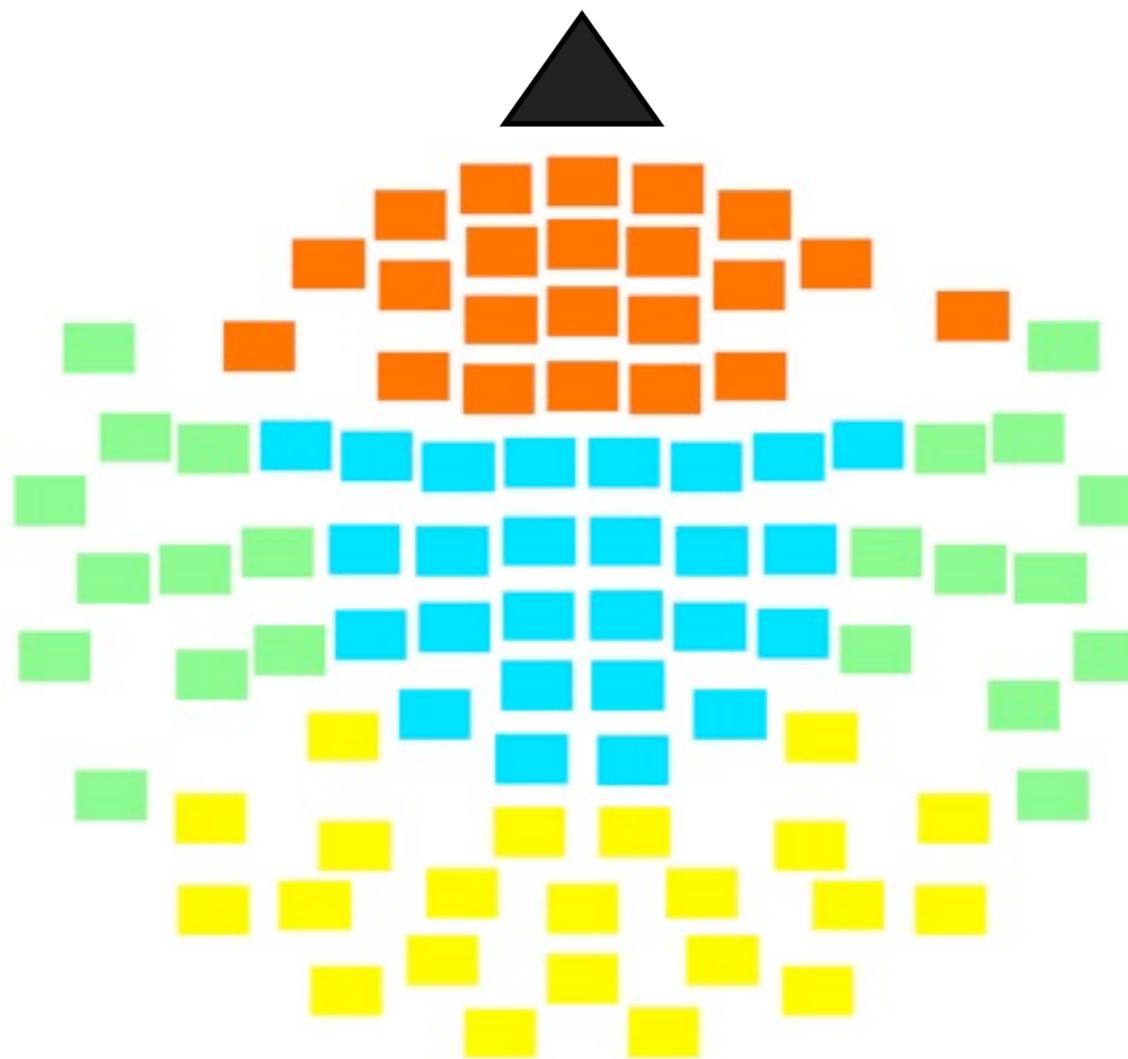
- fMRI
 - Measures blood oxygenation
 - Great spatial resolution
 - Poor time resolution
- MEG
 - Measures neuron's magnetic field
 - Great temporal resolution
 - Reduced spatial resolution

MEG machine

- 102 sensors in a helmet

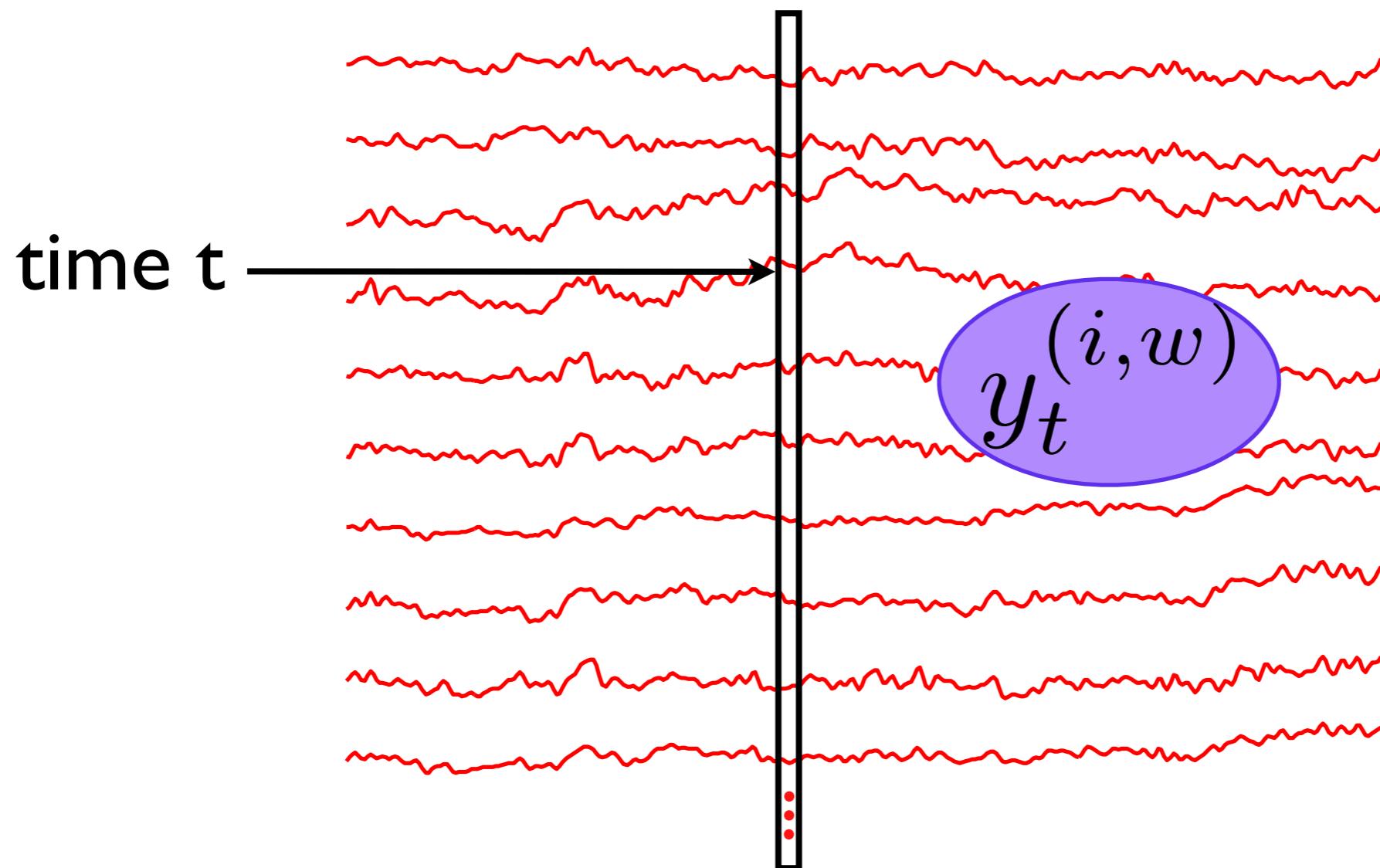


MEG Machine

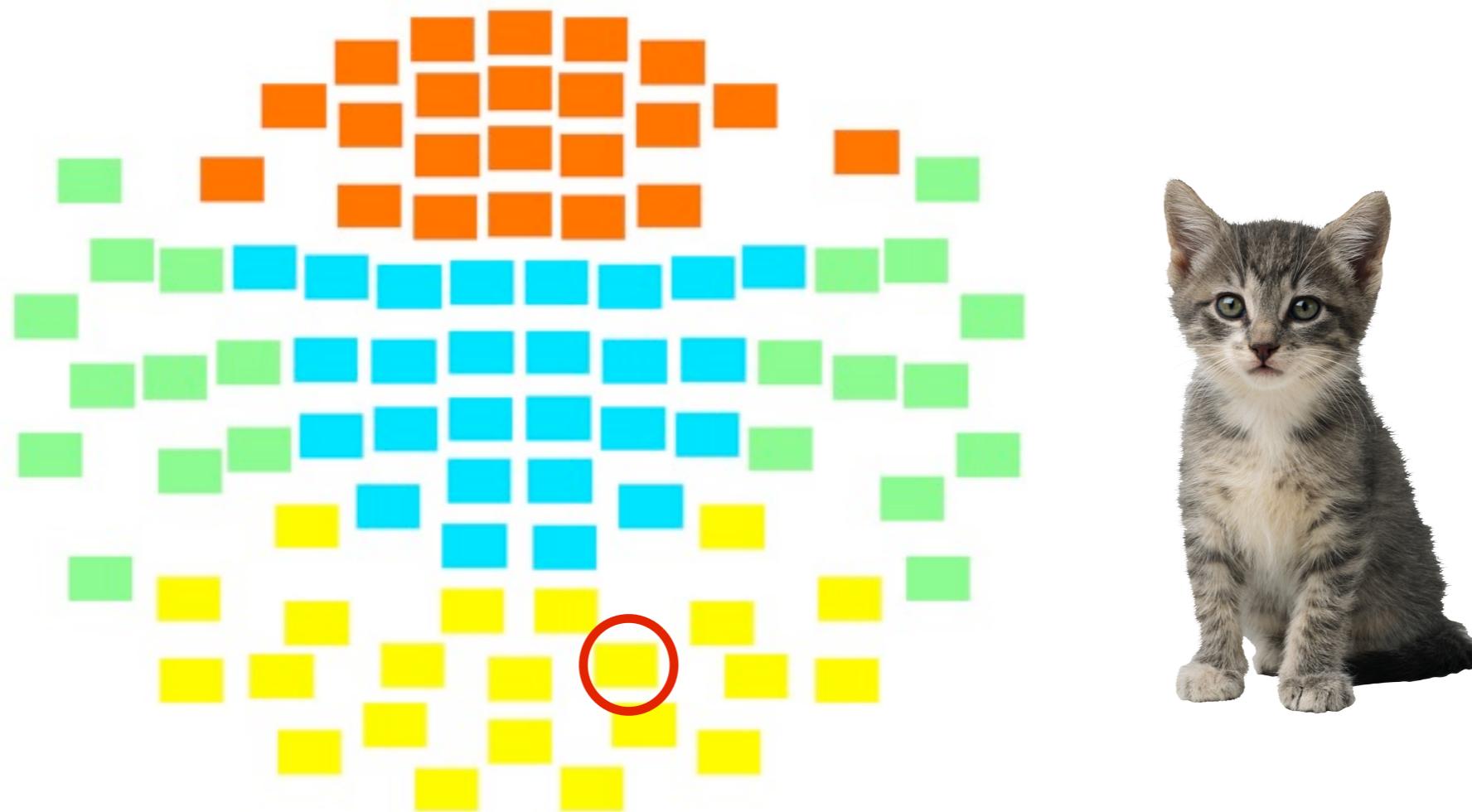


Terminology

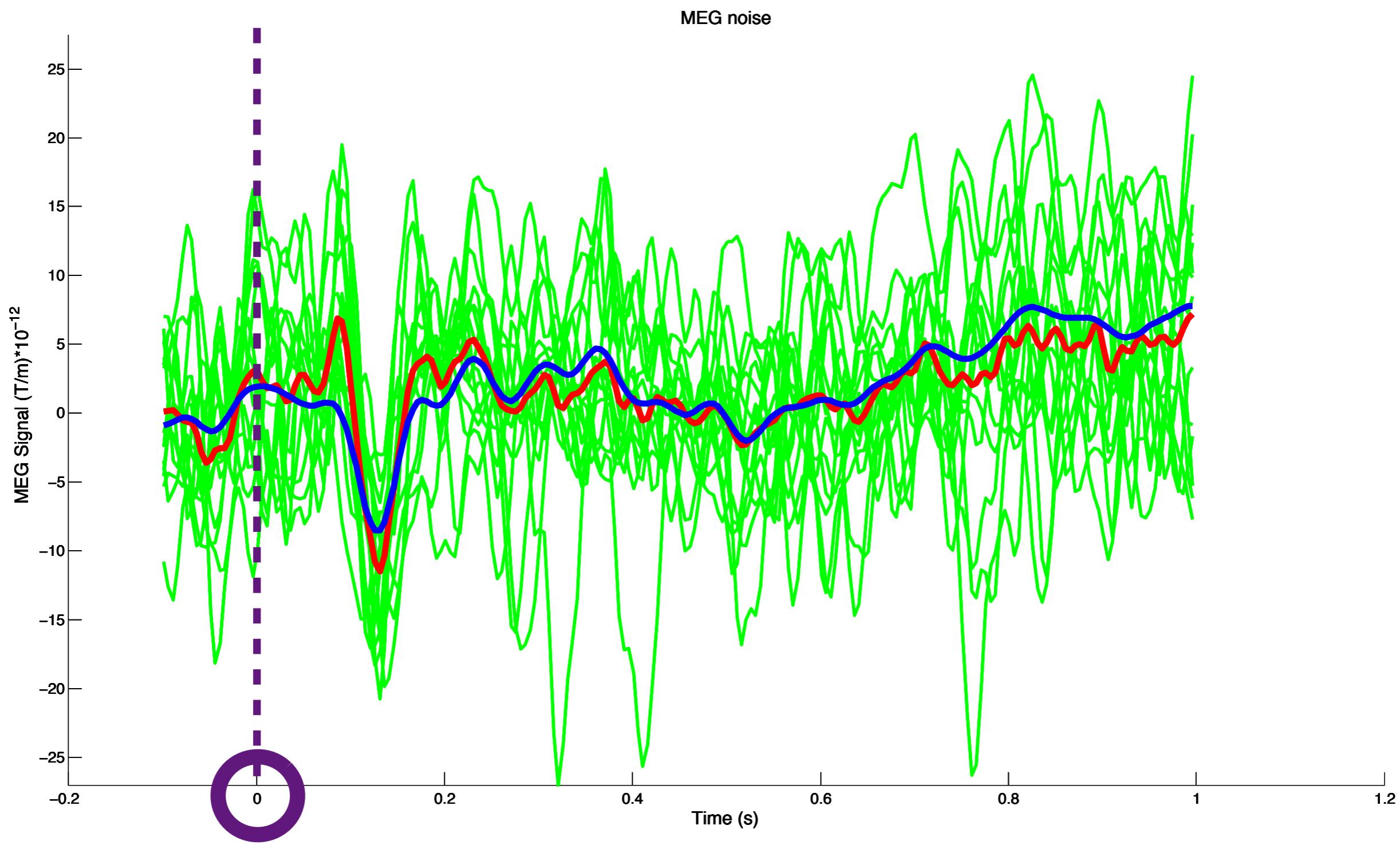
MEG recording for word w , trial i



MEG Noise



MEG Noise



Challenges

- 1) Noise
- 2) High dimension
- 3) Covariance changes with time
 - Functional connectivity
- 4) Trial-to-trial variation

Challenge I: **MEG data is noisy**

MEG data is noisy

- Traditionally:
 - Average repetitions (MLE)
 - assume channels are independent
- Can we do better than average?

