

# Learning in neuronal networks: Processing high-order statistics embedded in time series for classification tasks

Matthieu Gilson  
Chaire of Junior Prof  
<https://matthieugilson.eu>

- Sofia Lawrie
- Rubén Moreno-Bote

**upf.**

**Universitat  
Pompeu Fabra**  
*Barcelona*



- Sandra Nestler
- David Dahmen
- Moritz Helias

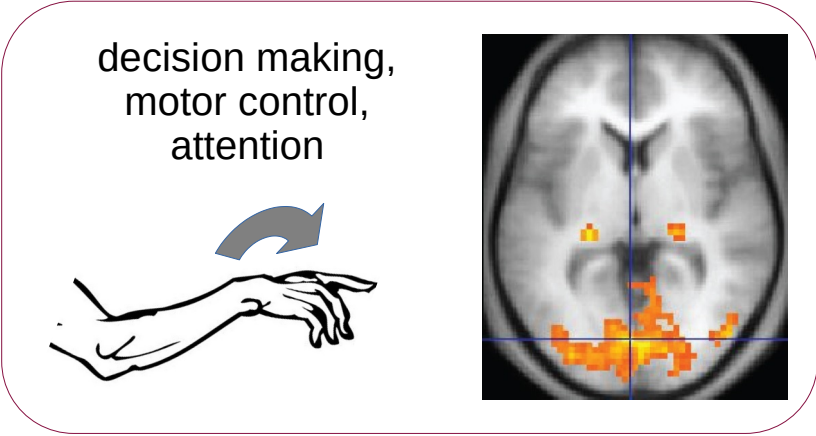
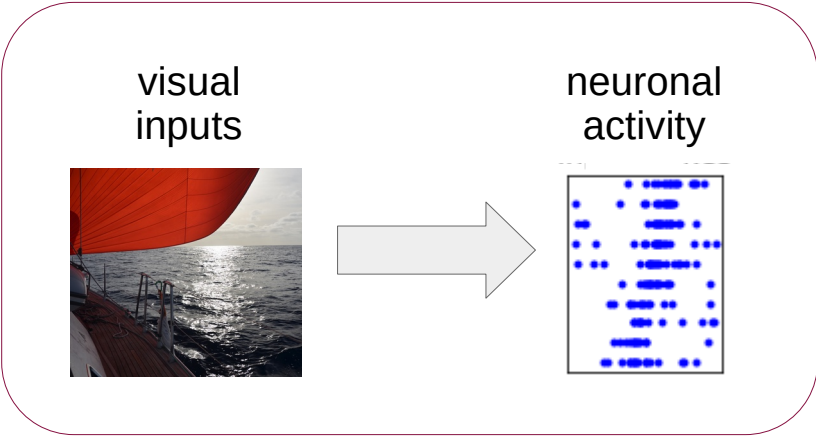
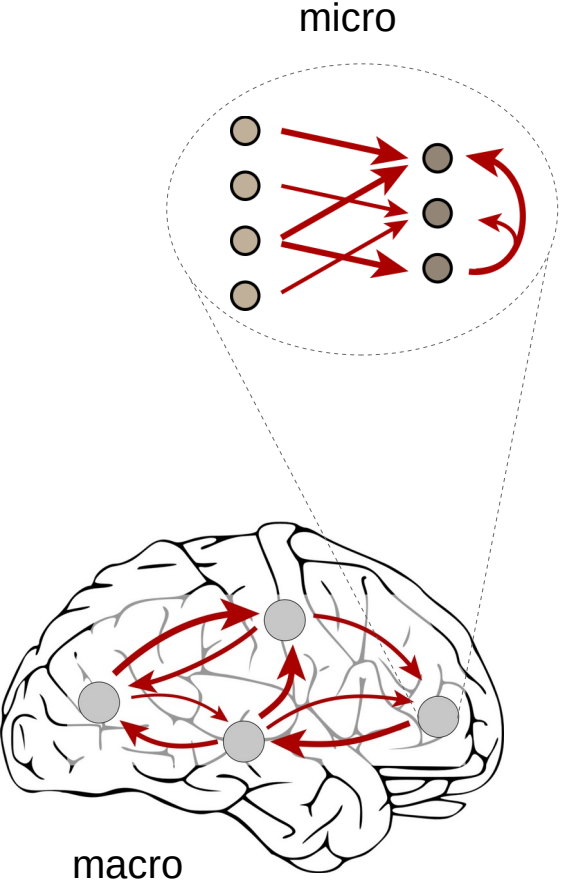


- Laurent Perrinet

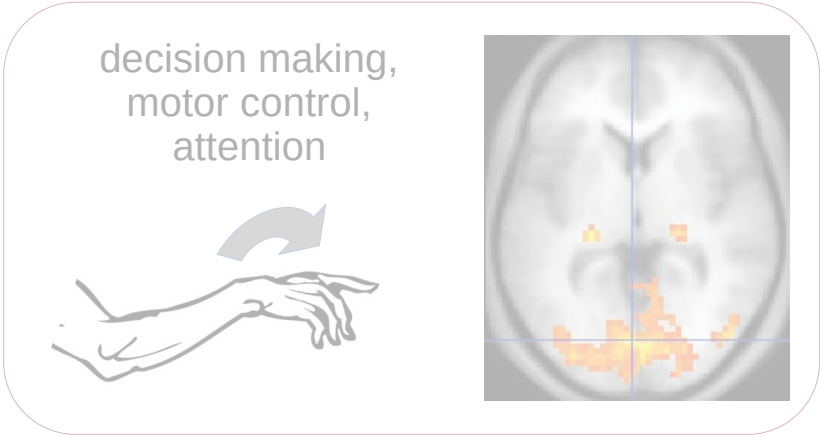
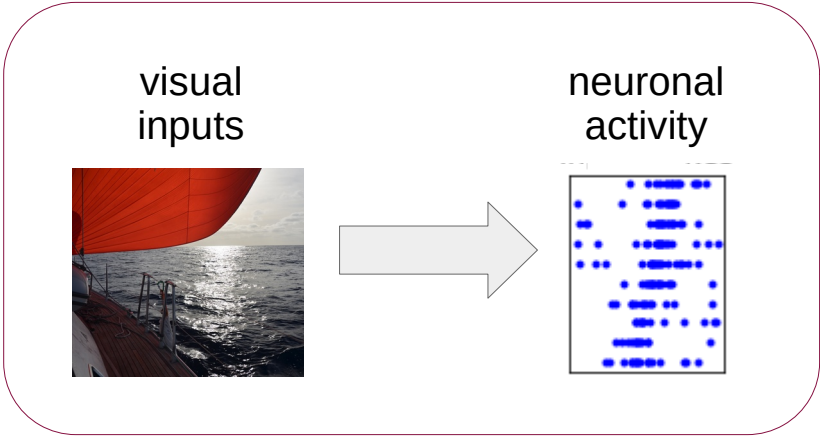
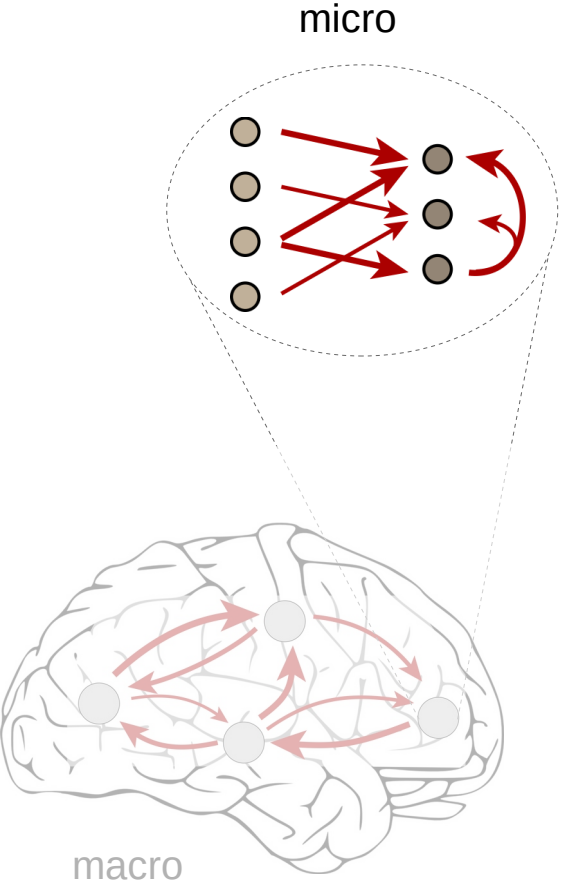


Institut de  
Neurosciences des  
Systèmes

# Neuronal processing in the brain: computations and communication

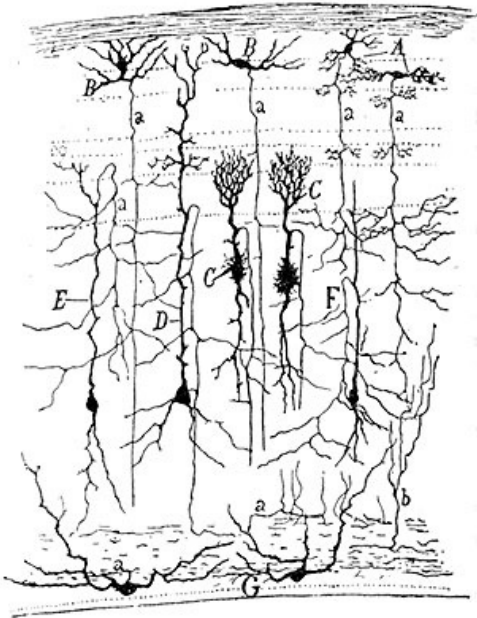


# Neuronal processing in the brain: computations and communication



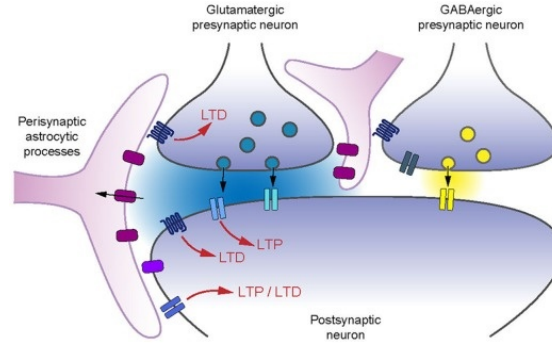
# How does the brain implement functions?

- Karl Lashley (1890-1958): storage of memory in brain regions (engram)



Ramon y Cajal 1905

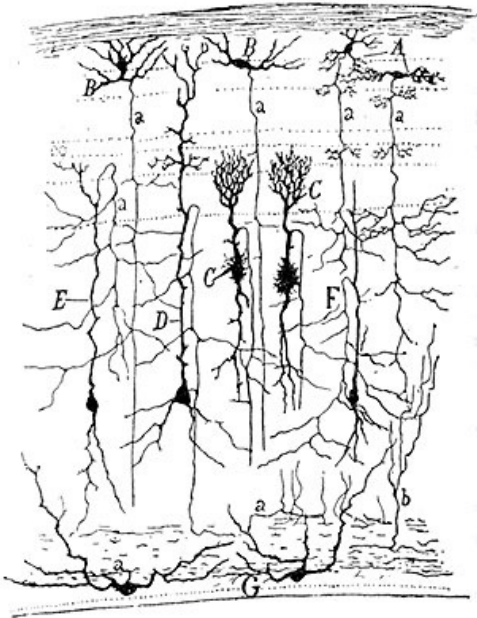
- Synaptic plasticity (learning) shapes the neuronal network dynamics



Valcheva et al. *Front Syn Neurosci* 2019

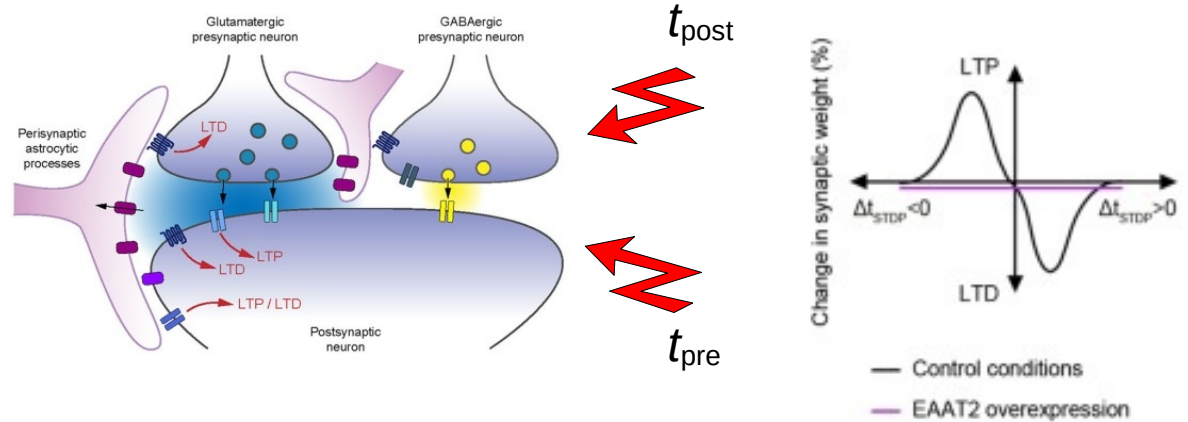
# How does the brain implement functions?

- Karl Lashley (1890-1958): storage of memory in brain regions (engram)



Ramon y Cajal 1905

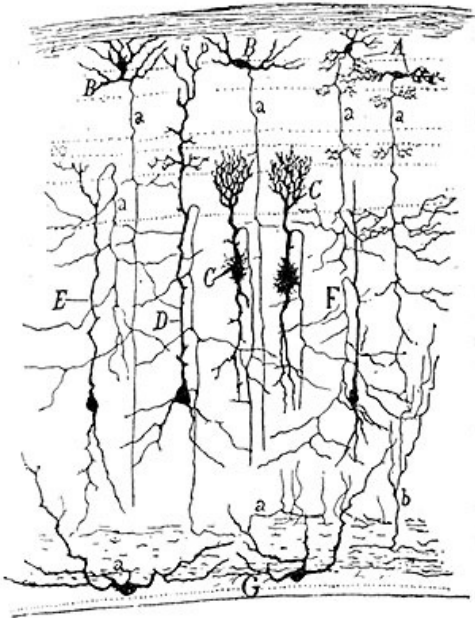
- Synaptic plasticity (learning) shapes the neuronal network dynamics



Valcheva et al. *Front Syn Neurosci* 2019

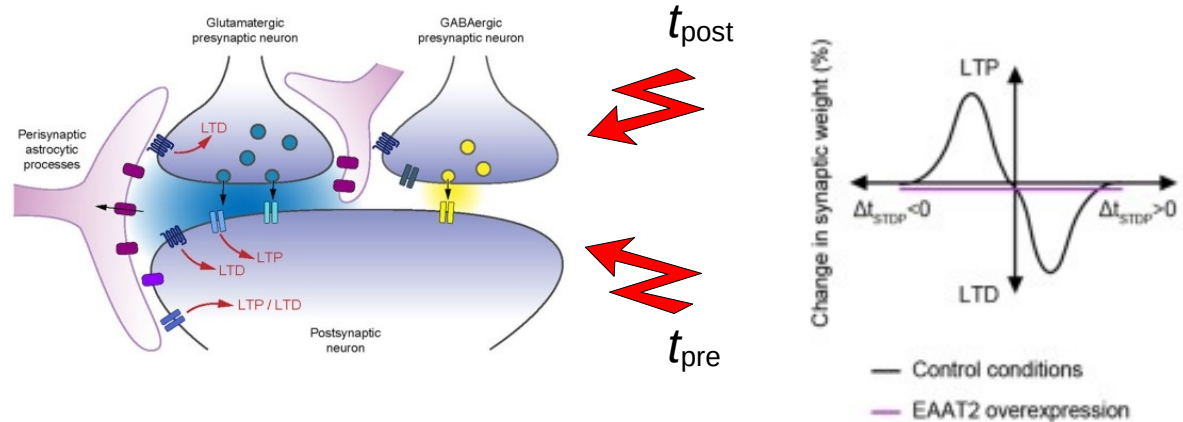
# How does the brain implement functions?