MEG data is noisy

Observation: underlying neural signal is smooth

Underlying signal is smooth

- Model with a Gaussian process (GP)
 - $f \sim GP(m, c)$
 - m : mean function
 - c : kernel function (covariance)
 - squared exponential

$$c_i(t, t') = d_i \exp(-\kappa \|t - t'\|_2^2)$$

Underlying signal is smooth

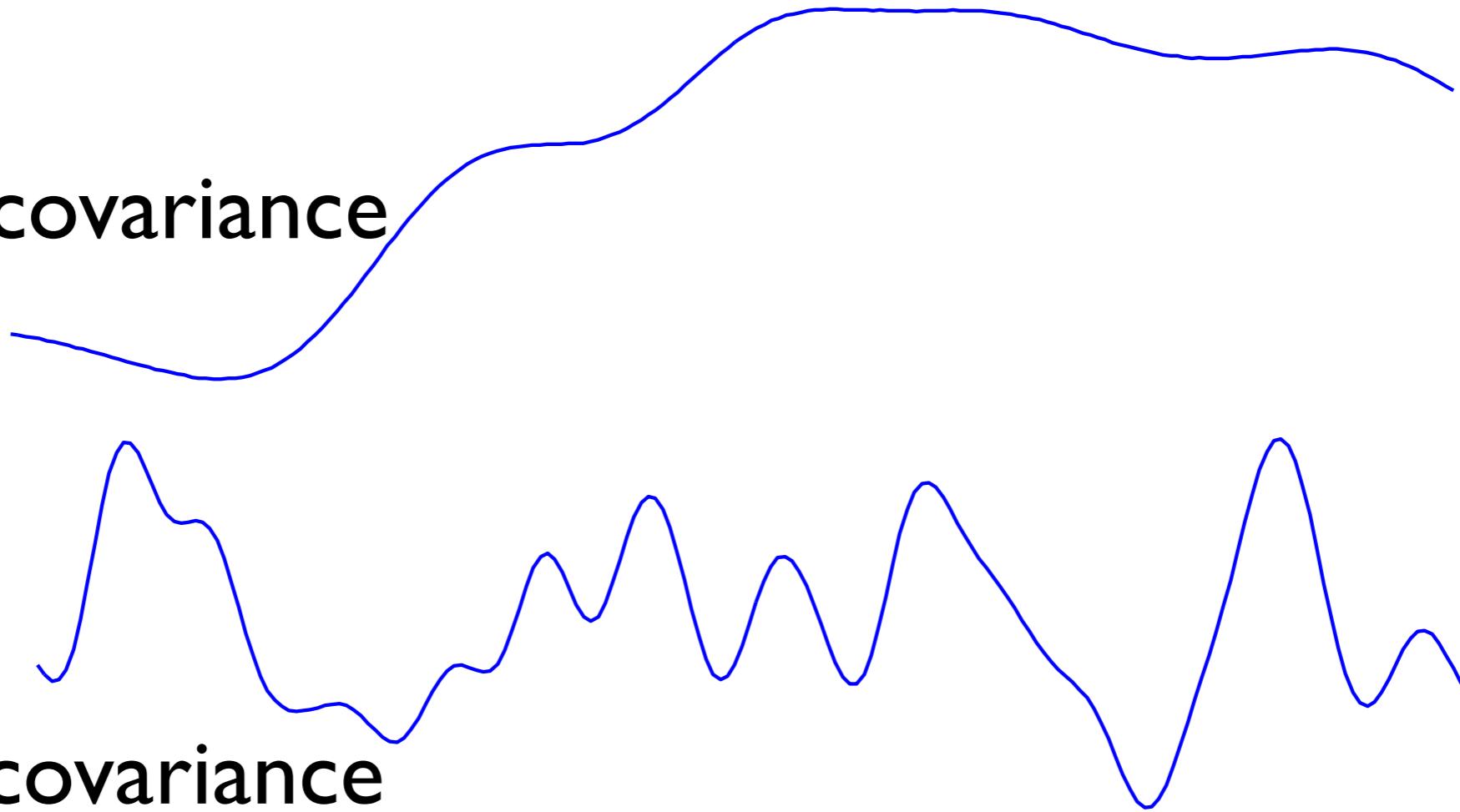
- $P(\mathbf{f}(t_1), \dots, \mathbf{f}(t_n)) \sim N_n(\mu, K)$
- $\mu = [\mathbf{m}(t_1), \dots, \mathbf{m}(t_n)]$
- $K_{ij} = c(t_i, t_j)$

Gaussian Processes

c: kernel function (covariance)

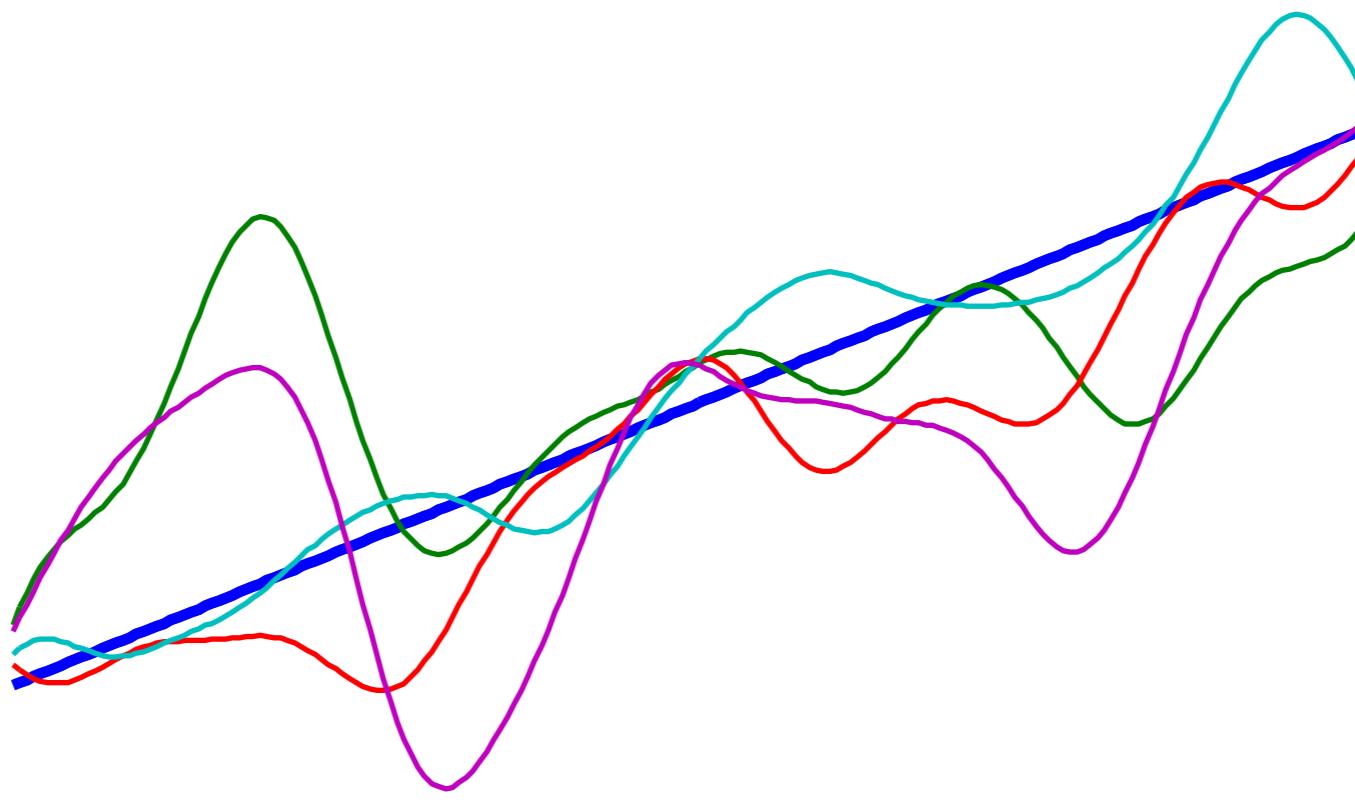
High covariance

Low covariance



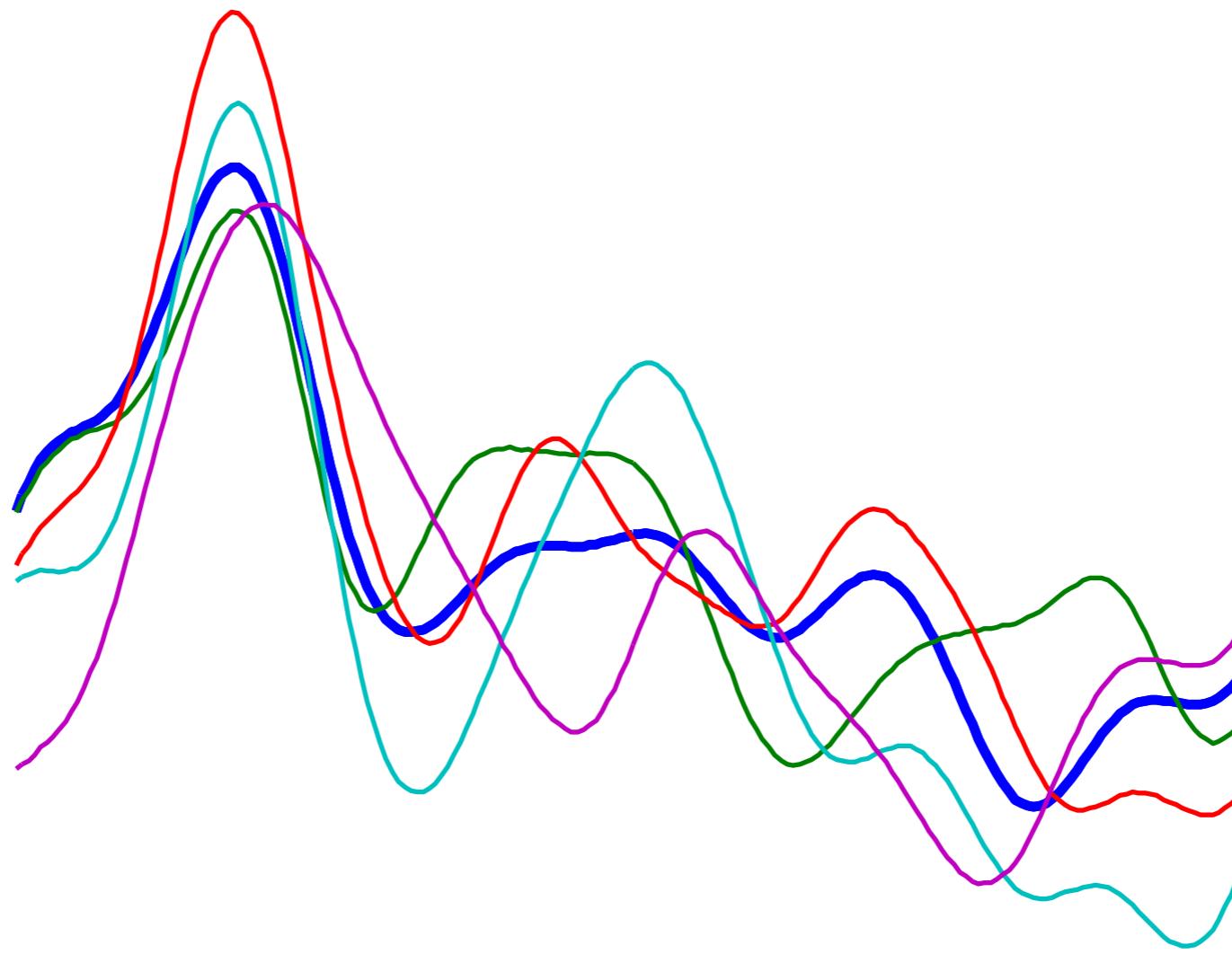
Gaussian Processes

m: mean function

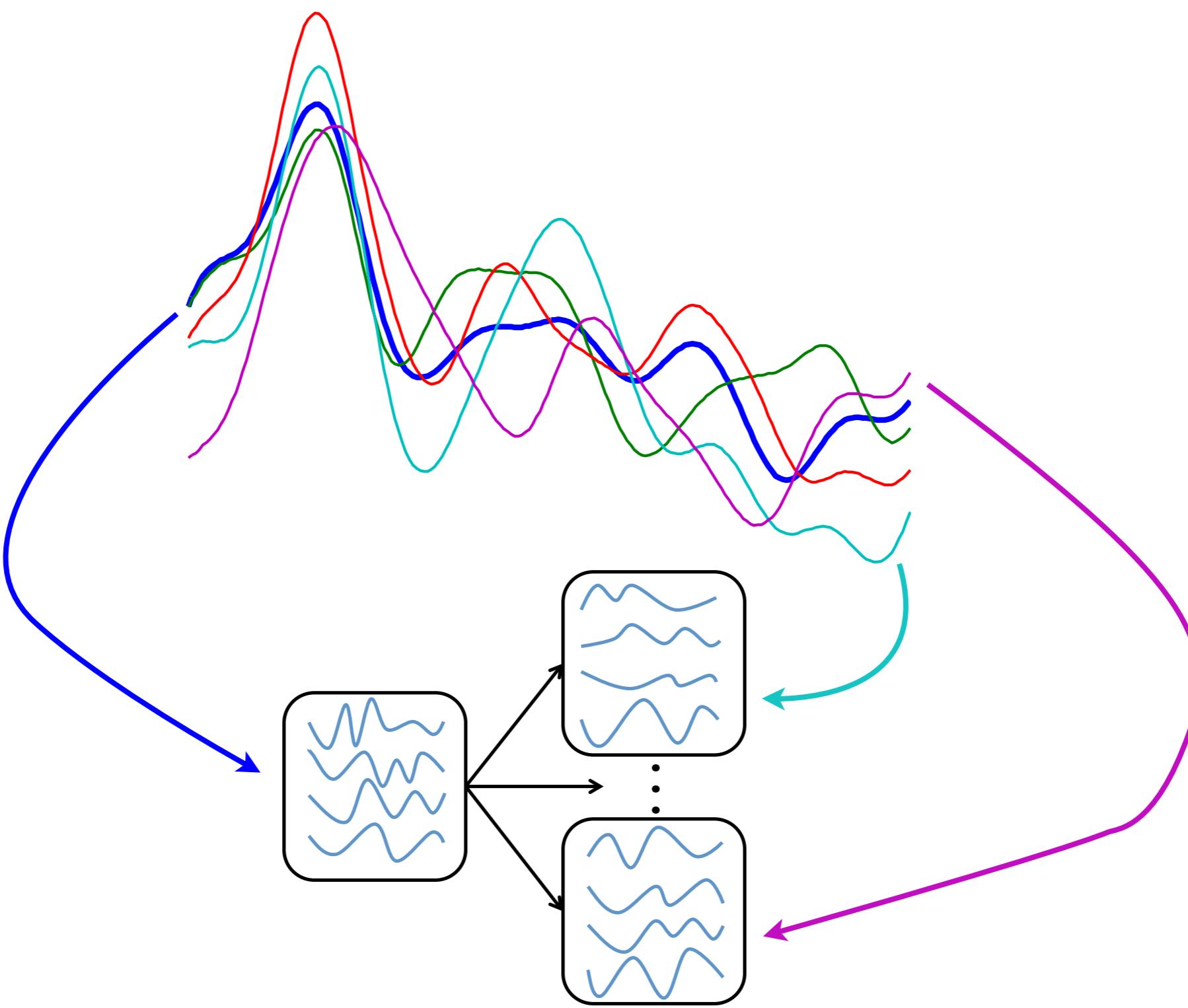


Gaussian Processes

m: mean function



Gaussian Processes



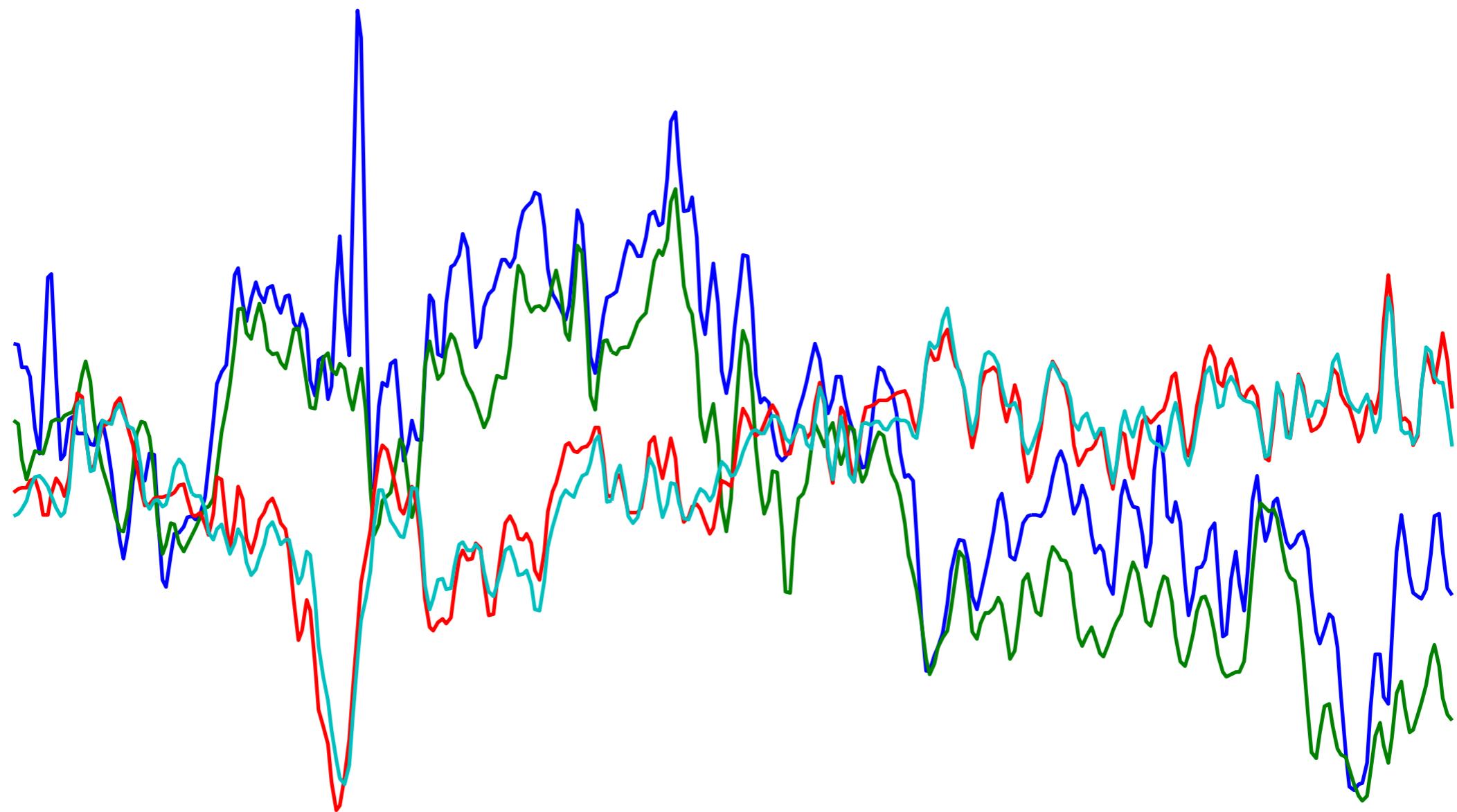
Challenge 2: Data is of high dimension

Data is of high dimension

Observation: Sensors are redundant

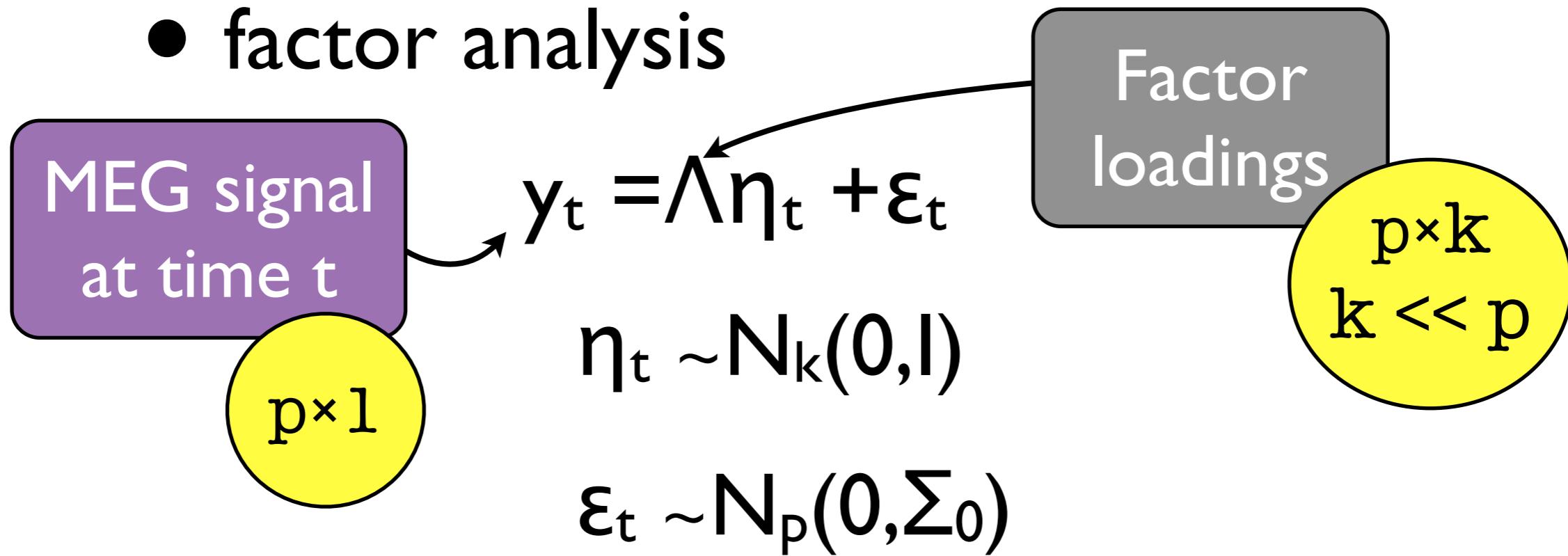
- spatially localized correlation
- long-range correlation