- Karl Lashley (1890-1958): storage of memory in brain regions (engram)



Ramon y Cajal 1905

- Synaptic plasticity (learning) shapes the neuronal network dynamics

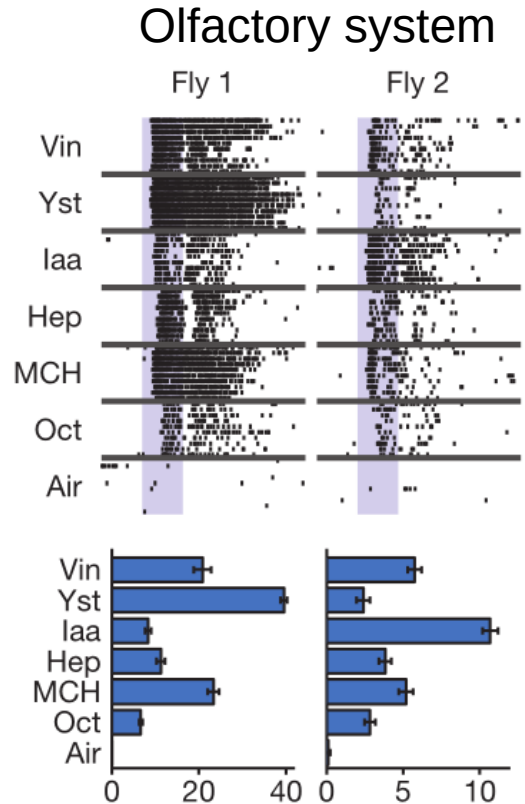


Valcheva et al. *Front Syn Neurosci* 2019

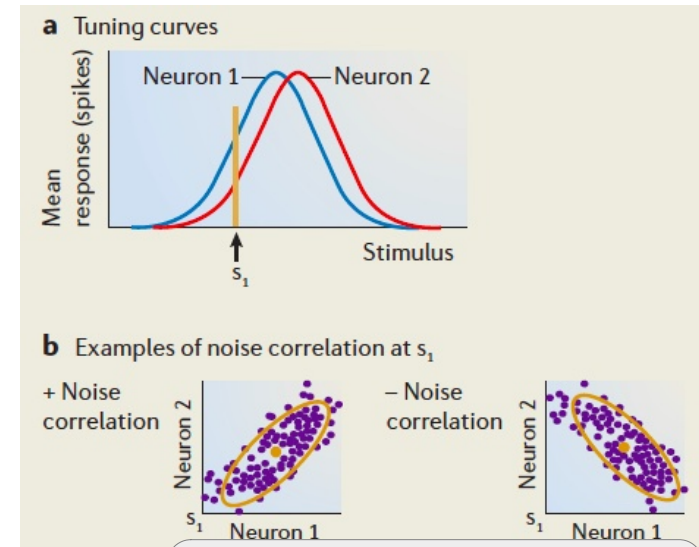
If plasticity depends on second-order (spiking) statistics, shouldn't the neuronal representations ("code") be based on them as well?



# Traditional view: neuronal representations based on spike counts



### Visual system tuning curves for orientation selectivity

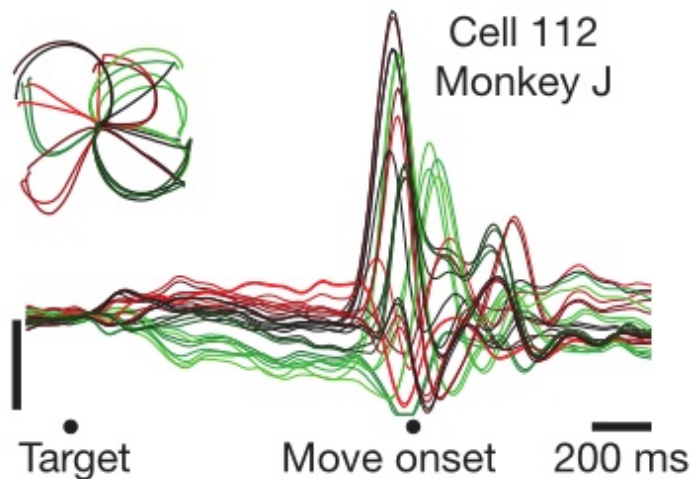


**noise correlations**  
spike count across trials



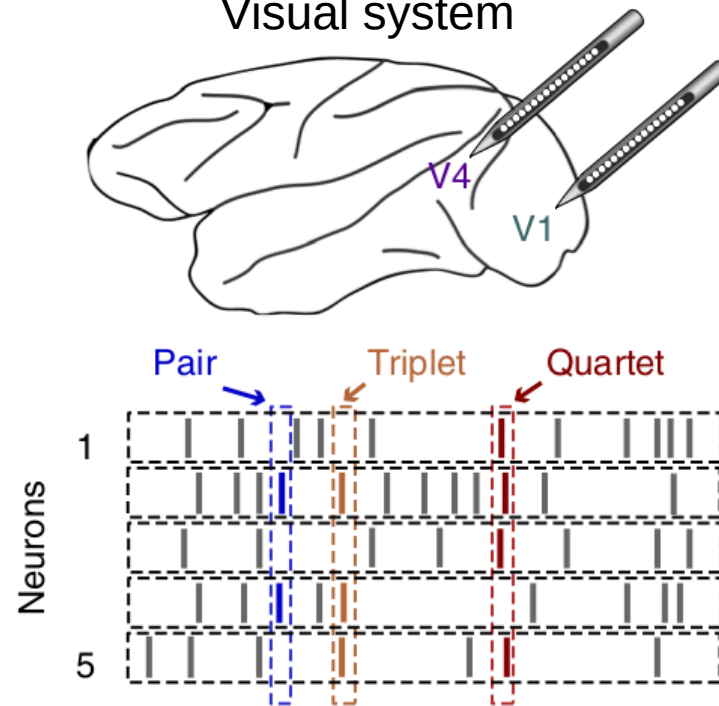
# Temporal structure matters! Rate fluctuations and spike synchrony

## Motor system



**fast rate correlation**  
within trial

## Visual system

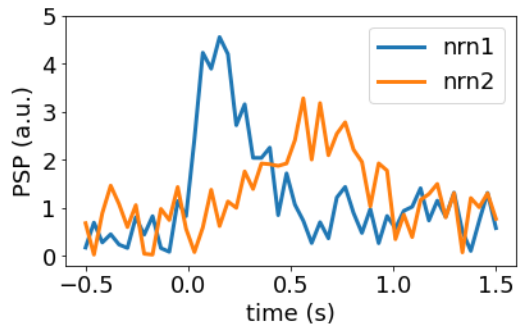


**spike correlations**  
within trial

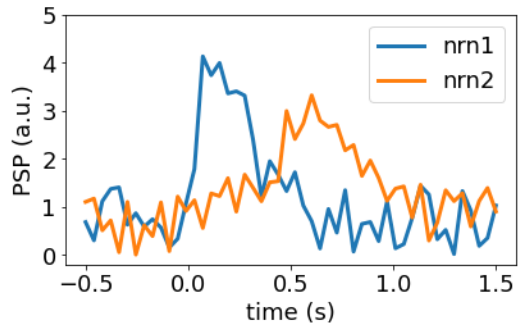
# Types of structured variability in temporal (neuronal) signals: which measure to apply to time series?

Stable profile

Trial 1



Trial 2



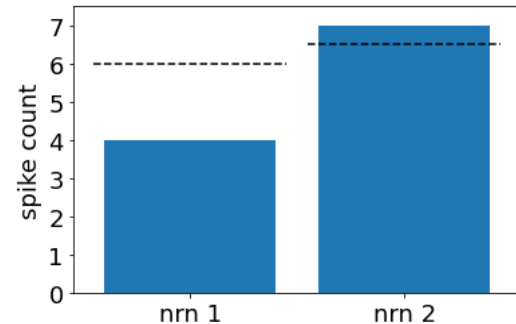
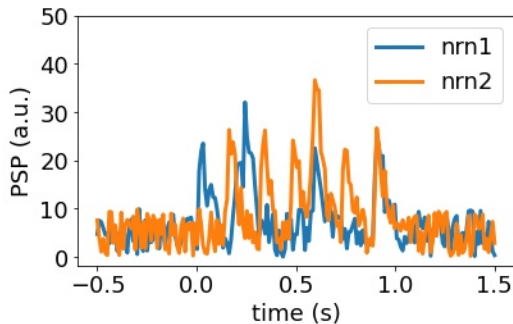
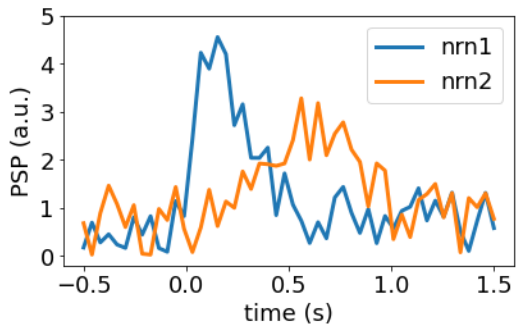
# Types of structured variability in temporal (neuronal) signals: which measure to apply to time series?

Stable profile

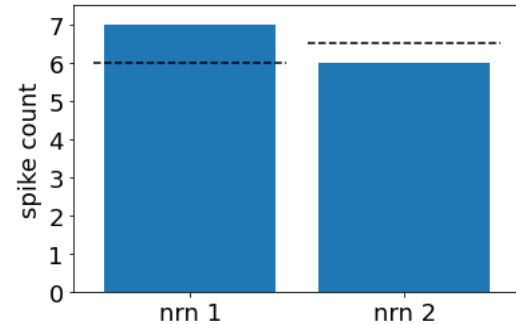
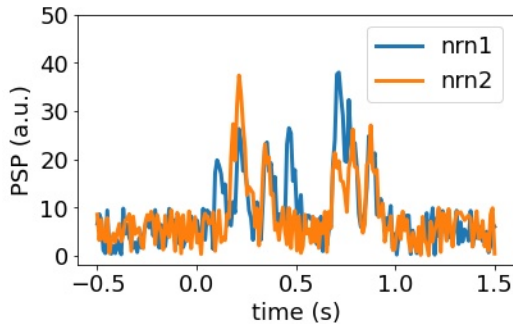
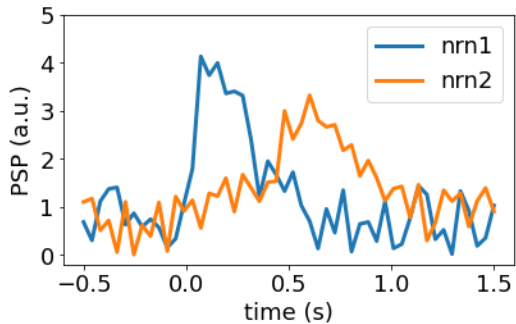
Cofluctuations