Redundant Sensors

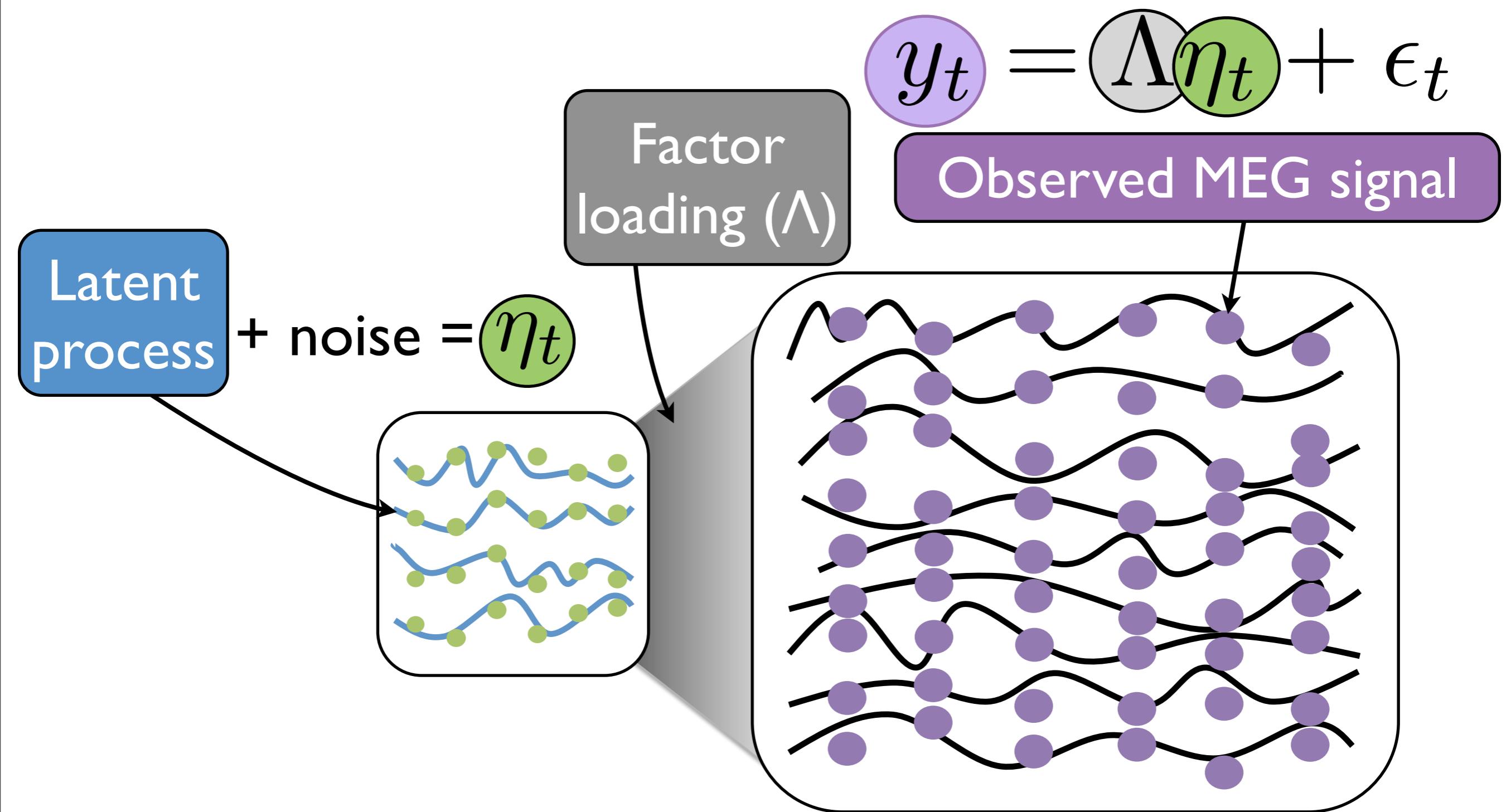


Redundant Sensors

- Model in a lower dimensional space
 - factor analysis

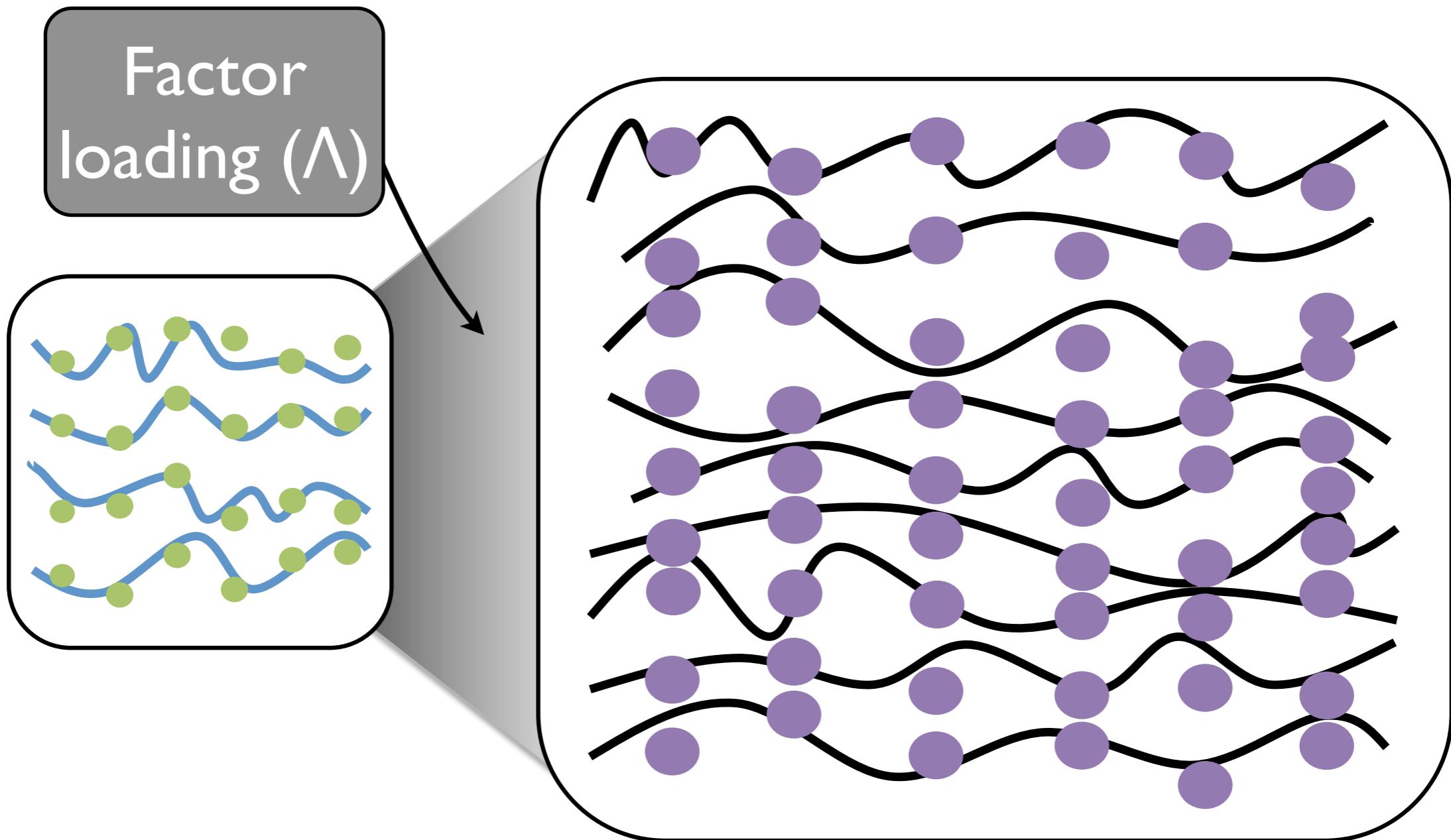


Redundant Sensors

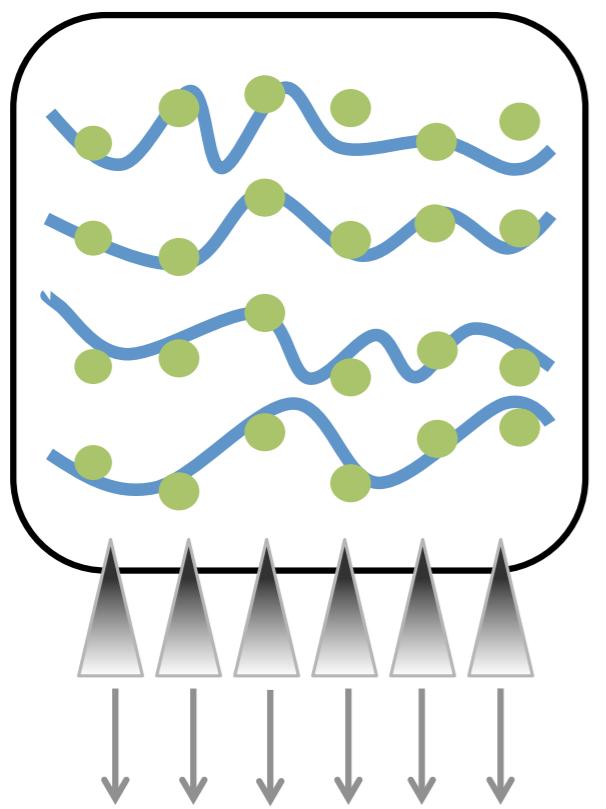
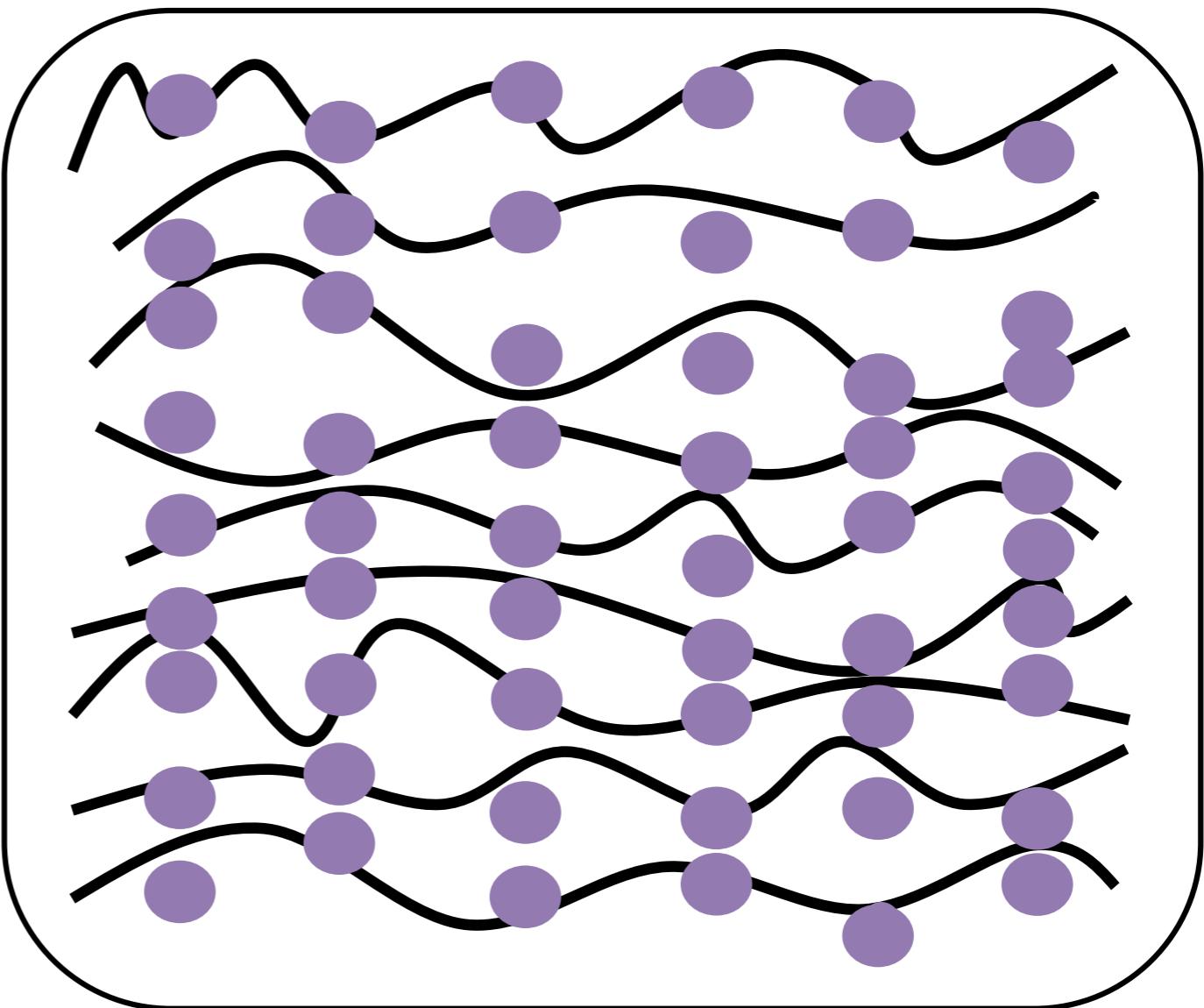


Challenge 3: Coordination changes with time (Heteroskedasticity)

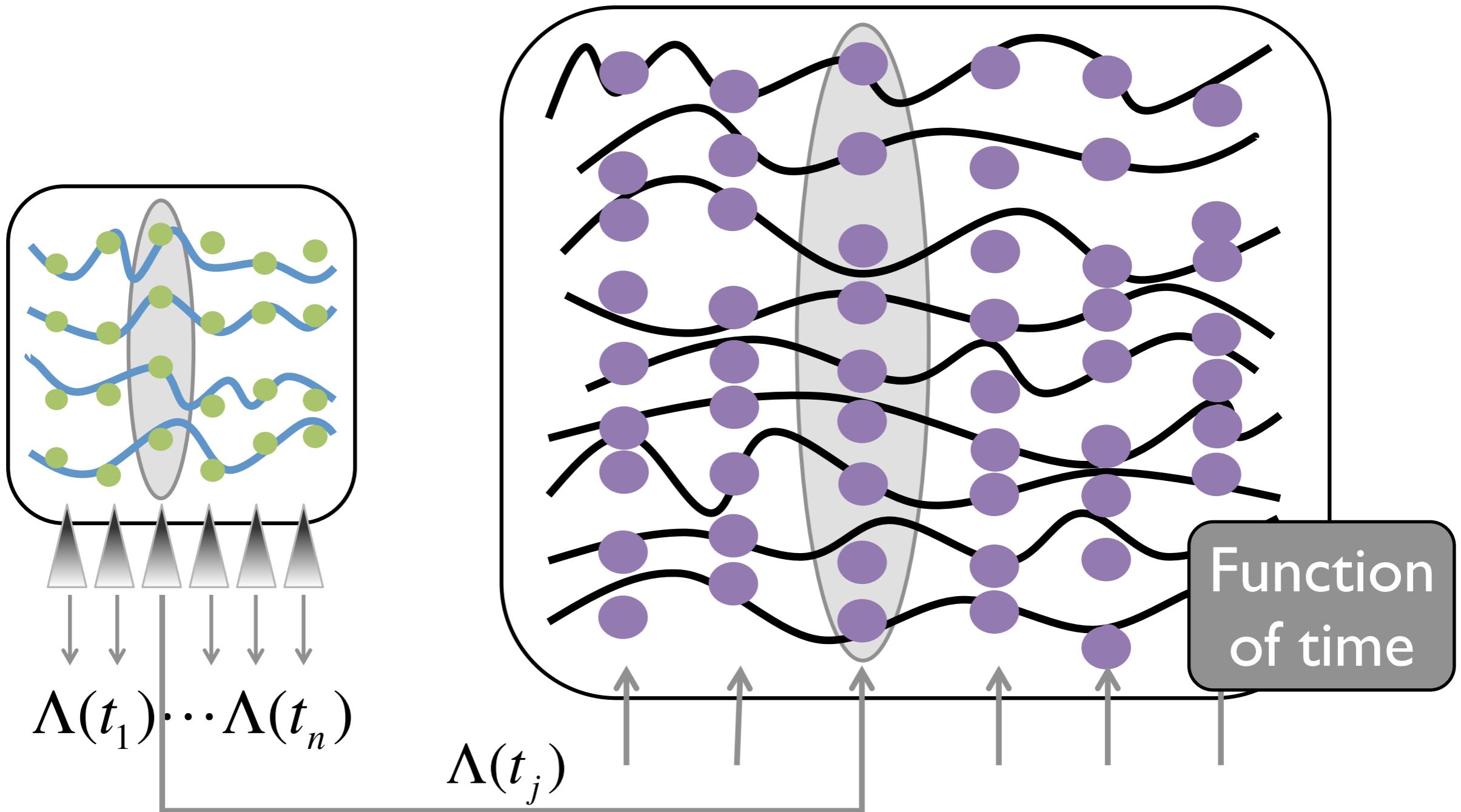
Coordination changes with time



Coordination changes with time


$$\Lambda(t_1) \cdots \Lambda(t_n)$$


Coordination changes with time



Coordination changes with time

MEG
signal

$$y_t^{(i,w)} = \Lambda^{(w)}(t) \eta_t^{(i,w)} + \epsilon_t^{(i,w)}, \quad \epsilon_t^{(i,w)} \sim N_p(0, \Sigma_0^{(w)})$$

$$\eta_t^{(i,w)} = \psi^{(i,w)}(t) + \nu_t^{(i,w)}, \quad \nu_t^{(i,w)} \sim N_k(0, I_k)$$