Spike count

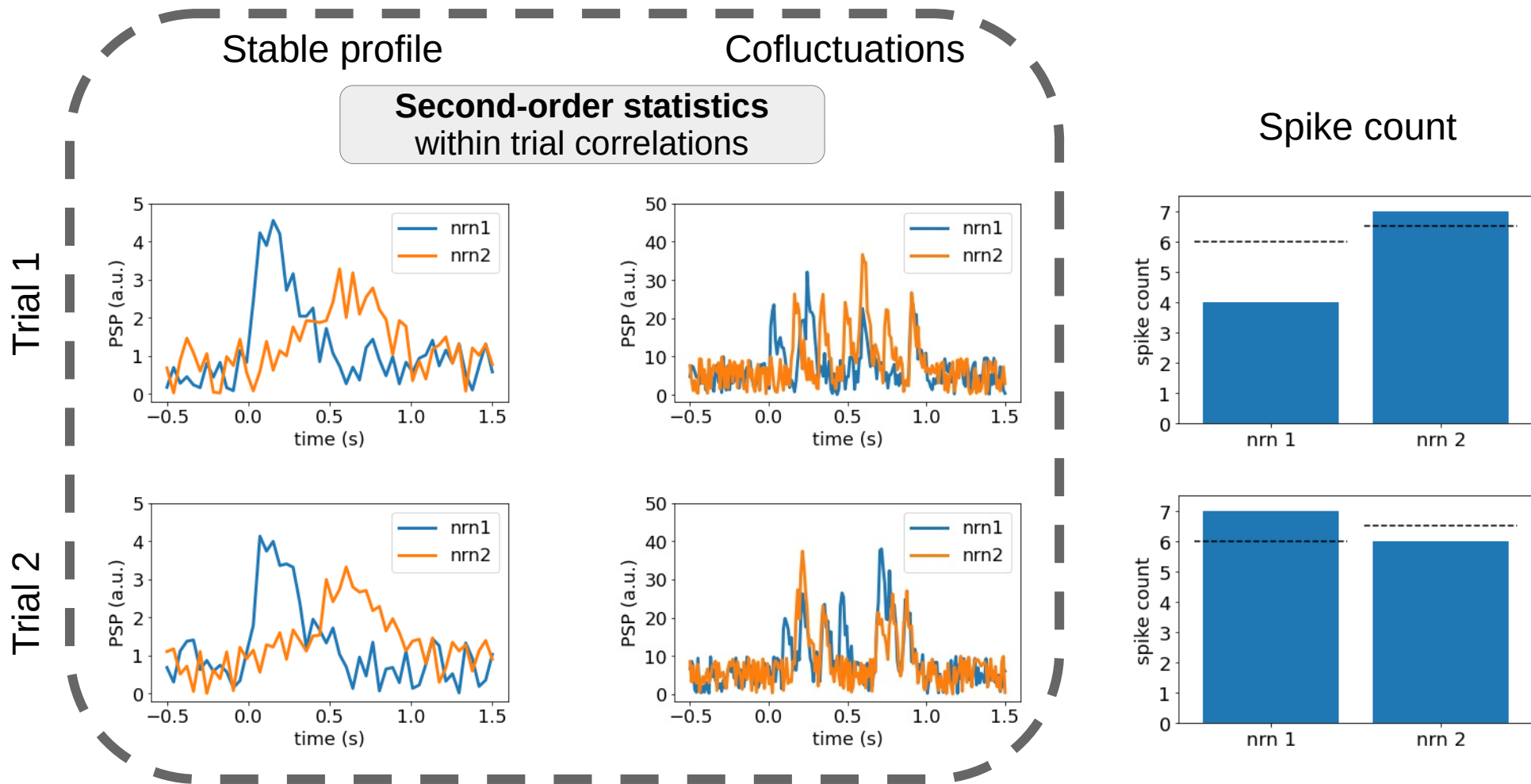
Trial 1



Trial 2

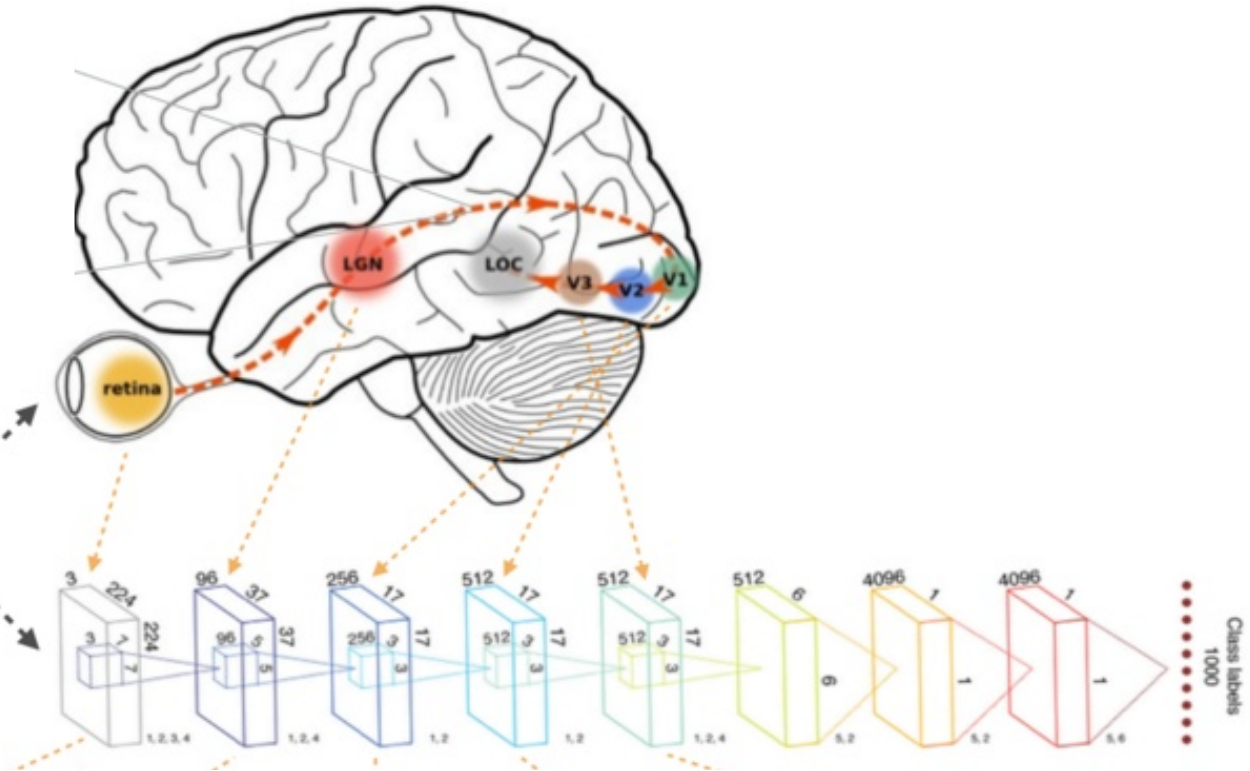


# Types of structured variability in temporal (neuronal) signals: which measure to apply to time series?



# Computations: biological processing versus machine learning

Beyond static images?  
**Time series**



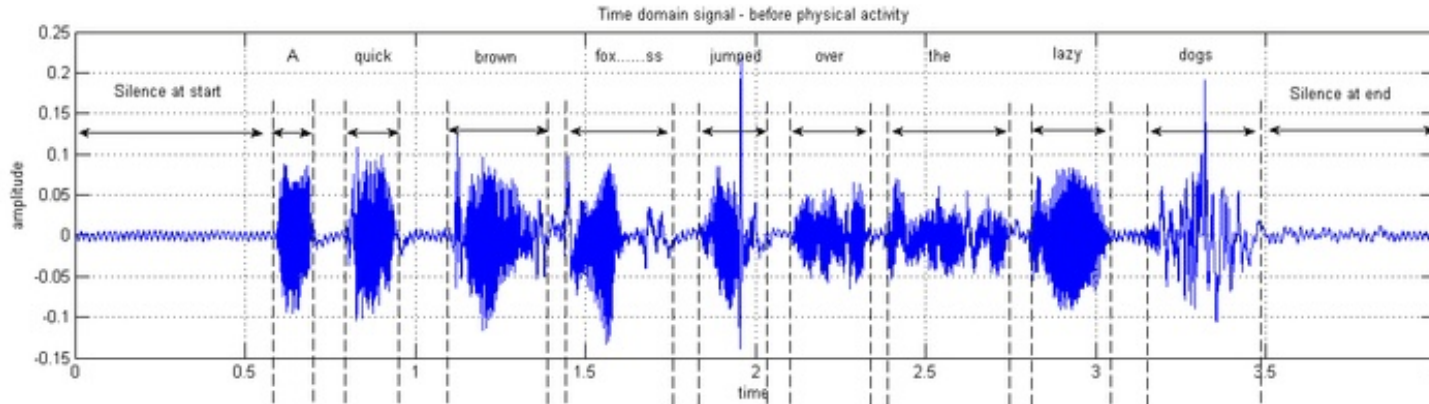
Deep learning networks: what is processed at each layer?  
**Interpretability**

# Outline

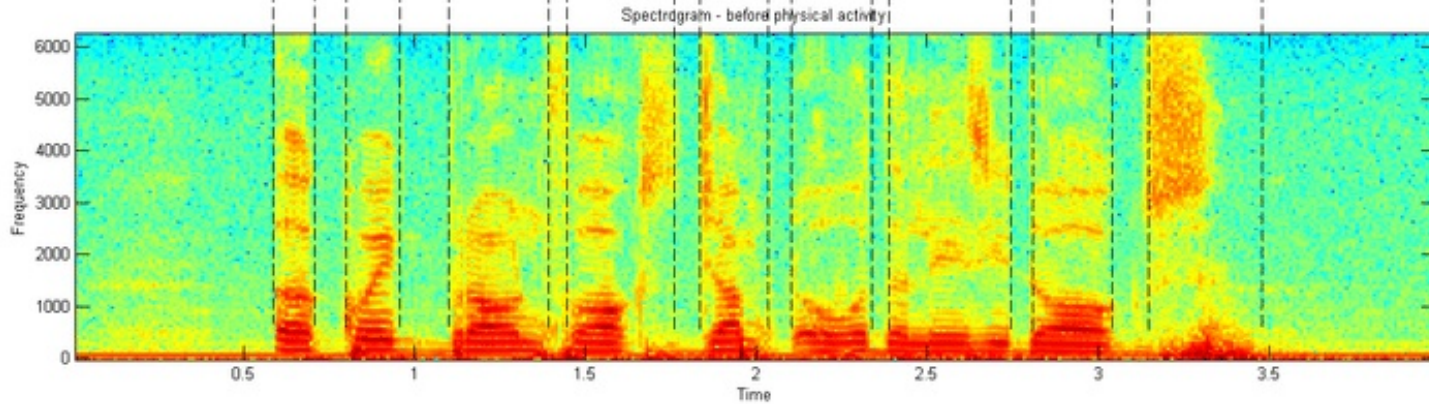
- **Statistical learning for time series**
  - **mean versus covariance decoding**
  - processing by neuronal reservoir
  - decoding by biological architecture
- Theory

# Phonemes in speech

soundwave



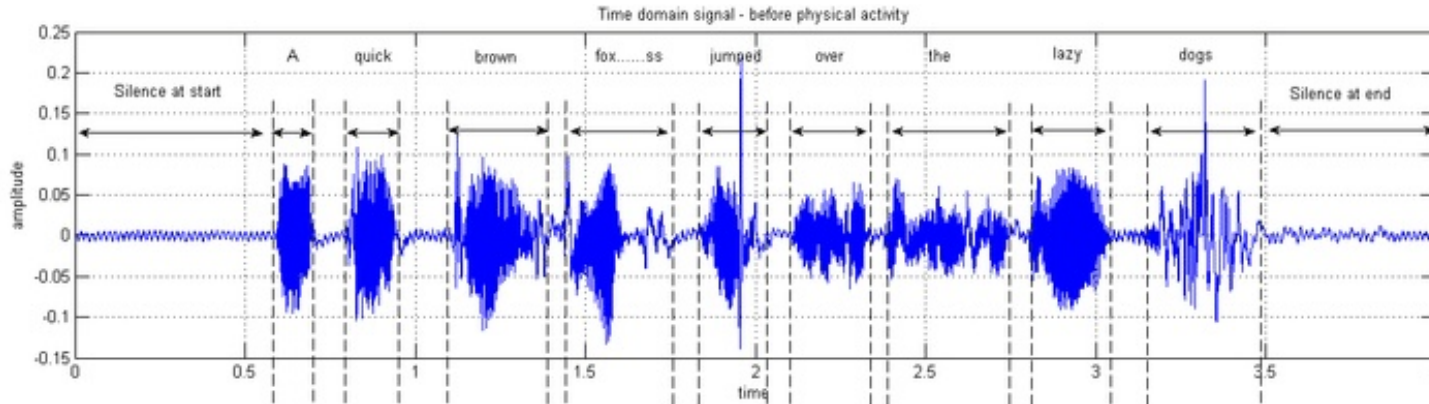
spectrogram



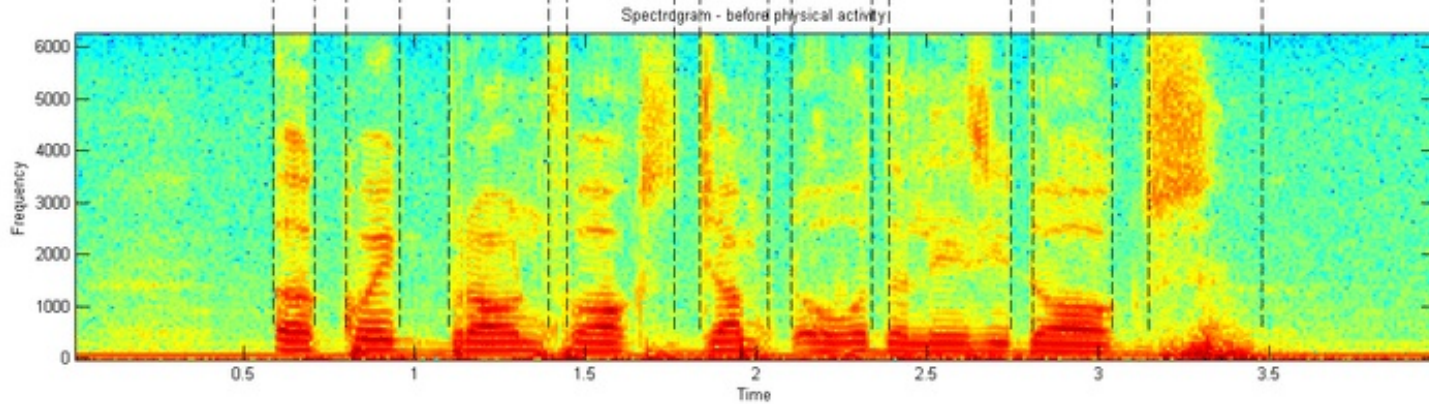


# Phonemes in speech

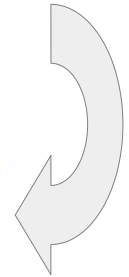
soundwave



spectrogram

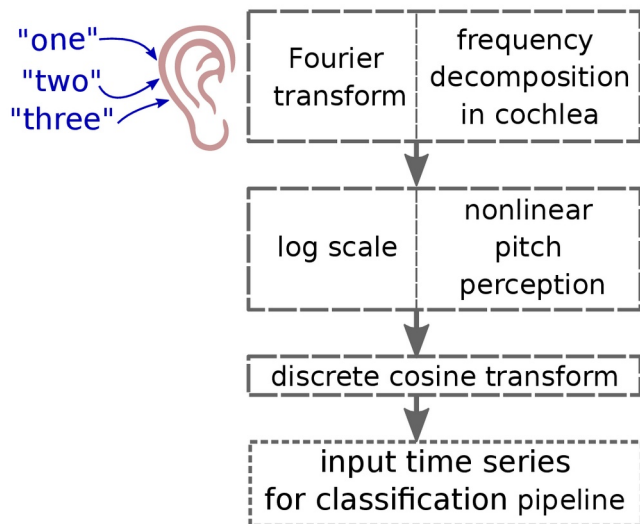


inner  
ear:  
cochlea

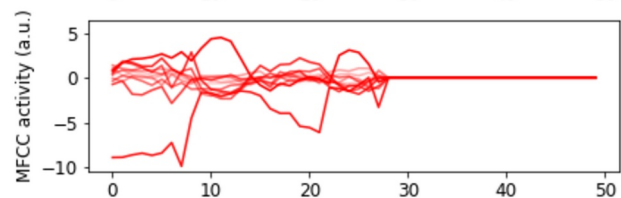


# Cochlear processing and spatio-temporal structure

## Preprocessing of spoken digits dataset



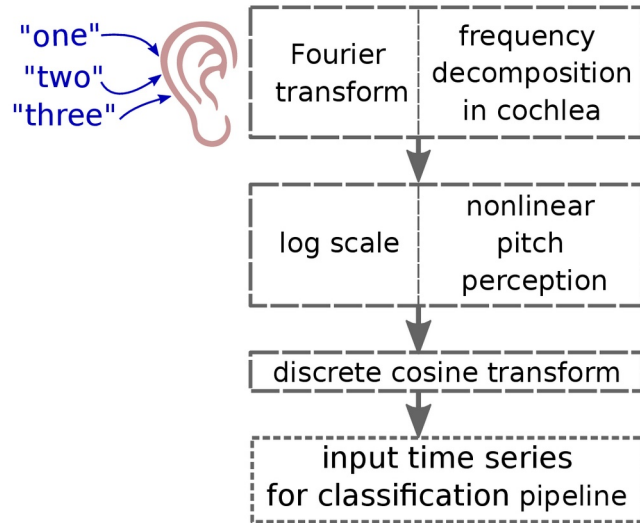
## Example time series of digits



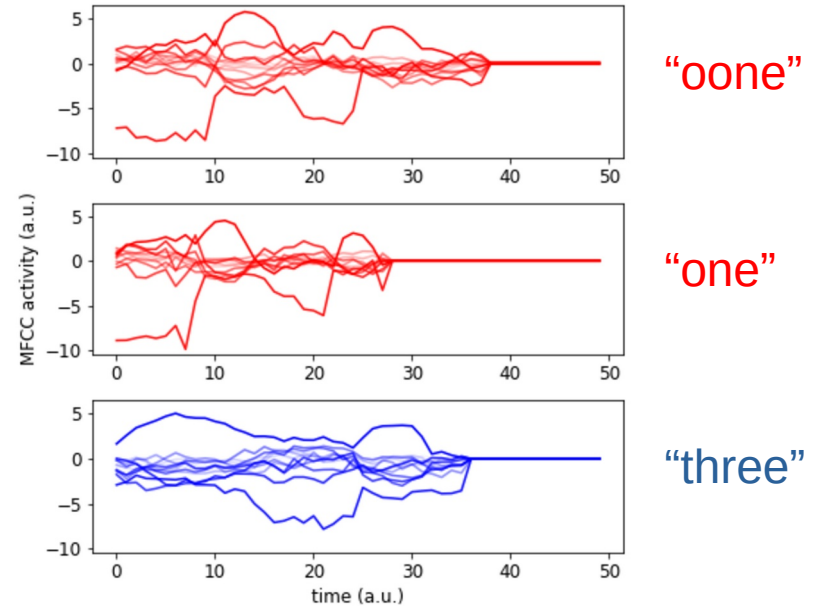
“one”

# Cochlear processing and spatio-temporal structure

## Preprocessing of spoken digits dataset

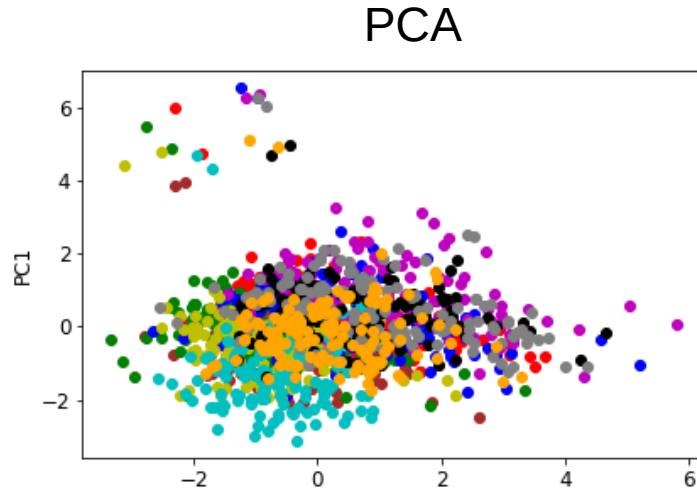


## Example time series of digits

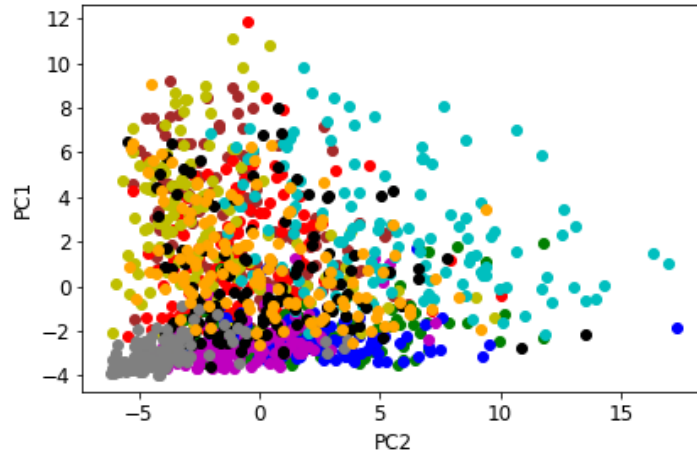


# Separability of spoken digits (10 classes)

mean



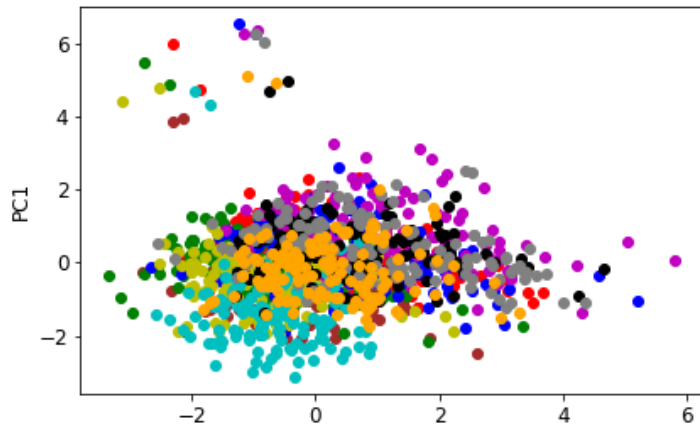
COV



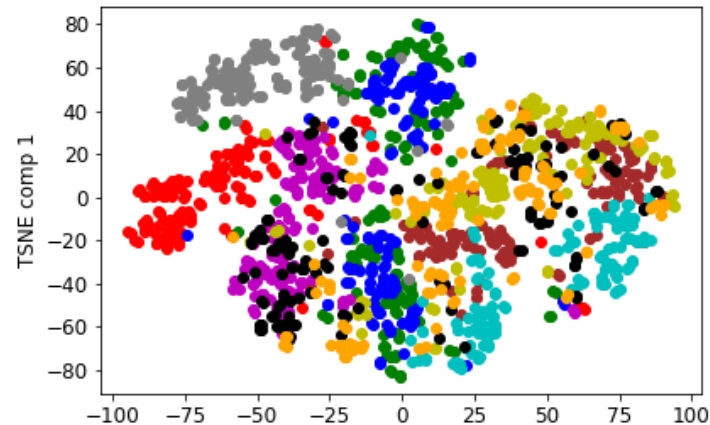
# Separability of spoken digits (10 classes)

mean

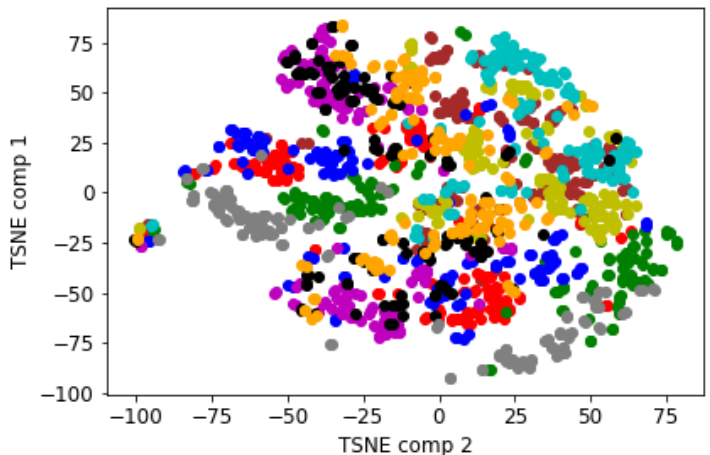
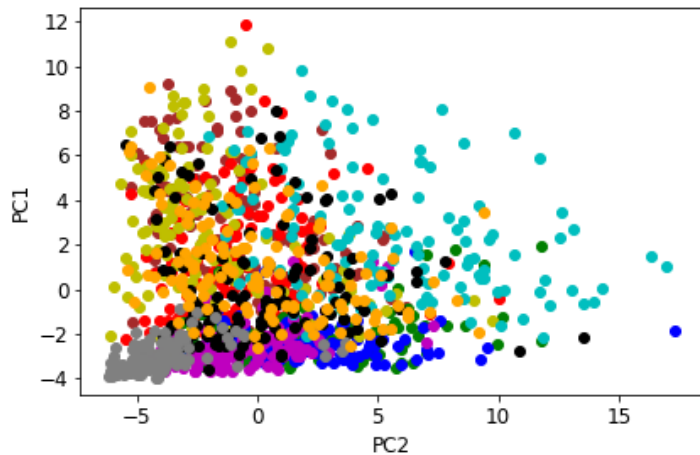
PCA



tSNE



COV



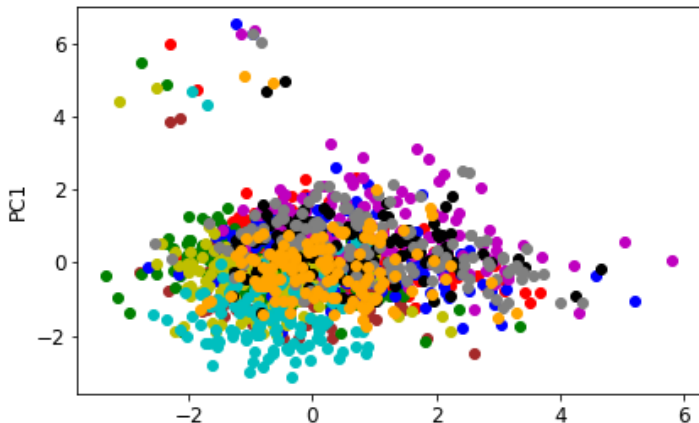
# Cross-validated classification to assess baseline separability

mean **69.7%**

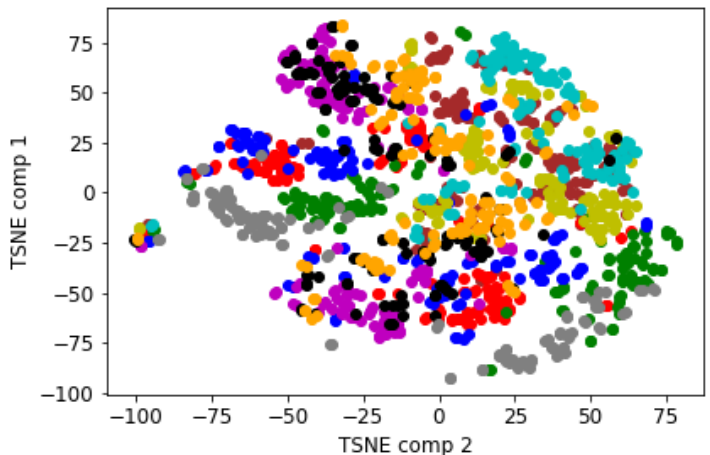
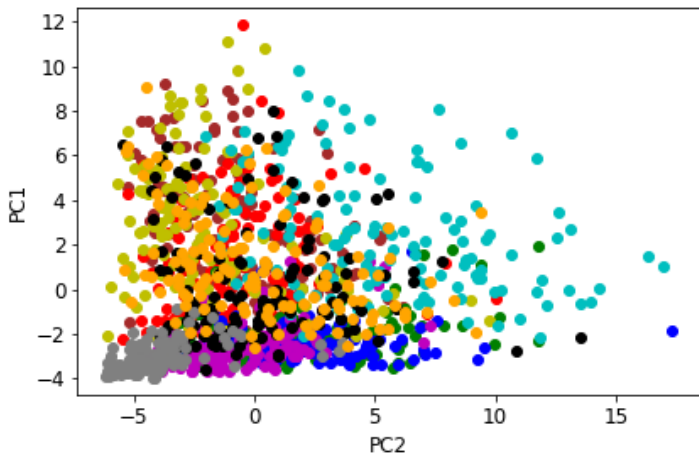
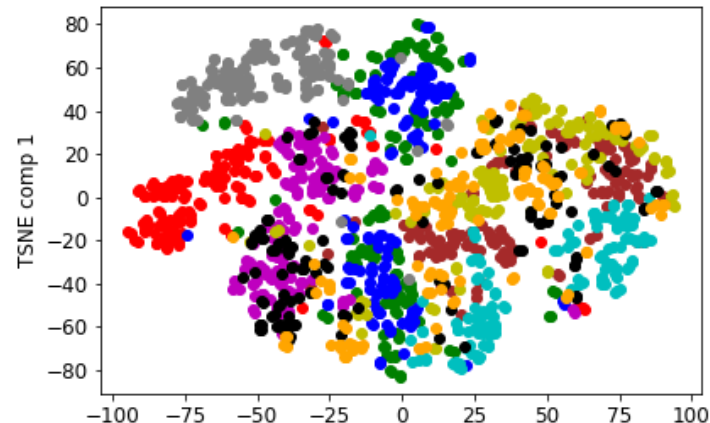
**Logistic  
regression**

COV **91.7%**

PCA



tSNE



# Outline

- **Statistical learning for time series**
  - mean versus covariance decoding
  - processing by neuronal reservoir
  - decoding by biological architecture
- Theory

- **Structured variability conveys information**



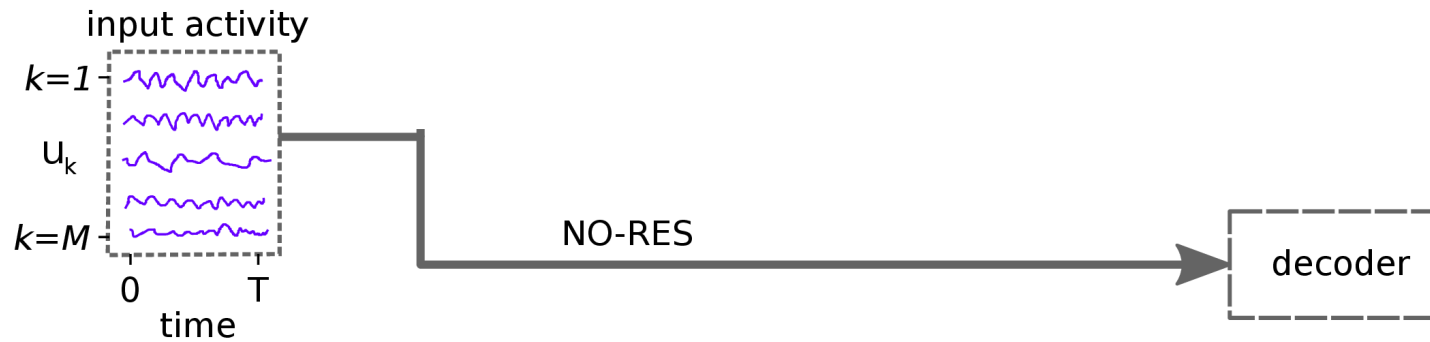
# Outline

- **Statistical learning for time series**
  - mean versus covariance decoding
  - **processing by neuronal reservoir**
  - decoding by biological architecture
- Theory

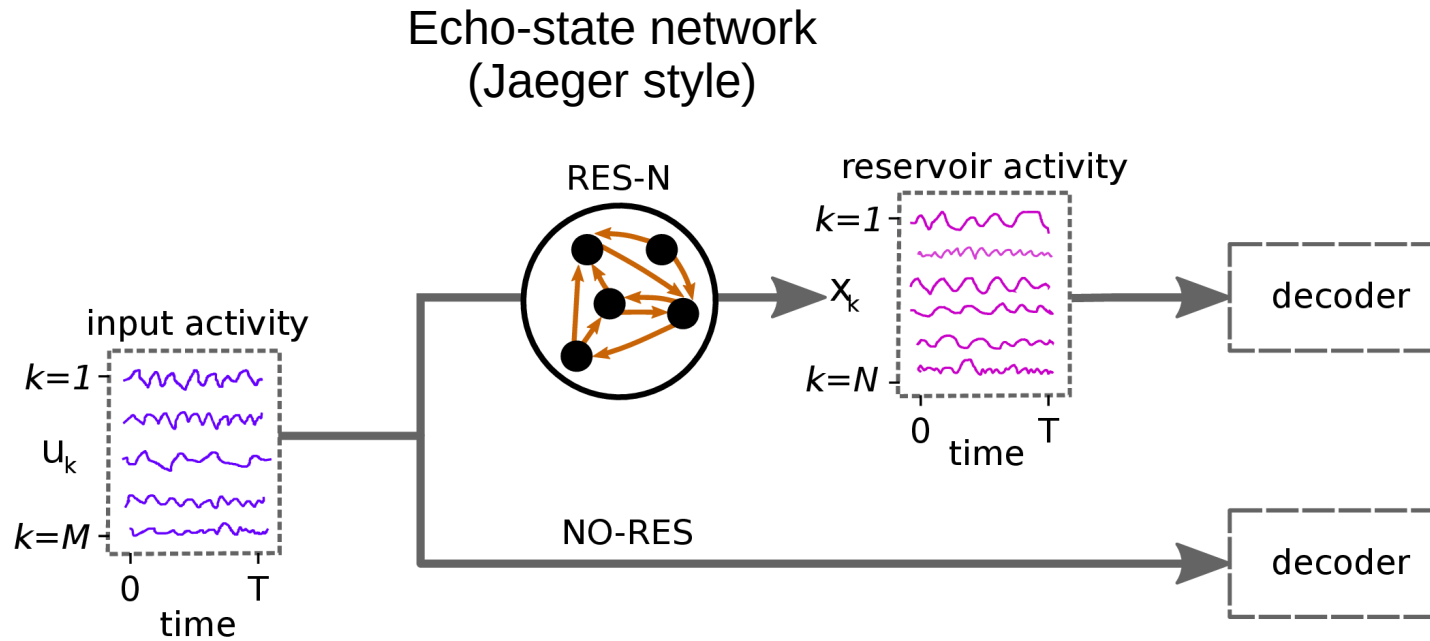
- Structured variability conveys information