$$\Lambda^{(w)}(t) = \Theta^{(w)} \xi^{(w)}(t)$$

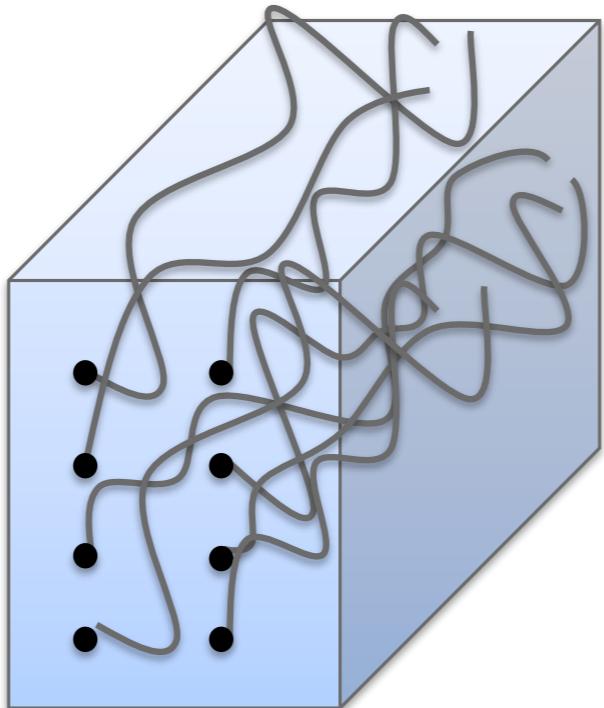
Factor
loadings

Latent
mean

Coordination changes with time

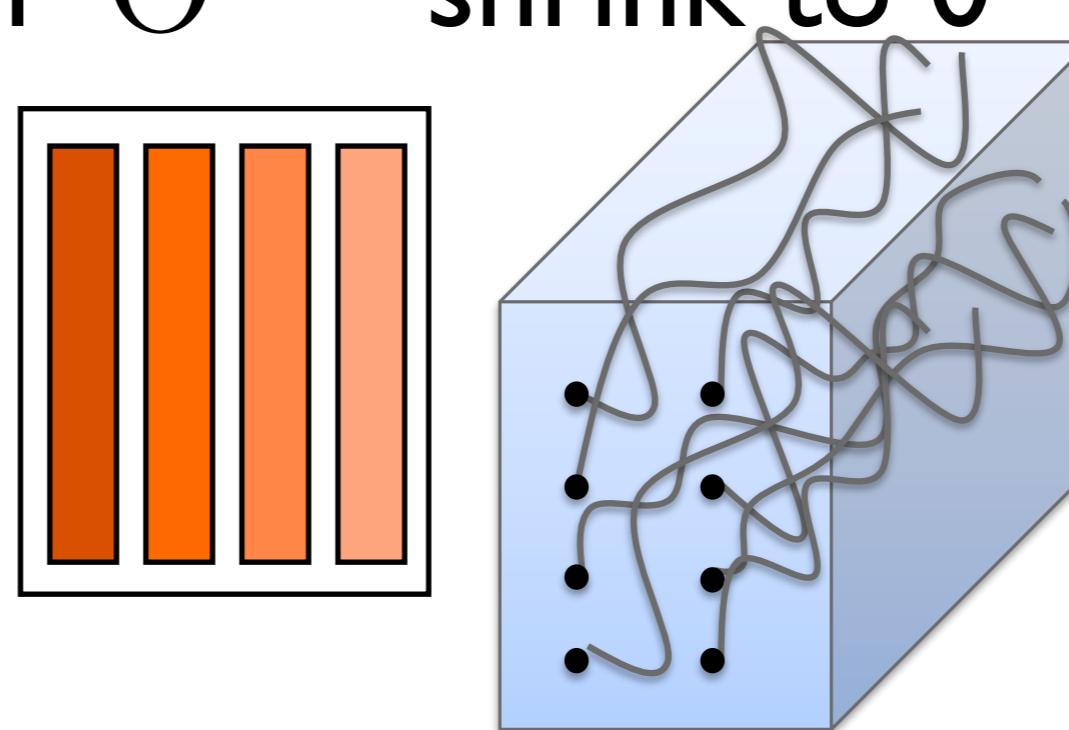
$$\begin{aligned}\Sigma^{(w)}(t) &= \Lambda^{(w)}(t)\Lambda^{(w)\prime}(t) \\ &= \Theta^{(w)} \boxed{\xi^{(w)}(t)} \boxed{\xi^{(w)}(t)'} \Theta^{(w)\prime}\end{aligned}$$

Heteroskedastic



Shrinkage Prior

- $\Lambda^{(w)}(t) = \Theta^{(w)}\xi^{(w)}(t)$
- $\xi^{(w)}(t)$ is a matrix $L \times K$
- Columns of $\Theta^{(w)}$ shrink to 0



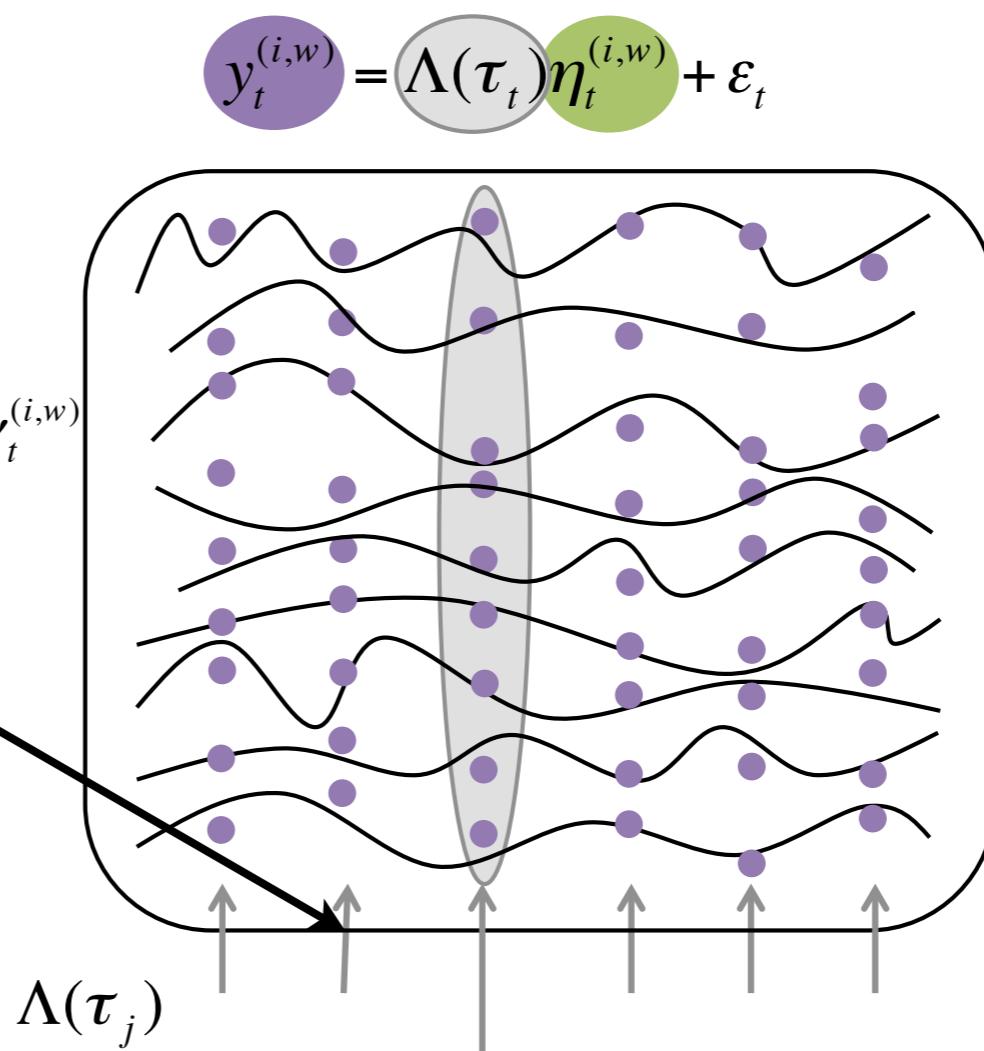
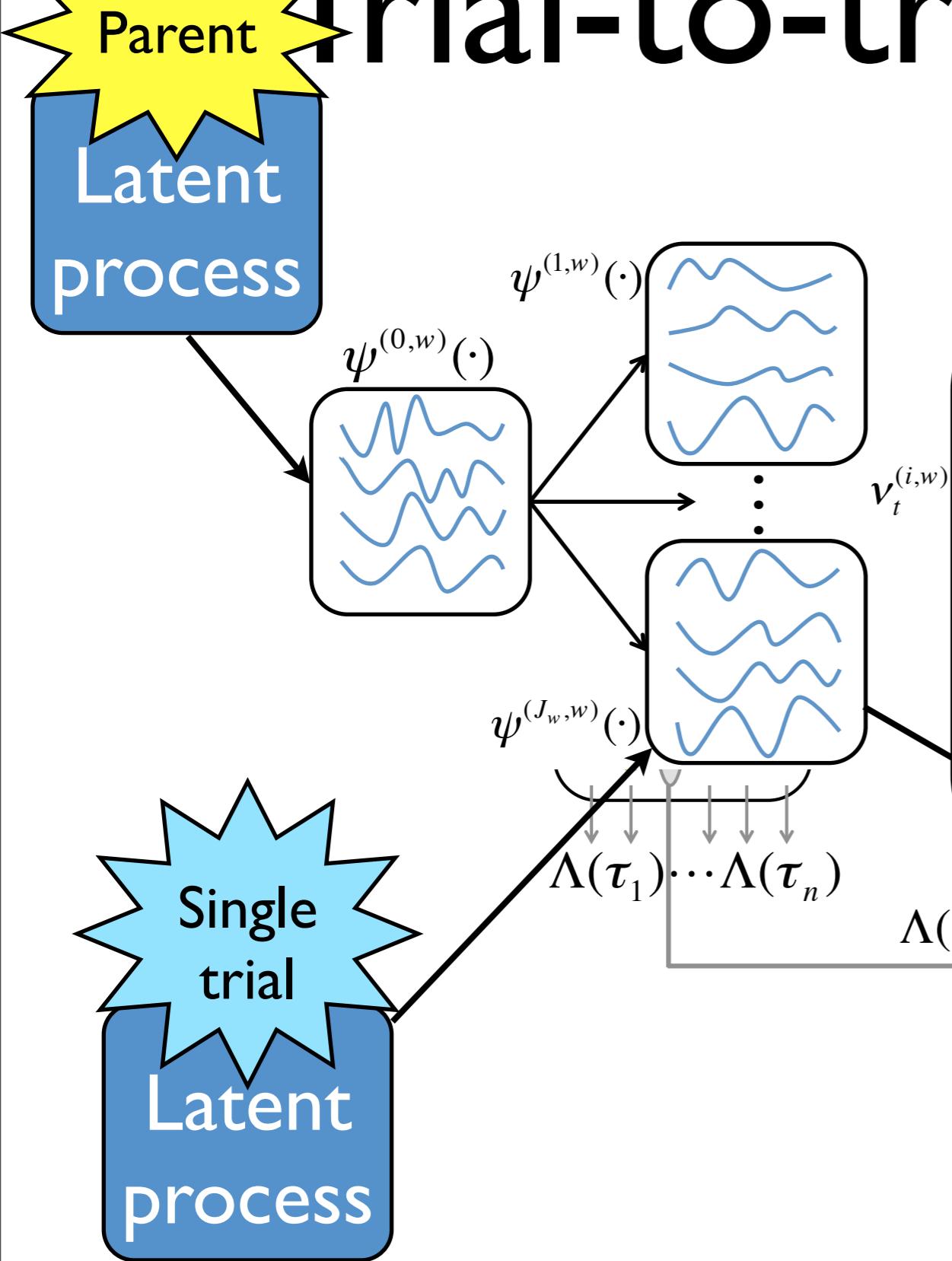
Challenge 4: Trial-to-trial variability