- **How to efficiently extract statistical patterns?**

# Reservoir computing to process input time series



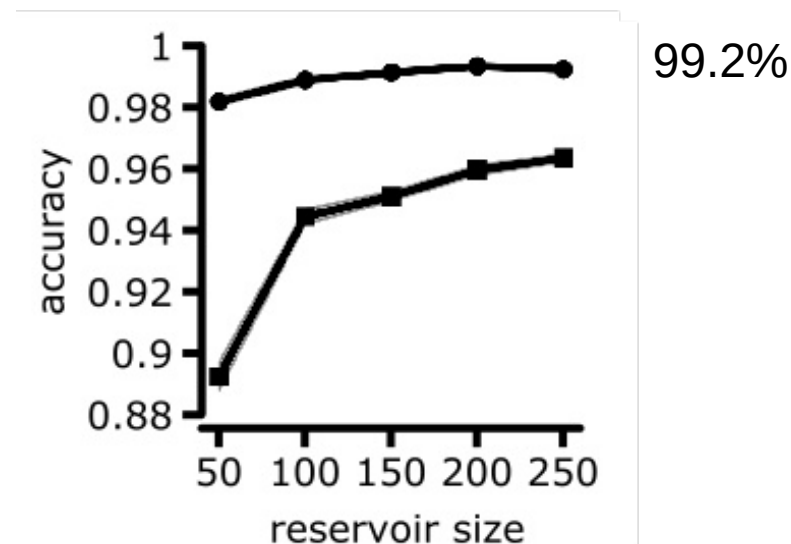
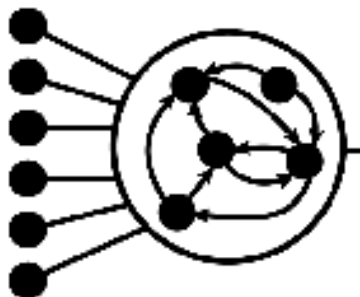
# Reservoir computing to process input time series



# Dimensionality expansion

Input size  
 $M=13$

Reservoir  
size  $N$



● covariance ■ mean

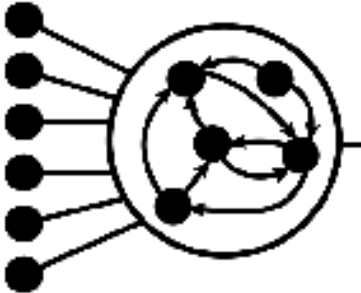
NO-RES ref: 91.7%

69.7%

# Dimensionality expansion

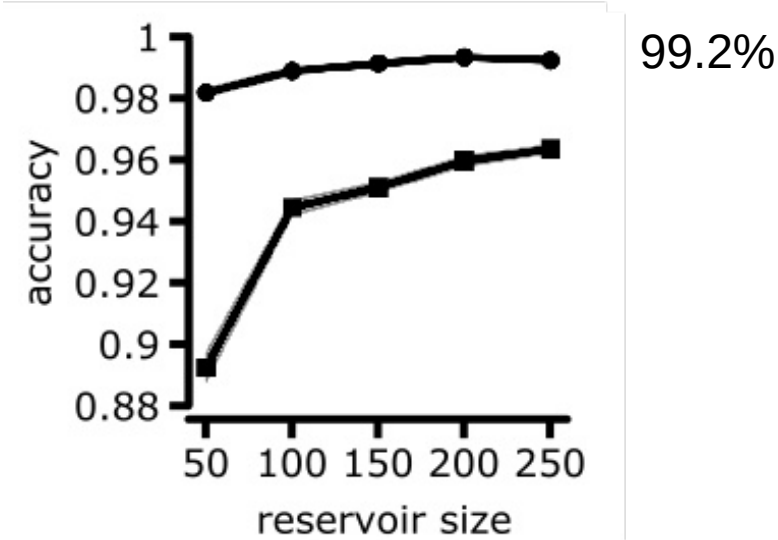
Input size  
 $M=13$

Reservoir  
size  $N$



Reservoir feature space:

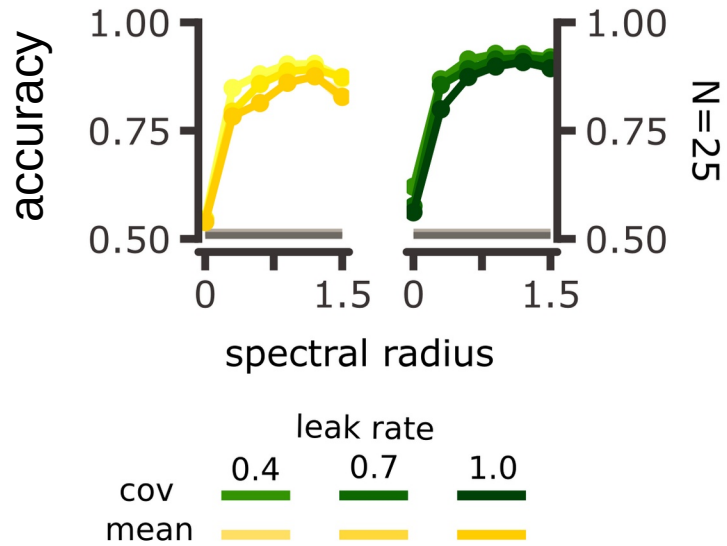
- mean vector  $\mathbf{N}$
- cov matrix  $\mathbf{N}^2$



● covariance ■ mean

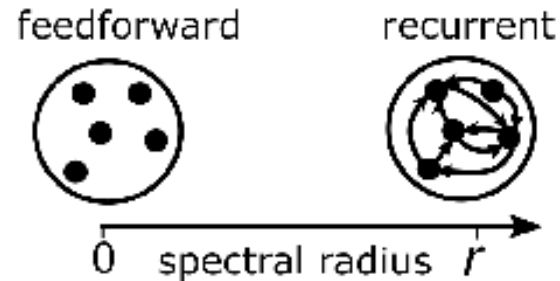
NO-RES ref: 91.7% 69.7%

# Influence of reservoir parameters on decoding accuracy?



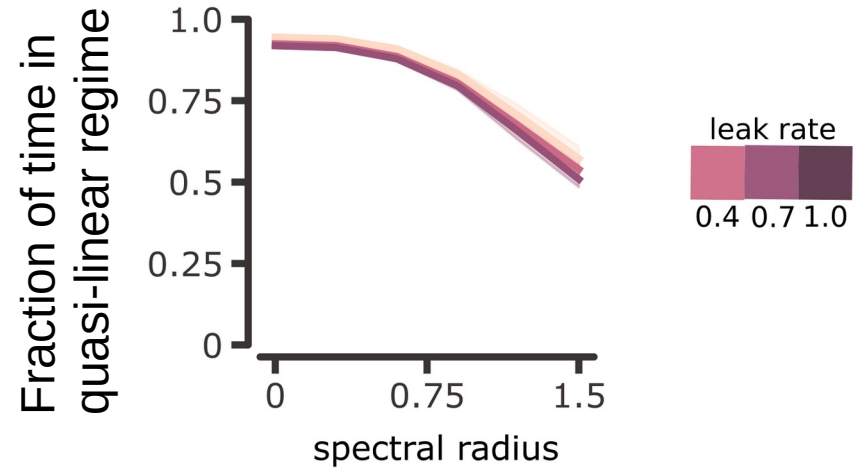
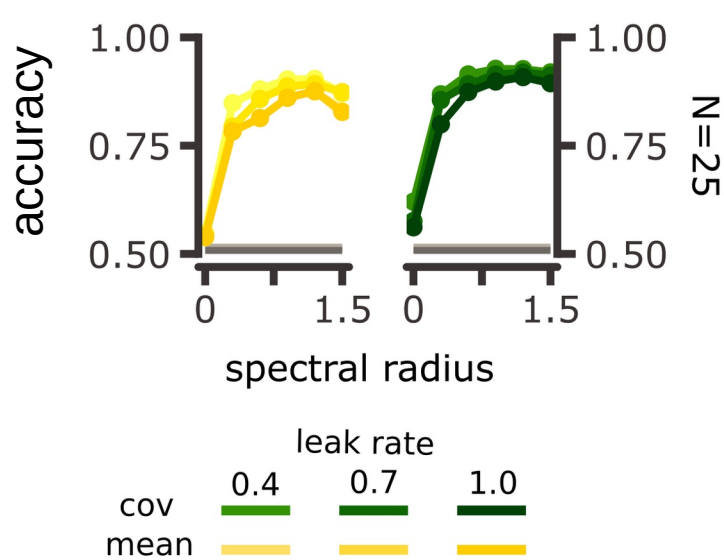
Parameter exploration:

- spectral radius
- nodal leak rate



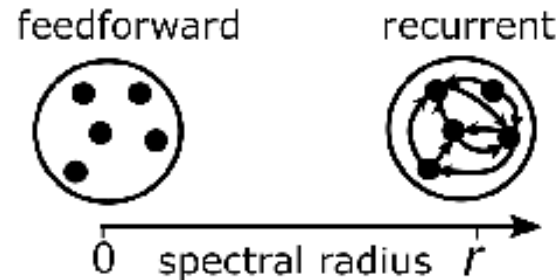
# Influence of reservoir parameters on decoding accuracy?

## Nonlinear reservoir dynamics support best decoding



Parameter exploration:

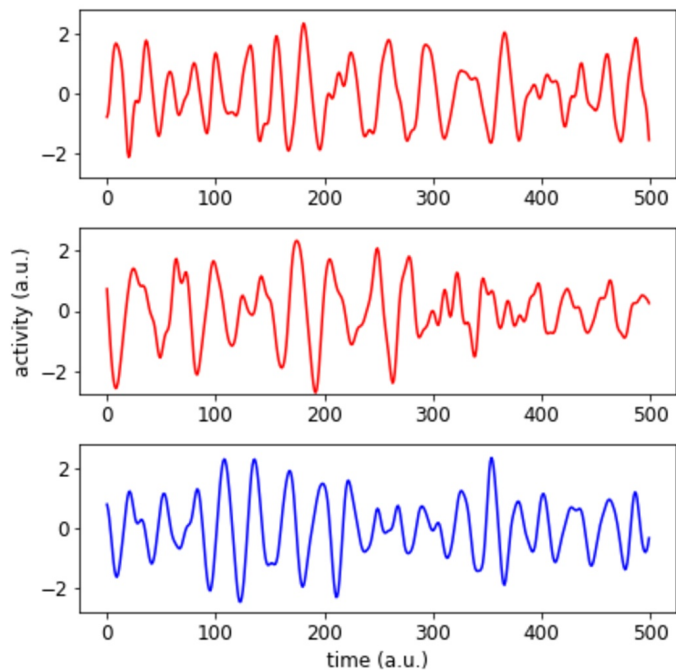
- spectral radius
- nodal leak rate





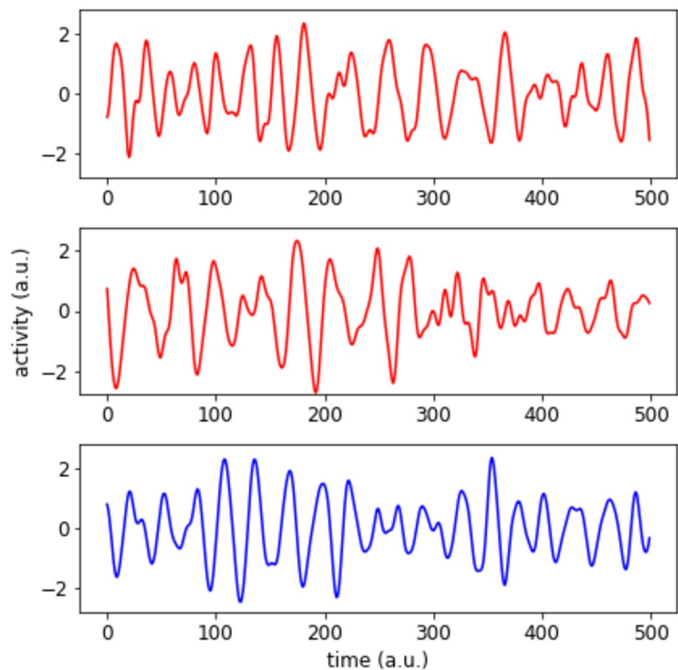
# Example dataset 2: faulty engine signals

Example time series for  
faulty engine signals

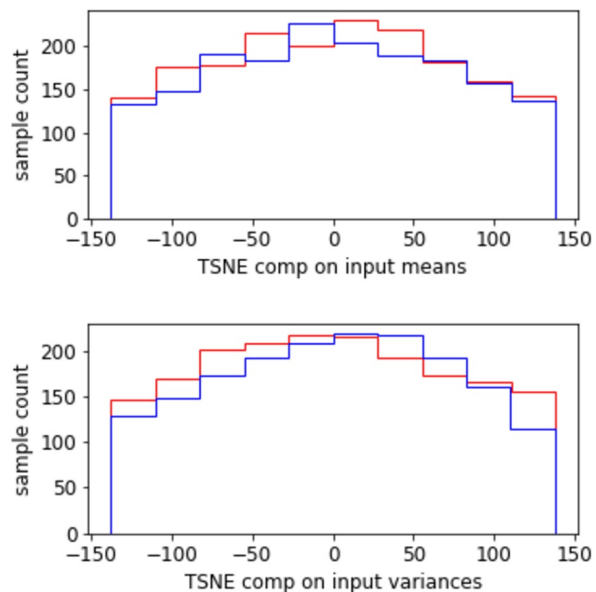


# Example dataset 2: faulty engine signals

## Example time series for faulty engine signals

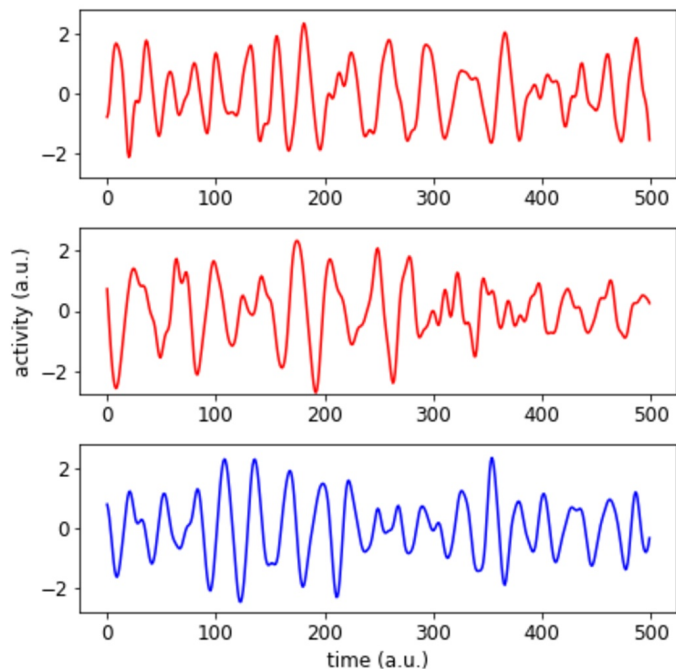


## TSNE on means/variances for faulty engine signals

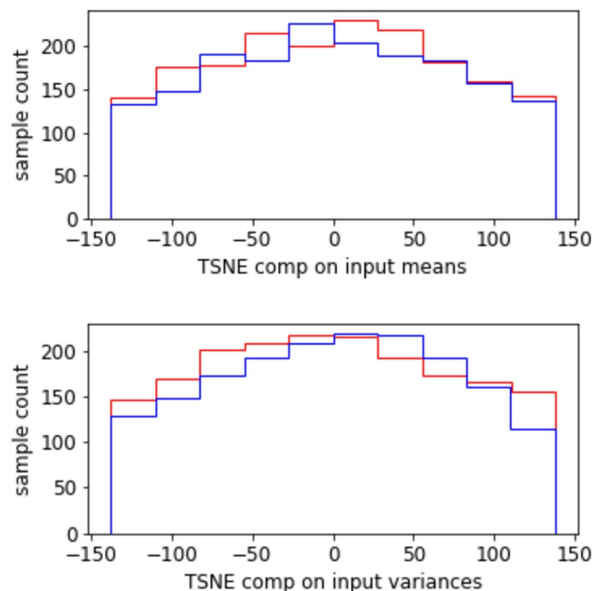


# Example dataset 2: faulty engine signals

## Example time series for faulty engine signals



## TSNE on means/variances for faulty engine signals



NO-RES

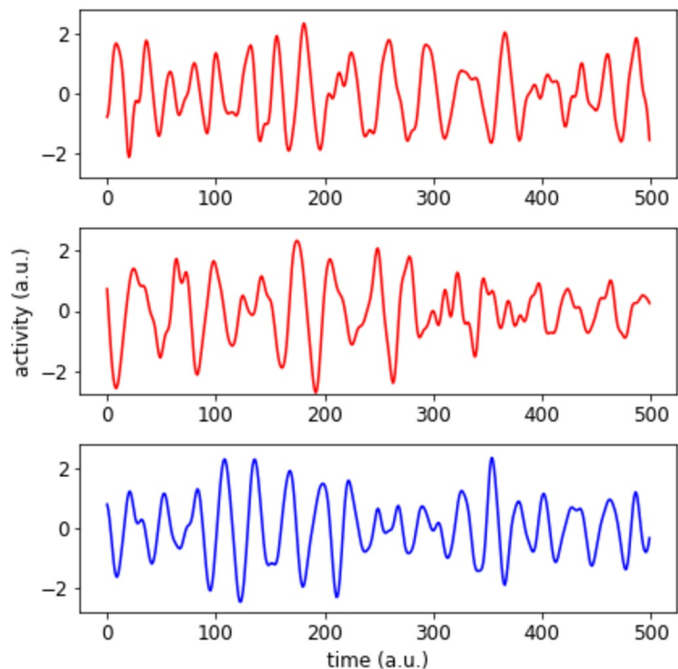
- **mean: ~50%**

- **cov: ~50%**

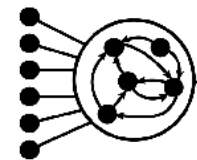
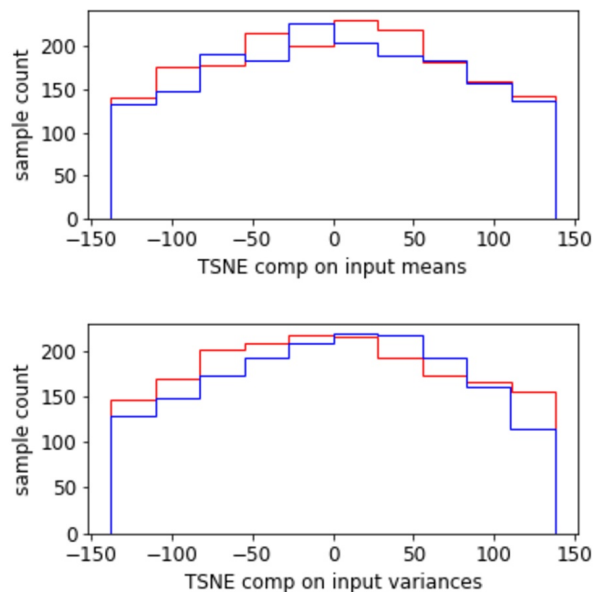
- **best: 96%**

# Example dataset 2: faulty engine signals

## Example time series for faulty engine signals



## TSNE on means/variances for faulty engine signals



(size N=50)

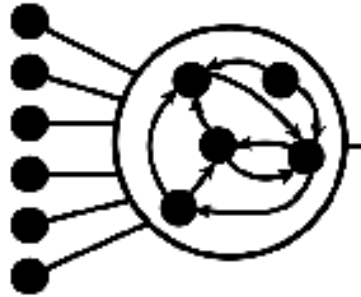
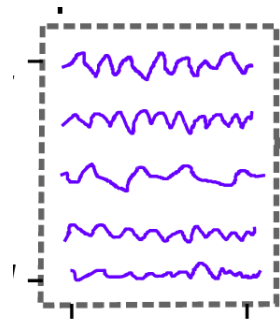
- **mean: 84%**

- **cov: 94%**

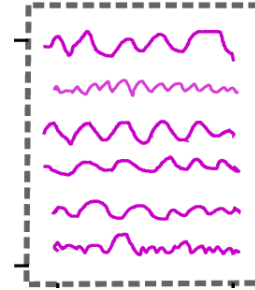
- **best: 96%**

# Neuronal reservoir for time series processing

Input time series

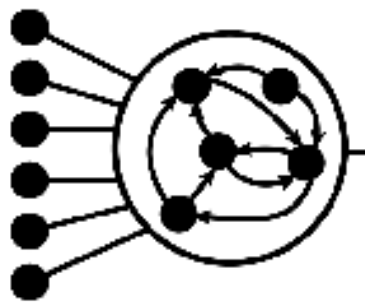
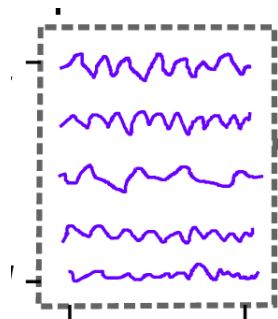


Reservoir time series

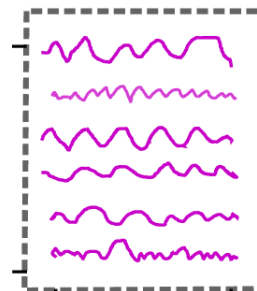


# Neuronal reservoir for time series processing

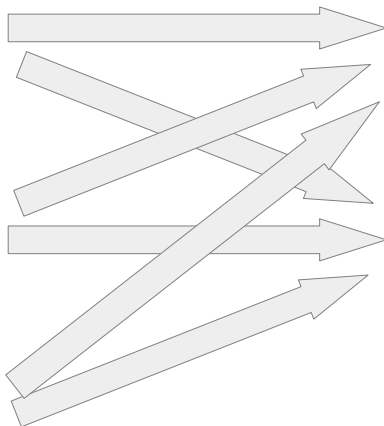
Input time series



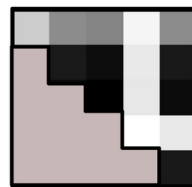
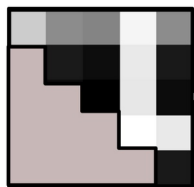
Reservoir time series



1st-order stats



2nd-order stats

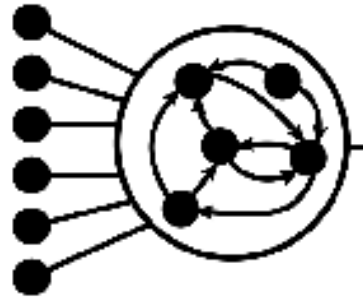
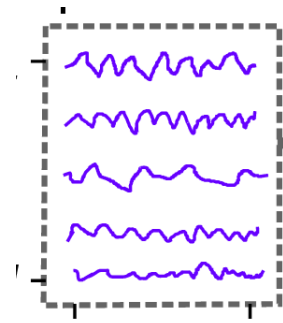


high-order stats

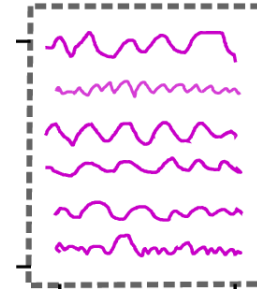
⋮

# Neuronal reservoir for time series processing

Input time series



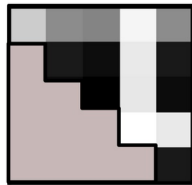
Reservoir time series



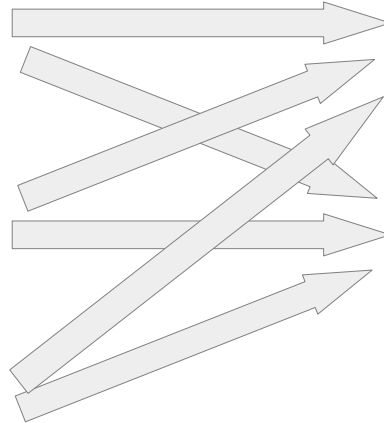
1st-order stats



2nd-order stats



high-order stats

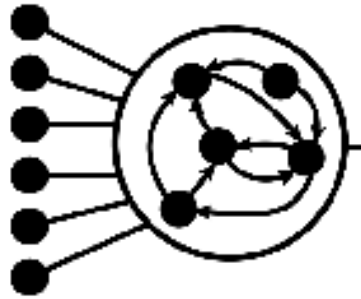
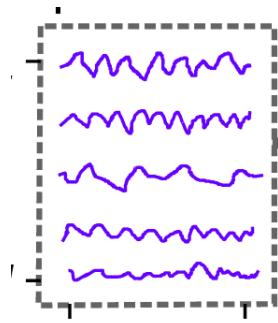


- recurrent connectivity
- nodal nonlinearity
- need to improve theory away from linear regime

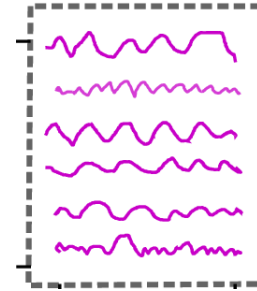


# Neuronal reservoir for time series processing

Input time series



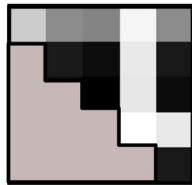
Reservoir time series



1st-order stats

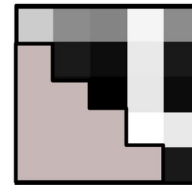
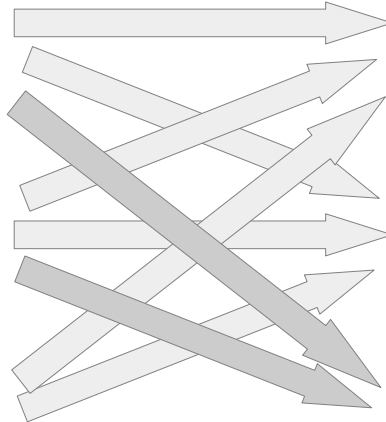


2nd-order stats



high-order stats

⋮



?

- recurrent connectivity
- nodal nonlinearity
- need to improve theory away from linear regime

# Outline

- **Statistical learning for time series**
  - mean versus covariance decoding
  - processing by neuronal reservoir
  - decoding by biological architecture
- Theory

• Structured variability  
conveys information

• **Recurrent connectivity +  
nonlinearity are key**

# Outline

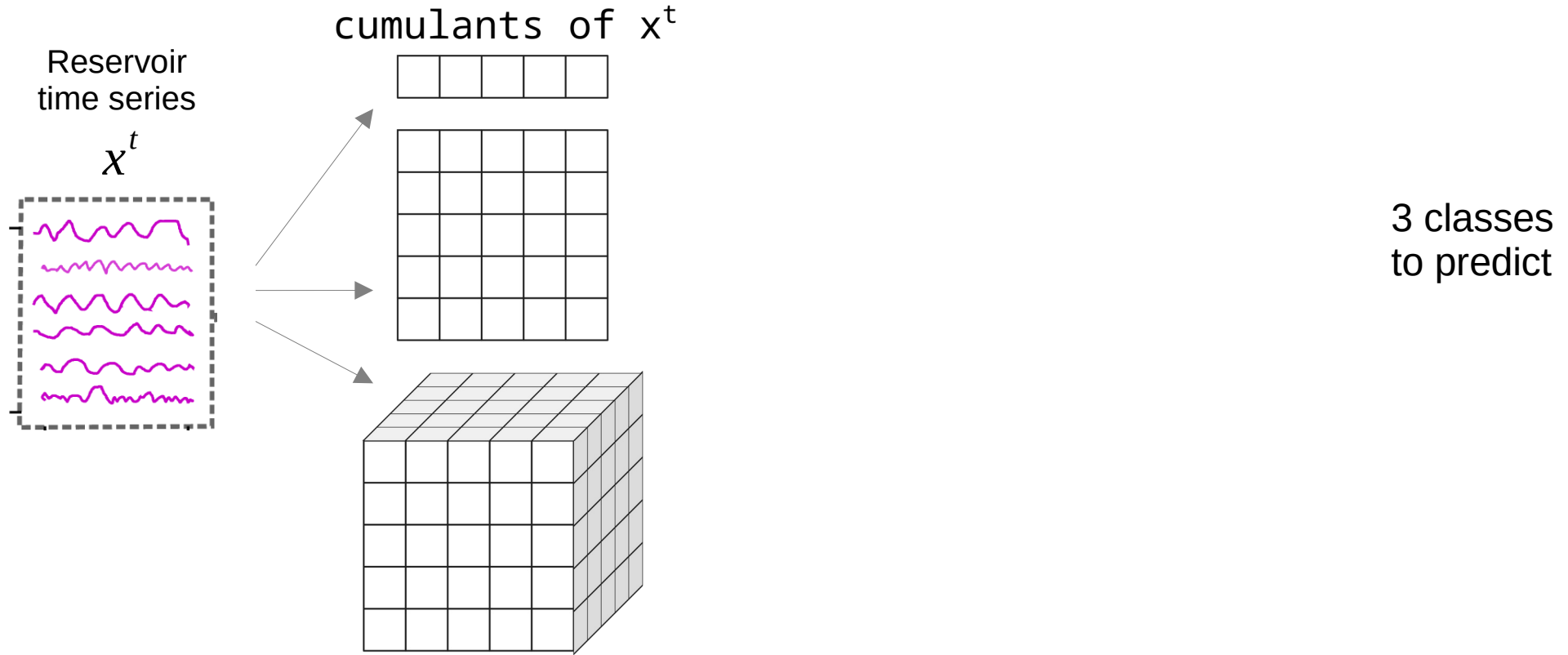
- **Statistical learning for time series**
  - mean versus covariance decoding
  - processing by neuronal reservoir
  - **decoding by biological architecture**
- Theory