Trial-to-trial variability

Observation: each trial is a (noisy) recording of the **same process**

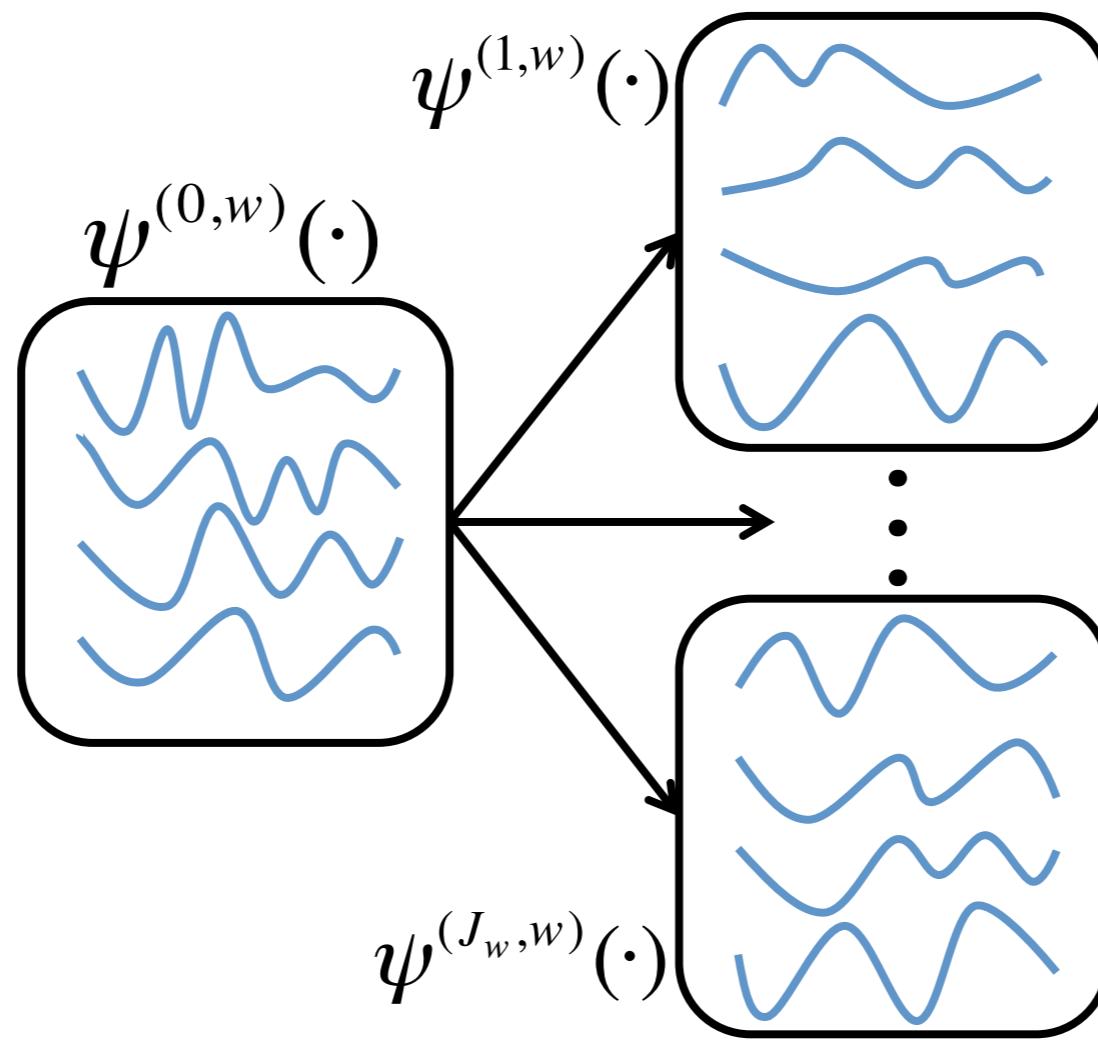
- Allow for multiple examples
- Each example **may deviate** somewhat

Trial-to-trial variability

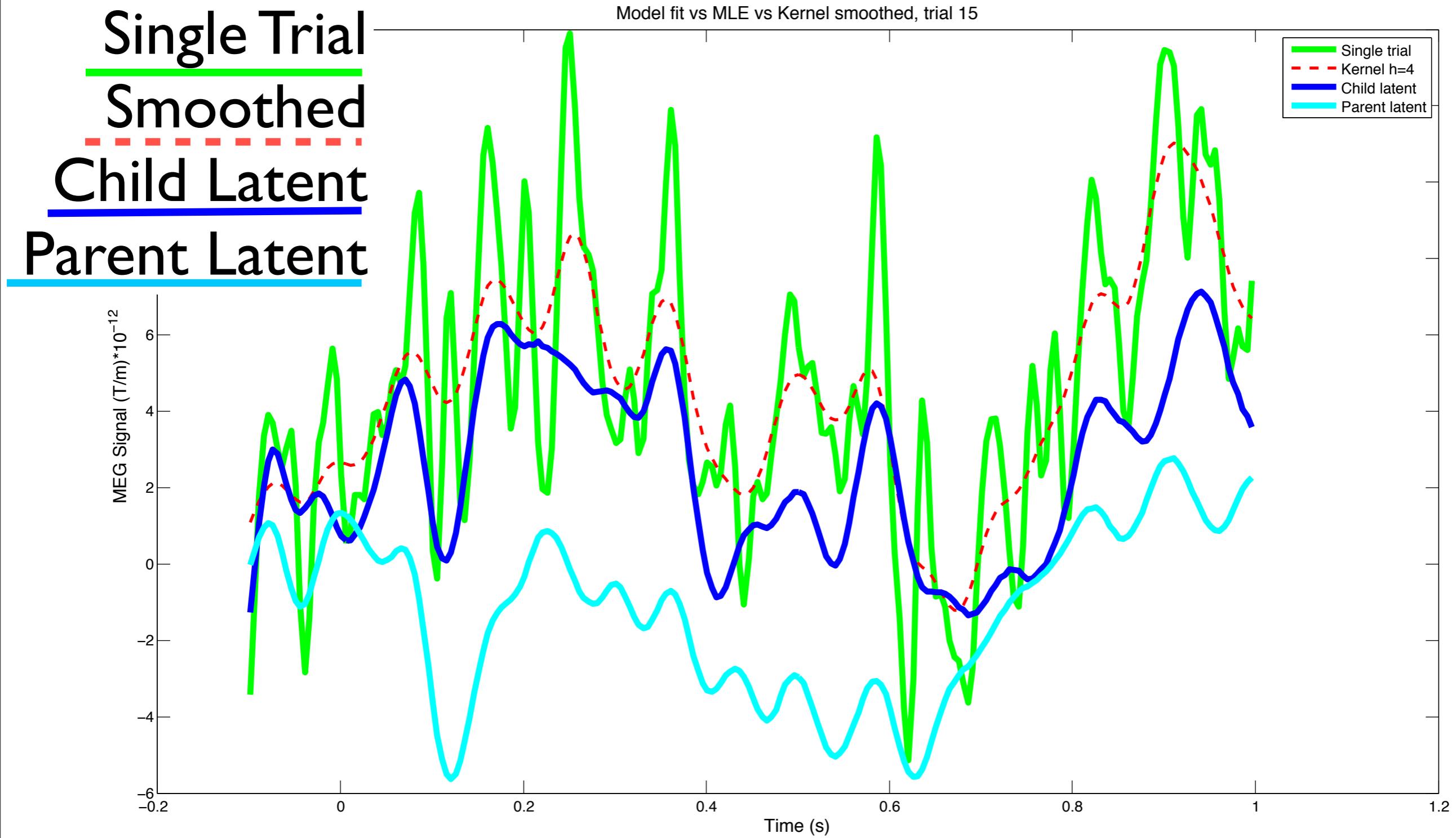


Trial-to-trial variability

- Hierarchy **shares information** between trials

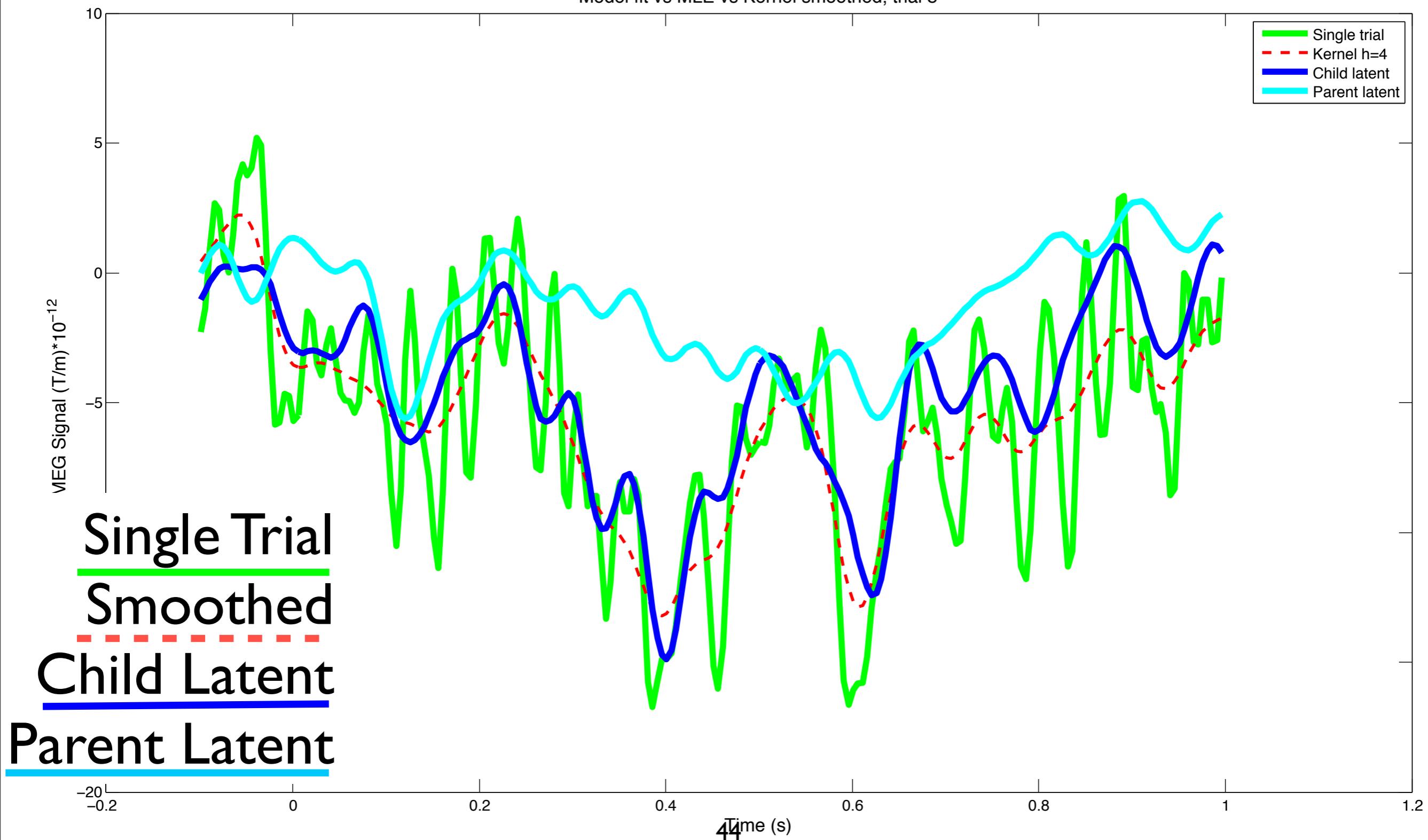


Trial-to-trial variability



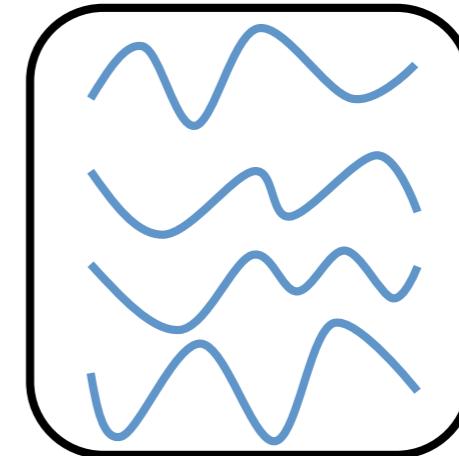
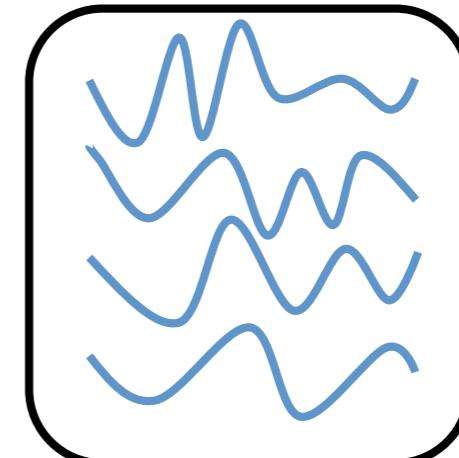
Trial-to-trial variability