• Structured variability conveys information

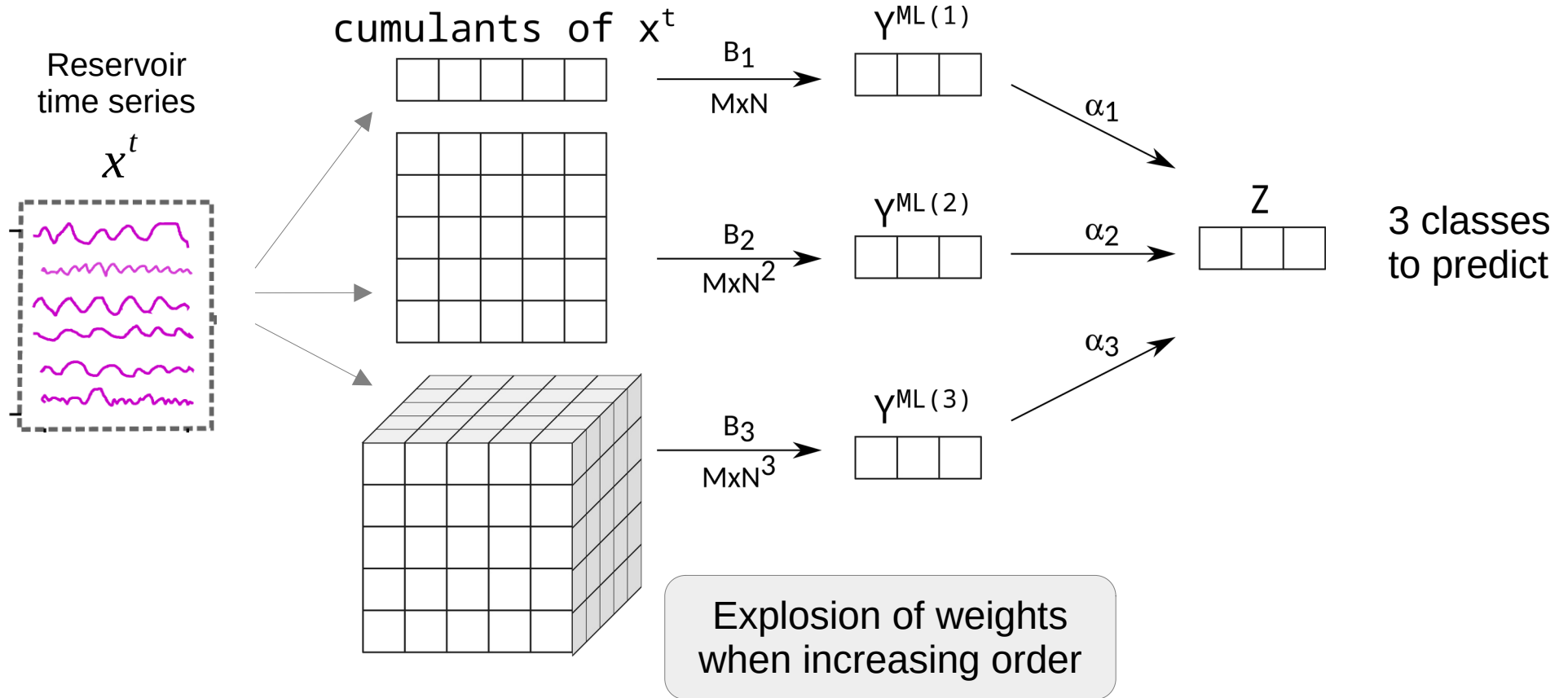
• Recurrent connectivity + nonlinearity are key

• **Extension to high-order statistics?**

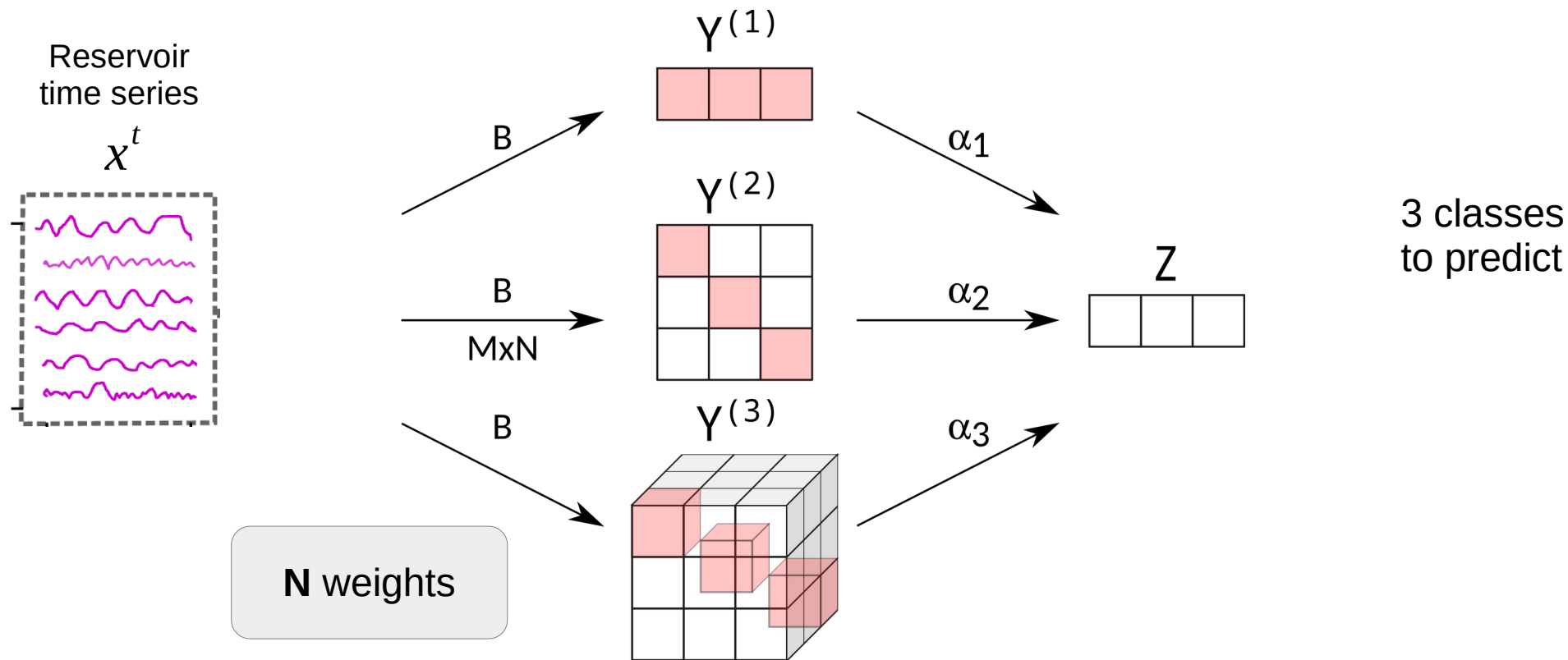
# Extension to statistical orders 1 to 3



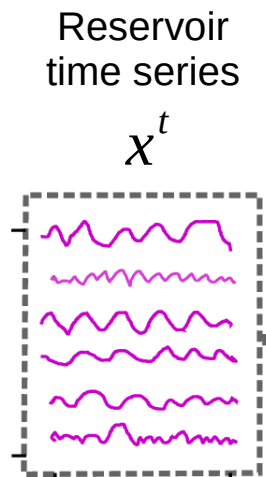
# Extension to statistical orders 1 to 3



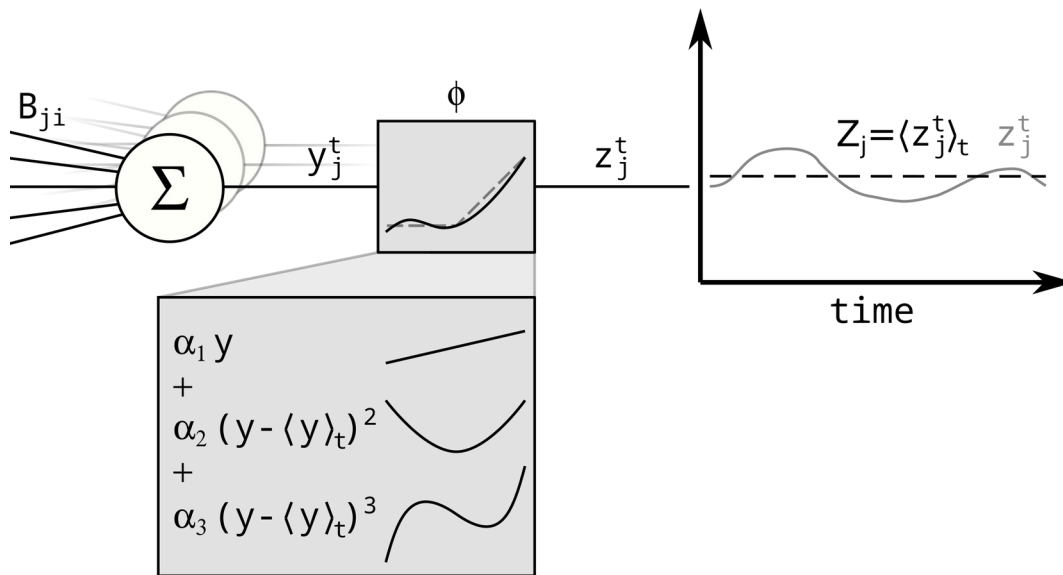
# Extension to statistical orders 1 to 3



# Order-selective perceptron (OSP)

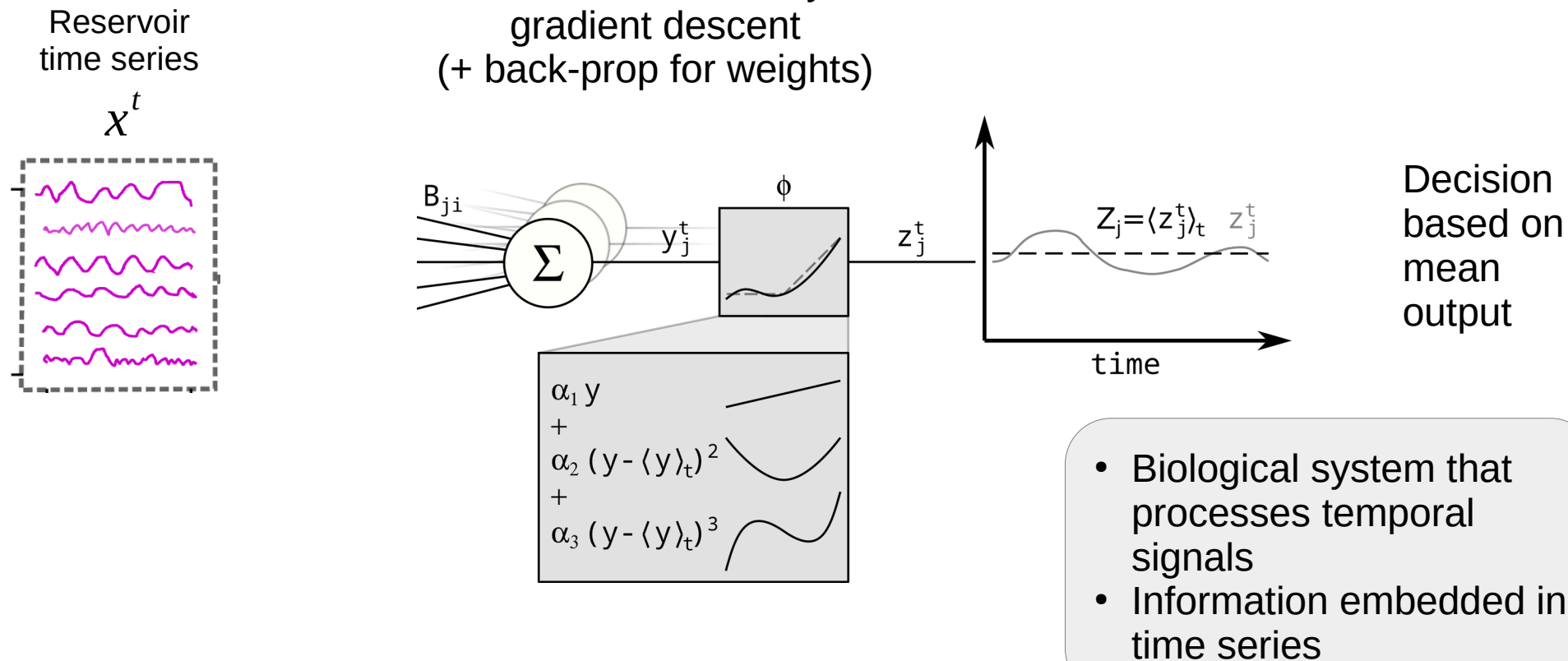


Tunable nonlinearity:  
gradient descent  
(+ back-prop for weights)