Model fit vs MLE vs Kernel smoothed, trial 8

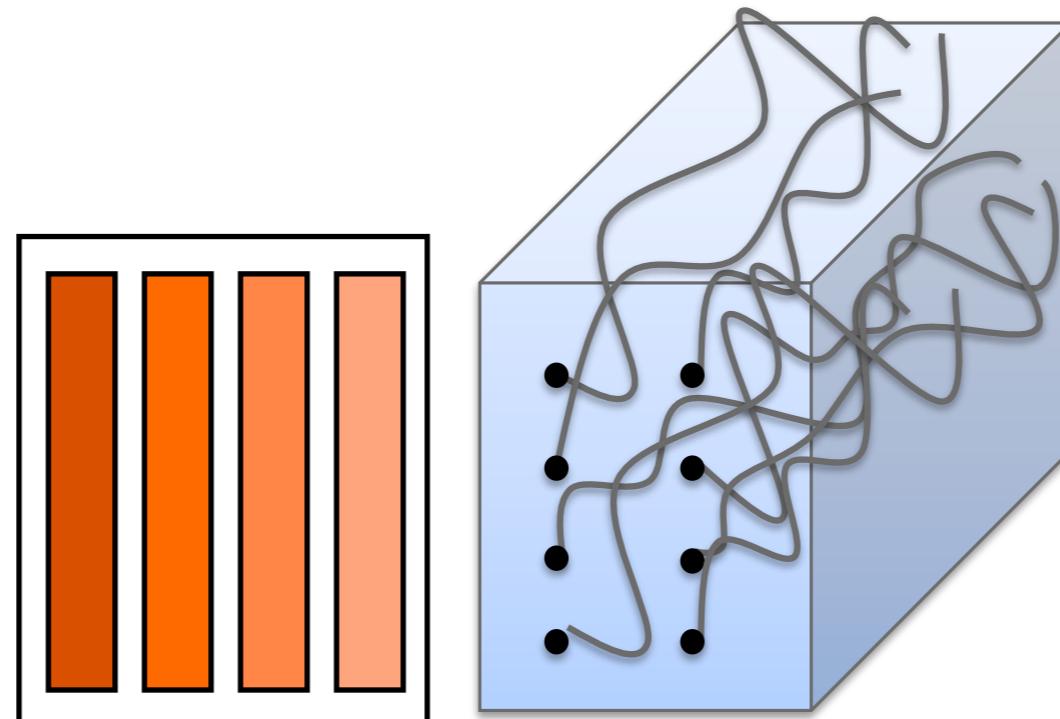


Generative Story

lettuce

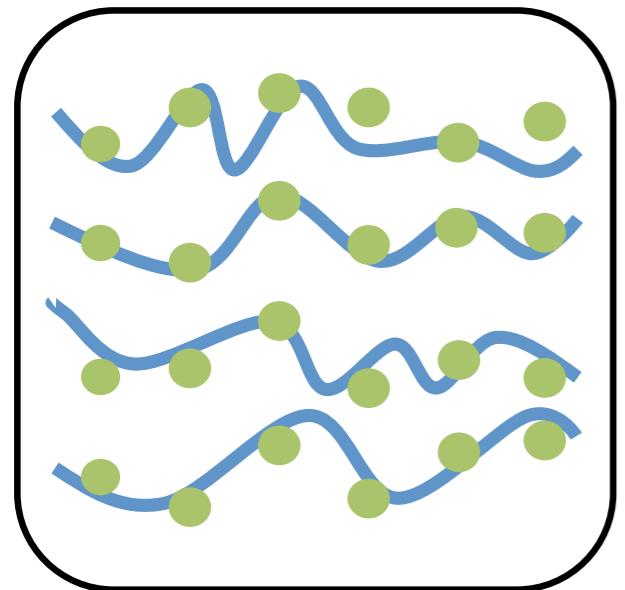


Generative Story

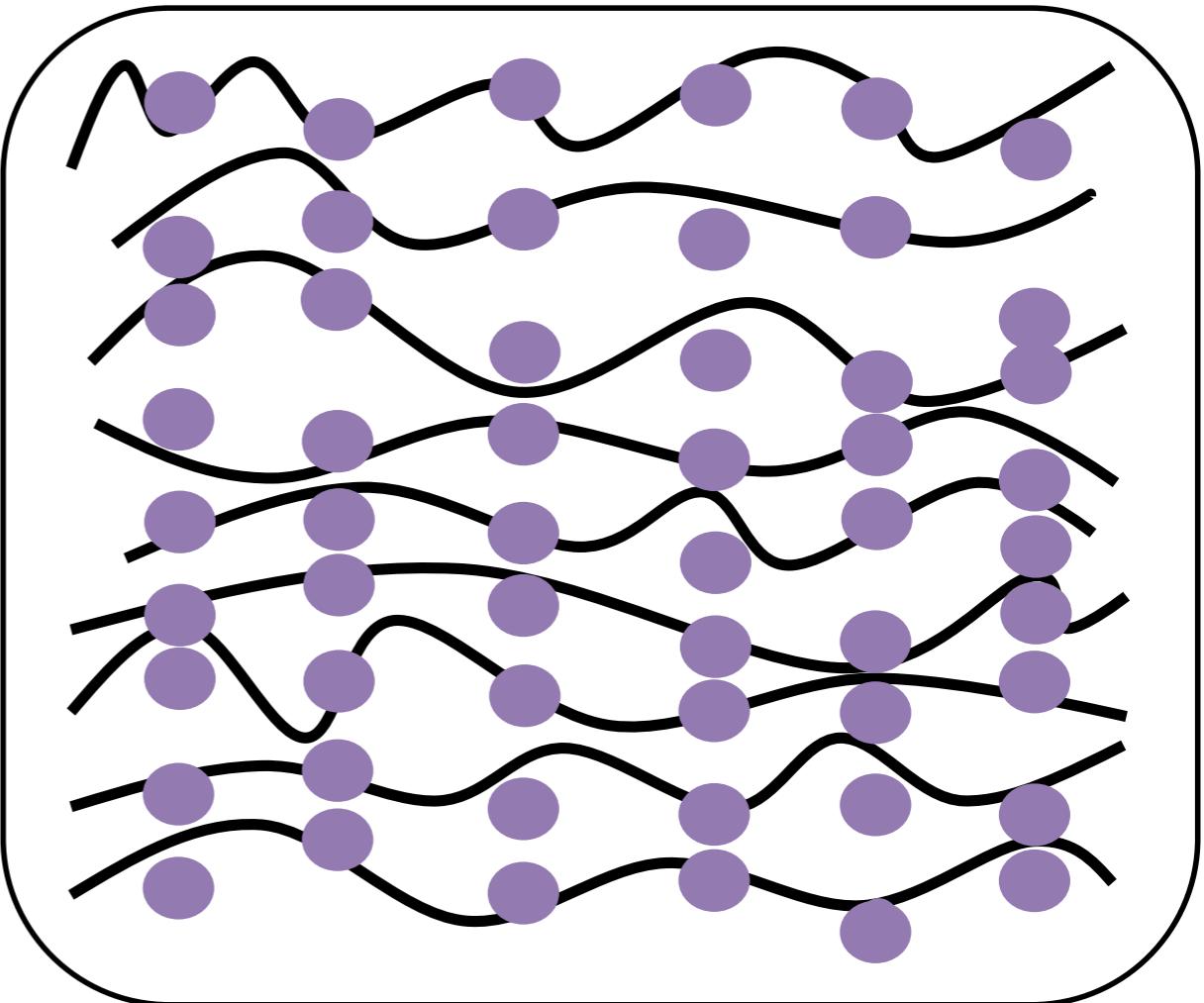


$$\Lambda^{(w)}(t) = \Theta^{(w)} \xi^{(w)}(t)$$

Generative Story



$$\Lambda^{(w)}(t) =$$

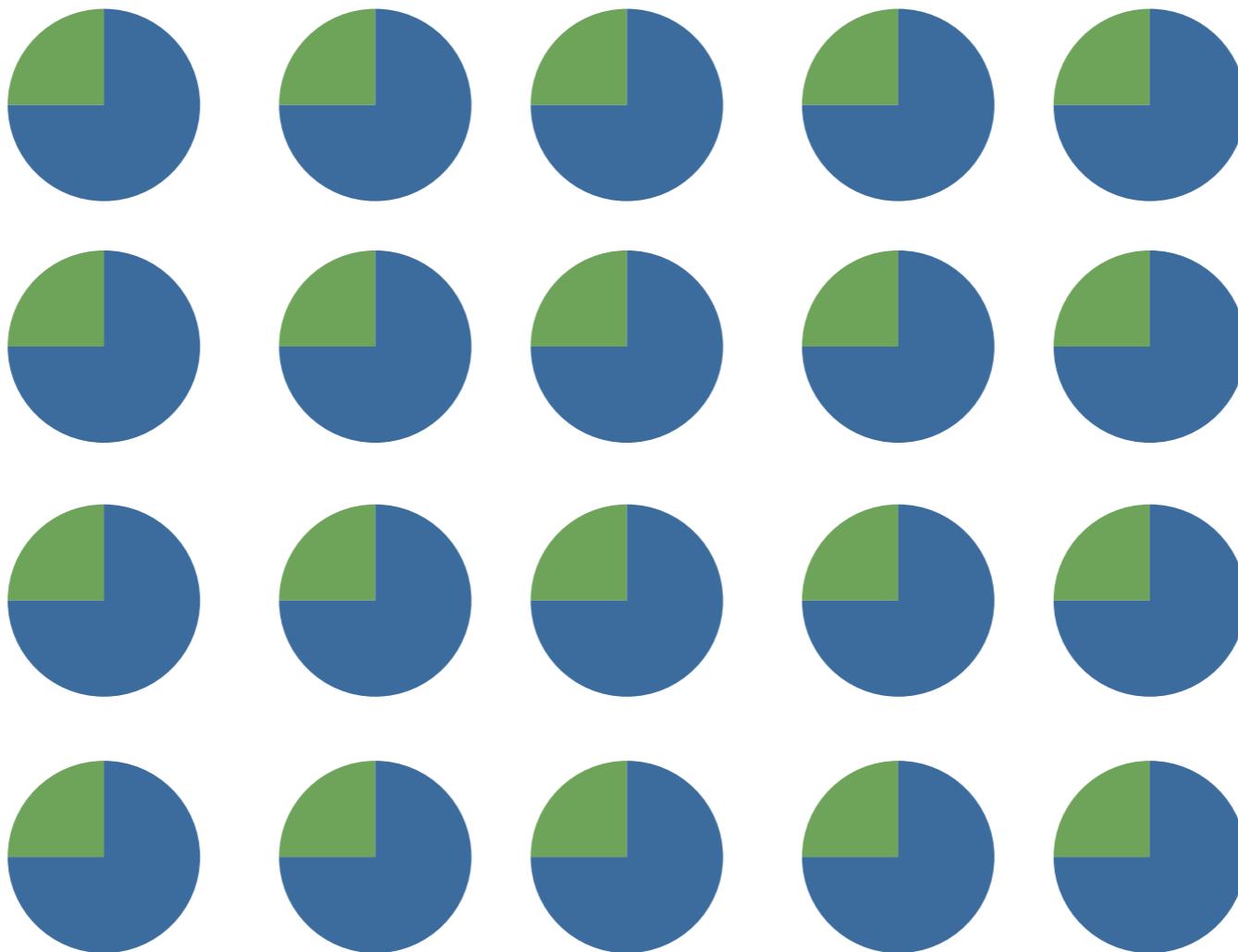


Evaluation

Recall the task

- Distinguish between 4 word categories
- Train on 15 single trials/word
 - 20 words
- Evaluate on 100 single trials
 - 5/word

Train & Test



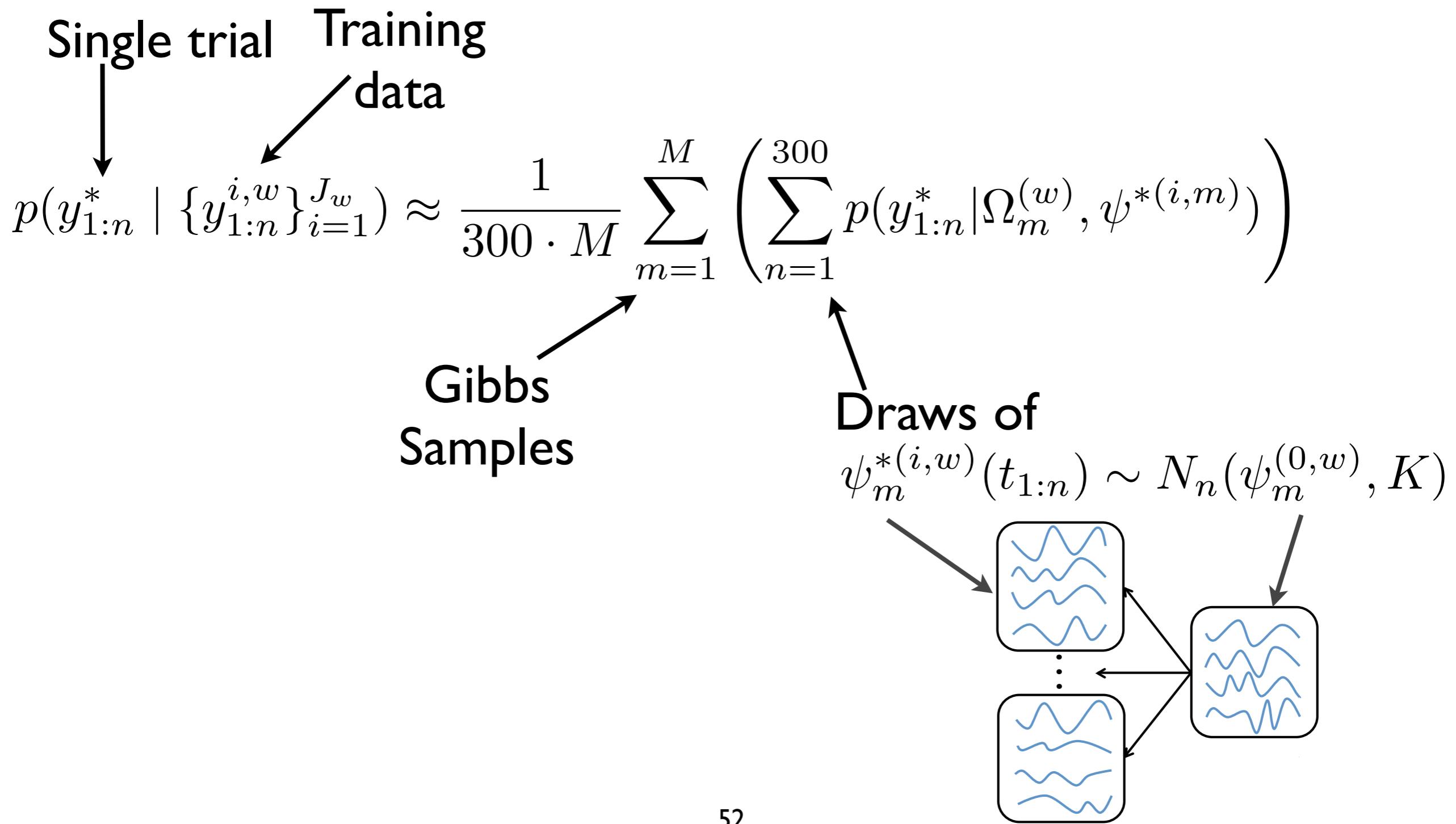
Fit Model

```

for  $i \leftarrow 1, J_w$  do
    for  $\ell \leftarrow 1, L$  do
         $\tilde{\Sigma}_{\psi}^{(w)} = K^{-1} + \text{diag}([\Lambda^{(w)}(\tau_1)]'_{.\ell} \Sigma^{-(w)}(\tau_1) [\Lambda^{(w)}(\tau_1)]_{.\ell}, \dots, [\Lambda^{(w)}(\tau_n)]'_{.\ell} \Sigma^{-(w)}(\tau_n) [\Lambda^{(w)}(\tau_n)]_{.\ell})$ 
         $\tilde{y}_t^{(i,w)} = y_t^{(i,w)} - \sum_{r \neq \ell} [\Lambda^{(w)}(\tau)]_{.r} \psi_r^{(i,w)}(\tau)$ 
         $\psi_{\ell}^{(i,w)}(\tau_{1:n}) \sim N_n^{-1} \left( \begin{bmatrix} [\Lambda^{(w)}(\tau_1)]'_{.\ell} \Sigma^{-(w)}(\tau_1) \tilde{y}_1^{(i,w)} + \psi_1^{(0,w)} \\ \vdots \\ [\Lambda^{(w)}(\tau_n)]'_{.\ell} \Sigma^{-(w)}(\tau_n) \tilde{y}_n^{(i,w)} + \psi_n^{(0,w)} \end{bmatrix}, \tilde{\Sigma}_{\psi}^{(w)} \right)$ 
    for  $t \leftarrow 1, n$  do
         $\tilde{y}_t^{(i,w)} = y_t^{(i,w)} - \Lambda^{(w)}(\tau) \psi^{(i,w)}(\tau)$ 
         $\nu_t^{(i,w)} \sim N_k^{-1} \left( \Lambda^{(w)'}(\tau) \Sigma_0^{-(w)} \tilde{y}_t^{(i,w)}, I + \Lambda^{(w)'}(\tau) \Sigma_0^{-(w)} \Lambda^{(w)}(\tau) \right)$ 
    for  $\ell \leftarrow 1, L$  do
         $\psi_{\ell}^{(0,w)}(\tau_{1:n}) \sim N_n^{-1} (K_1^{-1} \sum_{i=1}^{J_w} \psi_{\ell}^{(i,w)}(\tau_{1:n}), K_0^{-1} + J_w K_1^{-1})$ 
    for  $\ell \leftarrow 1, L$  do
        for  $m \leftarrow 1, K$  do
             $\tilde{\Sigma}_{\xi} = K_1^{-1} + \sum_{i=1}^{J_w} \text{diag} \left( \left( \eta_{1,m}^{(i,w)} \right)^2 \sum_{j=1}^p \left( \Theta_{j\ell}^{(w)} \right)^2 \sigma_{j,w}^{-2}, \dots, \left( \eta_{n,m}^{(i,w)} \right)^2 \sum_{j=1}^p \left( \Theta_{j\ell}^{(w)} \right)^2 \sigma_{j,w}^{-2} \right)$ 
             $\tilde{y}_{t,j}^{(i,w)} = y_{t,j}^{(i,w)} - \sum_{(r,s) \neq (\ell,m)} \Theta_{jr}^{(w)} \xi_{rs}^{(w)}(\tau).$ 
             $\xi_{\ell m}^{(w)}(\tau_{1:n}) \sim N_n^{-1} \left( \sum_{i=1}^{J_w} \begin{bmatrix} \eta_{1,m}^{(i,w)} \sum_{j=1}^p \Theta_{j\ell}^{(w)} \sigma_{j,w}^{-2} \tilde{y}_{1,j}^{(i,w)} \\ \vdots \\ \eta_{n,m}^{(i,w)} \sum_{j=1}^p \Theta_{j\ell}^{(w)} \sigma_{j,w}^{-2} \tilde{y}_{n,j}^{(i,w)} \end{bmatrix}, \tilde{\Sigma}_{\xi} \right)$ 
        for  $j \leftarrow 1, p$  do
             $\Theta_{j\cdot}^{(w)} = \begin{bmatrix} \Theta_{j1}^{(w)} & \dots & \Theta_{jL}^{(w)} \end{bmatrix}, \eta_t^{(i,w)} = \psi^{(i,w)}(\tau) + \nu_t^{(i,w)}$ 
             $\sigma_{w,j}^{-2} \sim \text{Ga} \left( a_{\sigma} + \frac{n J_w}{2}, b_{\sigma} + \frac{1}{2} \sum_{i=1}^{J_w} \sum_{t=1}^n (y_{t,j}^{(i,w)} - \Theta_{j\cdot}^{(w)} \xi^{(w)}(\tau) \eta_t^{(i,w)})^2 \right)$ 
    for  $j \leftarrow 1, p$  do
         $y_{\cdot,j} = [y_{(1,j)}(t_{1:n}) \dots y_{(J_w,j)}(t_{1:n})]'$ 
         $\tilde{\eta}^{(w)'} = \begin{bmatrix} \xi^{(w)}(\tau_1) \eta_1^{(1,w)} & \dots & \xi^{(w)}(\tau_n) \eta_n^{(1,w)} & \dots & \xi^{(w)}(\tau_1) \eta_1^{(J_w,w)} & \dots & \xi^{(w)}(\tau_n) \eta_n^{(J_w,w)} \end{bmatrix}$ 
         $\tilde{\Sigma}_{\Theta}^{-(w)} = \sigma_{w,j}^{-2} \tilde{\eta}^{(w)'} \tilde{\eta}^{(w)} + \text{diag}(\phi_{j1}^{(w)} \zeta_1^{(w)}, \dots, \phi_{jL}^{(w)} \zeta_L^{(w)})$ 
         $\Theta_{j\cdot}^{(w)} \sim N_L \left( \sigma_{w,j}^{-2} \tilde{\Sigma}_{\Theta}^{(w)} \tilde{\eta}^{(w)'} y_{\cdot,j}, \tilde{\Sigma}_{\Theta}^{(w)} \right)$ 

```

Evaluating Models



Testing

- MLE-based models (likelihood)
 - Homoskedastic $\mu_{MLE}^{(w)}$ $\hat{\Sigma}^{(w)}$
 - Heteroskedastic $\mu_{MLE}^{(w)}$ $\hat{\Sigma}_t^{(w)}$
 - Our model $\psi^*(t) \sim \text{GP}(\psi^{(0,w)}(t), K)$
$$\mu^*(t) = \Lambda^{(w)}(t)\psi^*(t)$$
$$\Sigma^{(w)}(t) = \Lambda^{(w)}(t)\Lambda^{(w)}(t)'$$

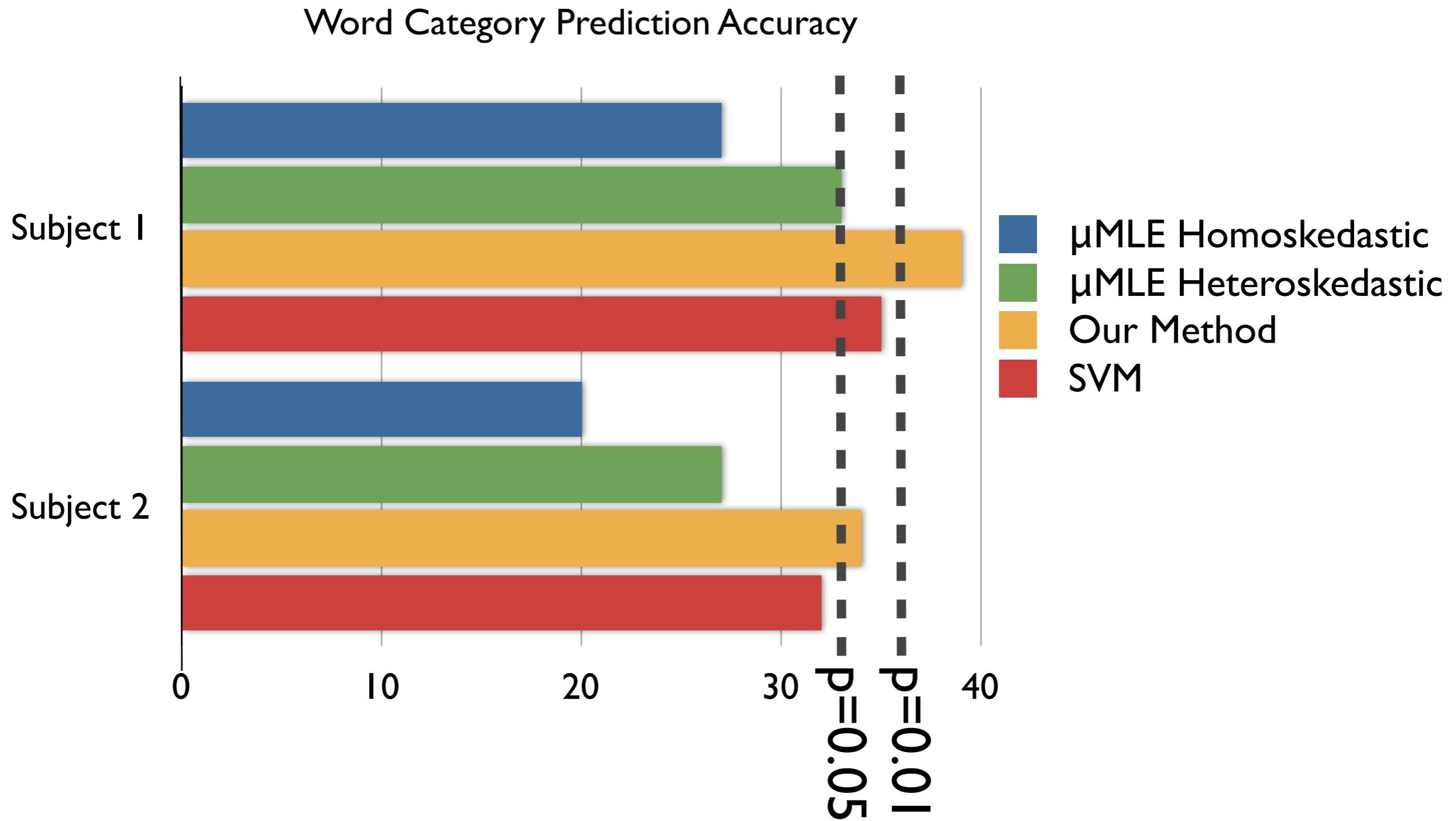
Testing

- Choose category of word that maximizes log likelihood

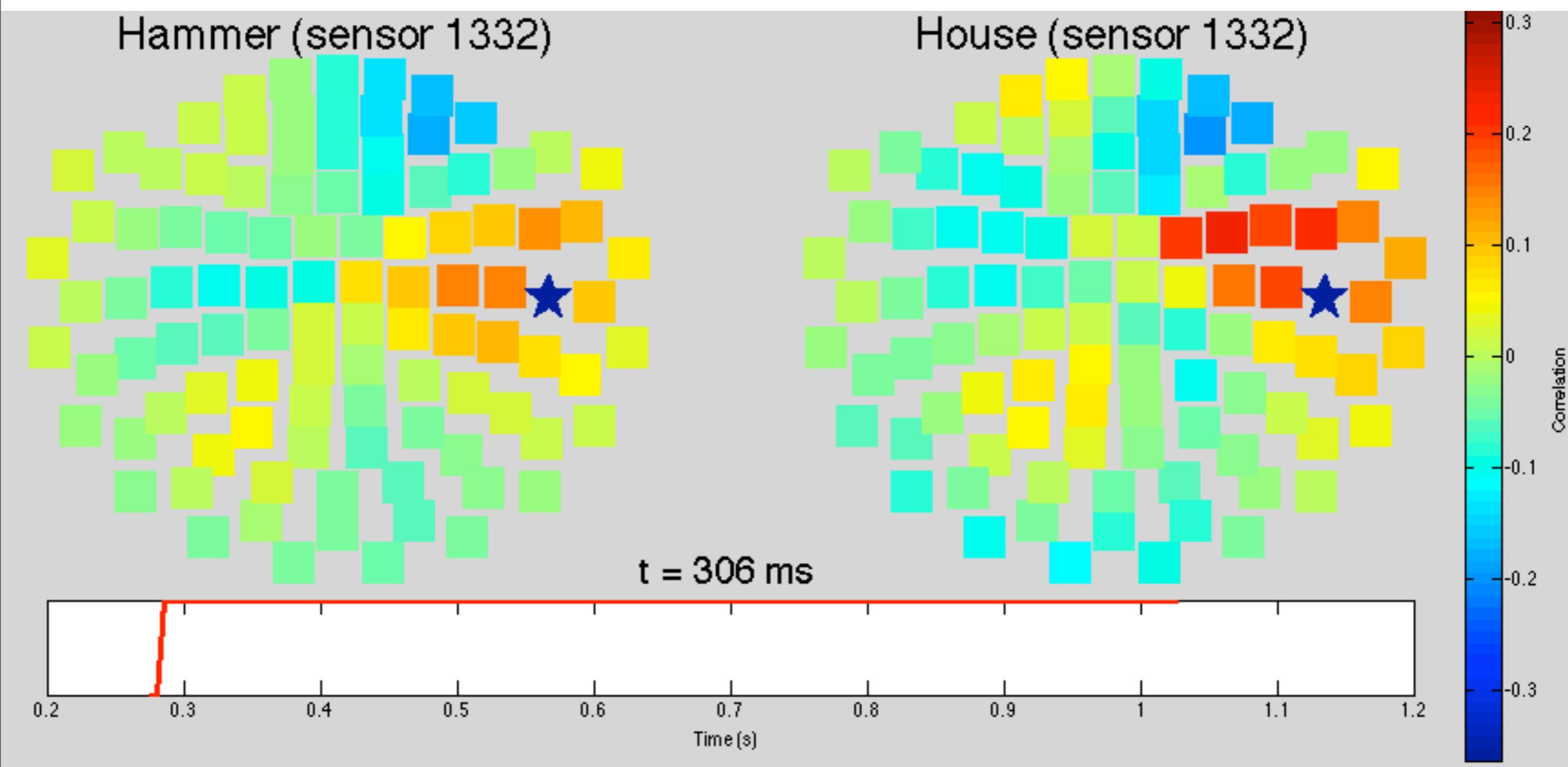
$$\underset{w}{\operatorname{argmax}} \left\{ ll(y_{1:n}^{(*)} | \mu, \Sigma) \right\}$$

- Evaluate Linear SVM for comparison
 - 4 one vs all SVMs
 - ties broken by distance to hyperplane

Results



Results



Summary

- Heteroskedastic model
- Hierarchy of latent processes

- 1) Noise
- 2) High Dimension
- 3) Brain coordination changes over time
- 4) Trial-to-Trial Variability

Future Work

- Extend hierarchy
 - over words in categories
 - all words
 - subjects

Thanks!

Contact:

Alona Fyshe afyshe@cs.cmu.edu

Code (soon):

<http://www.cs.cmu.edu/~afyshe/papers/aistats2012>

Funding:

- Natural Sciences and Engineering Research Council of Canada (NSERC)
- W. M. Keck Foundation
- National Science Foundation (NSF)
- National Institute of Environmental Health Sciences