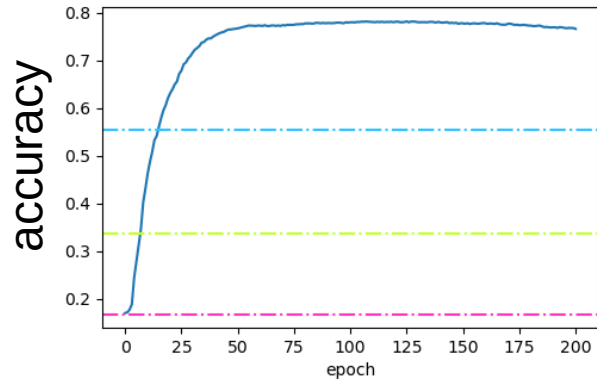
Decision  
based on  
mean  
output

# Order-selective perceptron (OSP)





# Interpretable learning

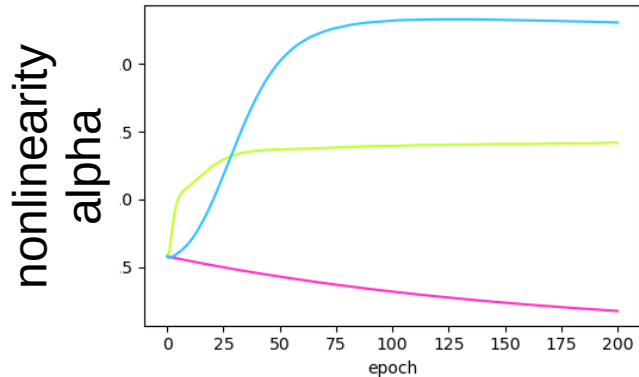


3<sup>rd</sup> order only

2<sup>nd</sup> order only

1<sup>st</sup> order only

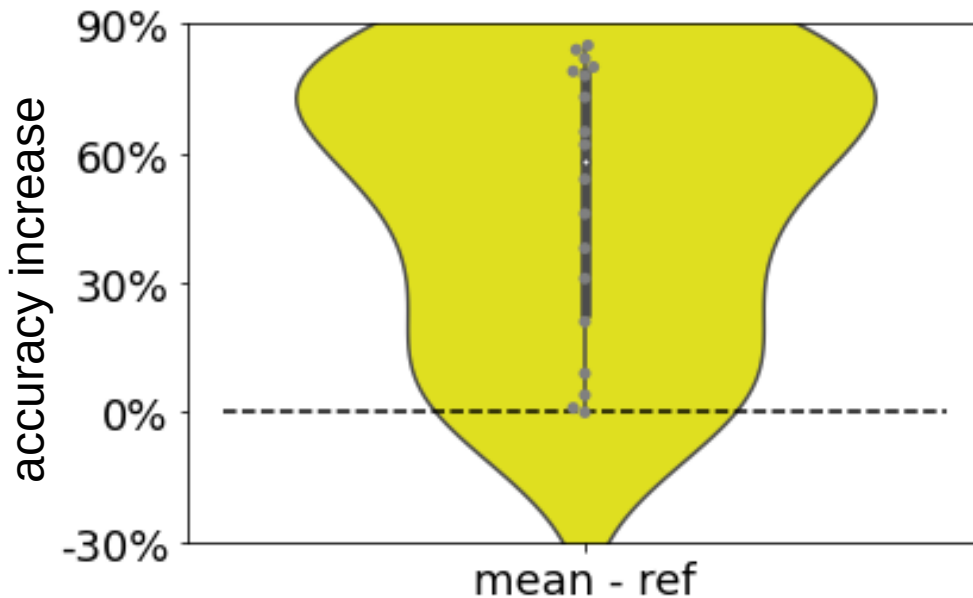
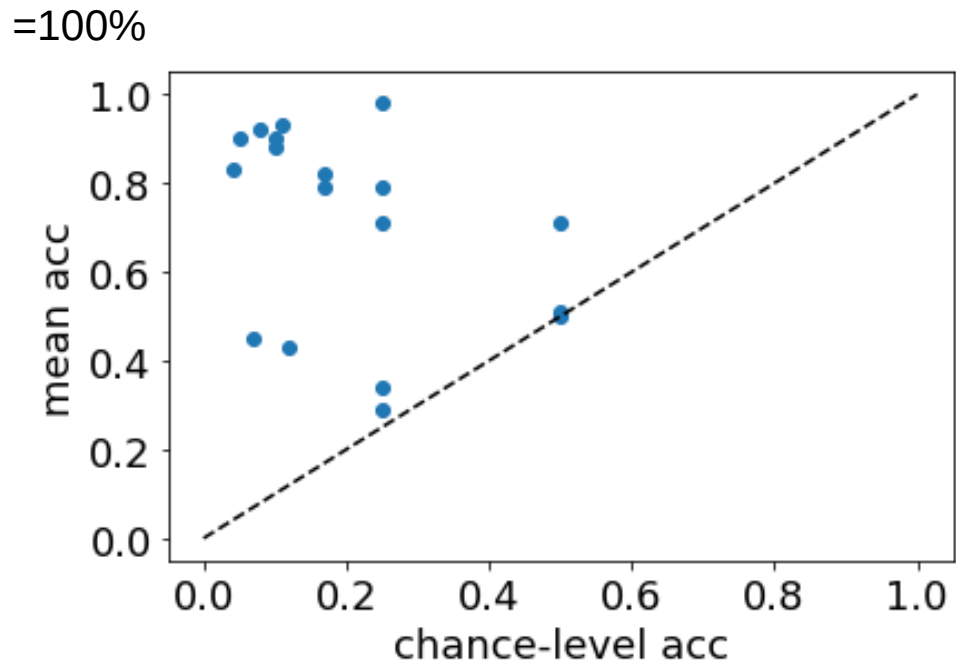
**Several orders are effectively used**



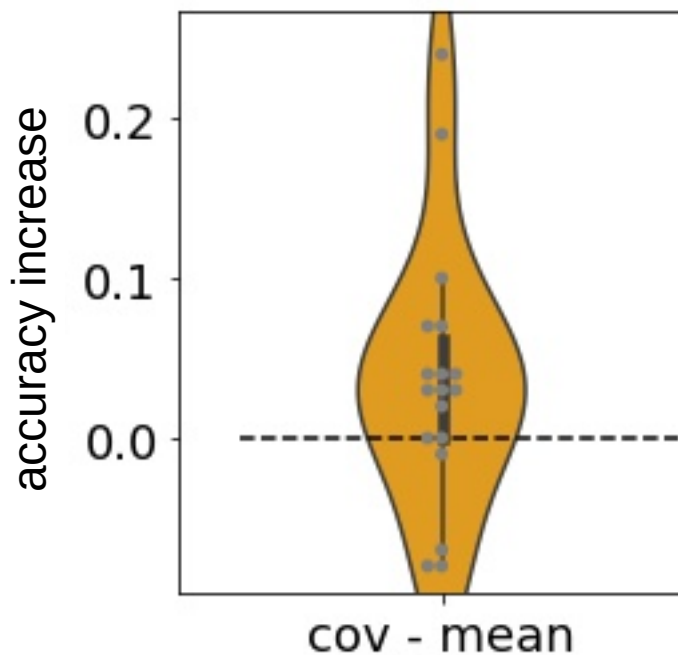
**Which order is relevant?**

# Comparison on real multivariate datasets

Hand movements (video), epilepsy (EEG), sensors, etc.

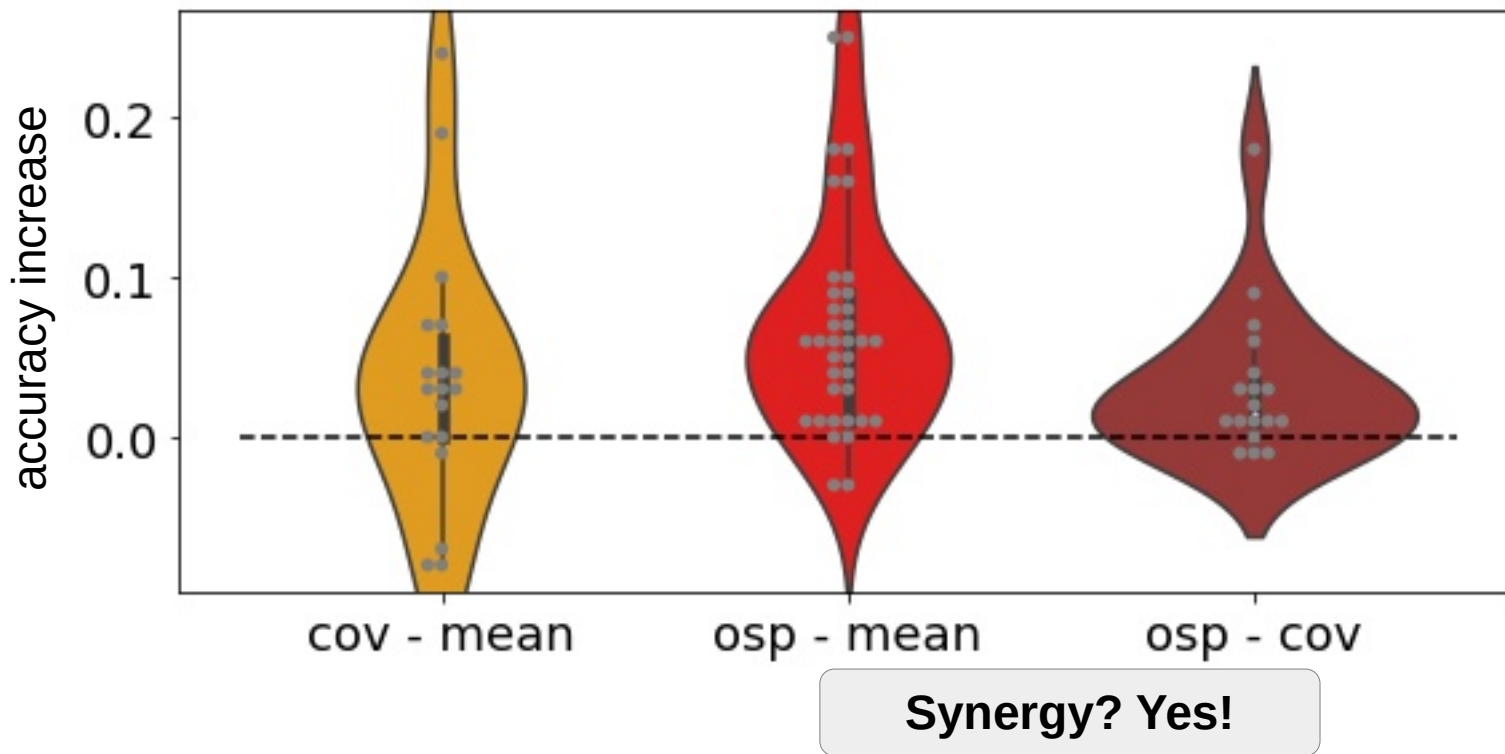


# Comparison on real multivariate datasets (reservoir + decoding)

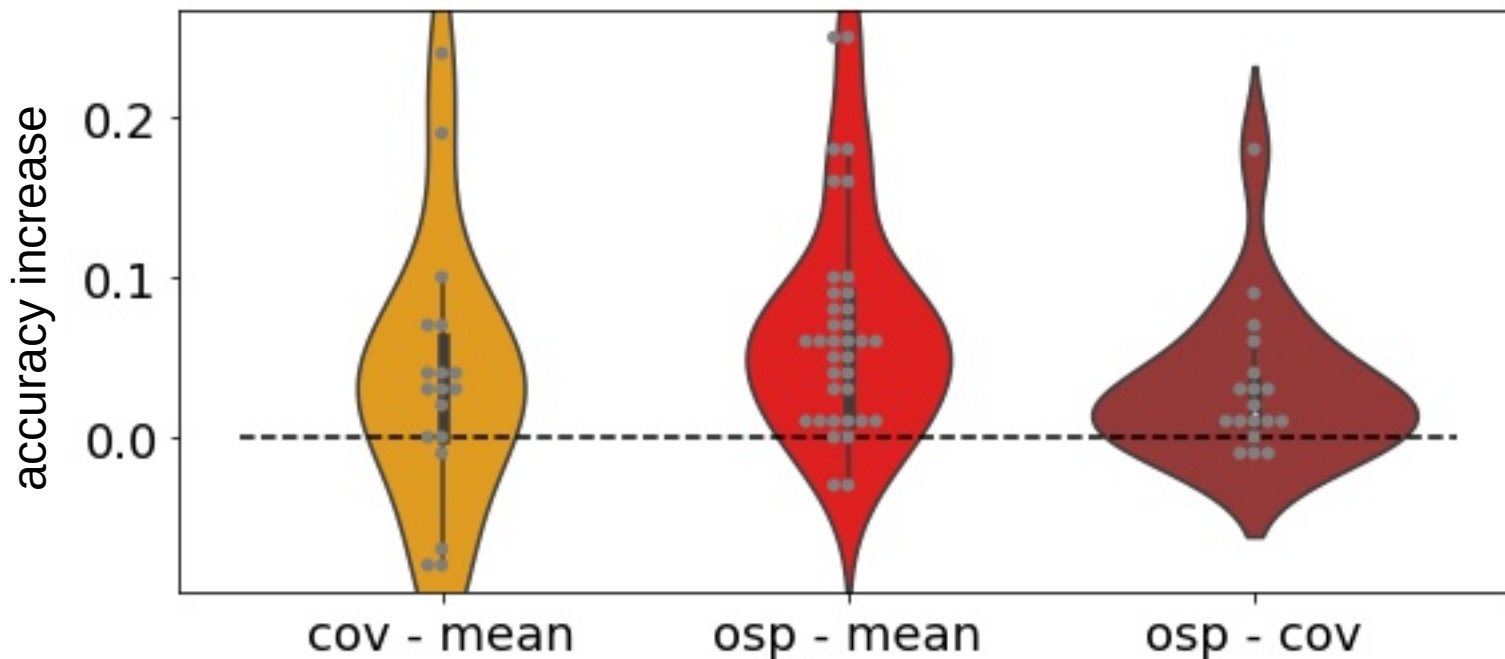


Covariance decoding  
outperforms  
mean decoding  
in most cases

# Comparison on real multivariate datasets (reservoir + decoding)



# Comparison on real multivariate datasets (reservoir + decoding)



**Synergy? Yes!**

**N versus  $N^2$  weights!**

# Outline

- **Statistical learning for time series**
  - mean versus covariance decoding
  - processing by neuronal reservoir
  - decoding by biological architecture
- Theory

- Structured variability conveys information

- Recurrent connectivity + nonlinearity are key

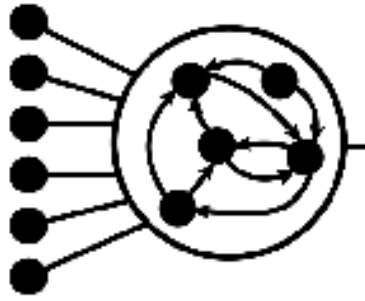
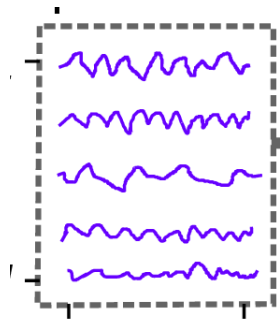
- **Robust decoding with limited resources**

**Next steps:**

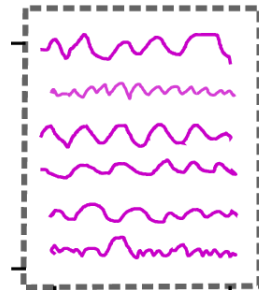
- **Trainable reservoir**
- **Deep net architecture**

# Conclusion on biological system

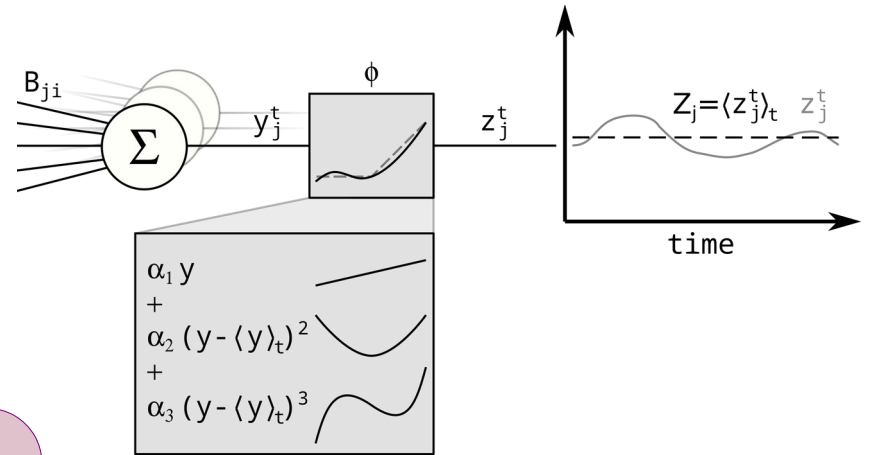
Input time series



Reservoir time series



Readout



- Rich representation
- Expansion and cross-talk across statistical orders
- Recurrent connectivity + nonlinearity

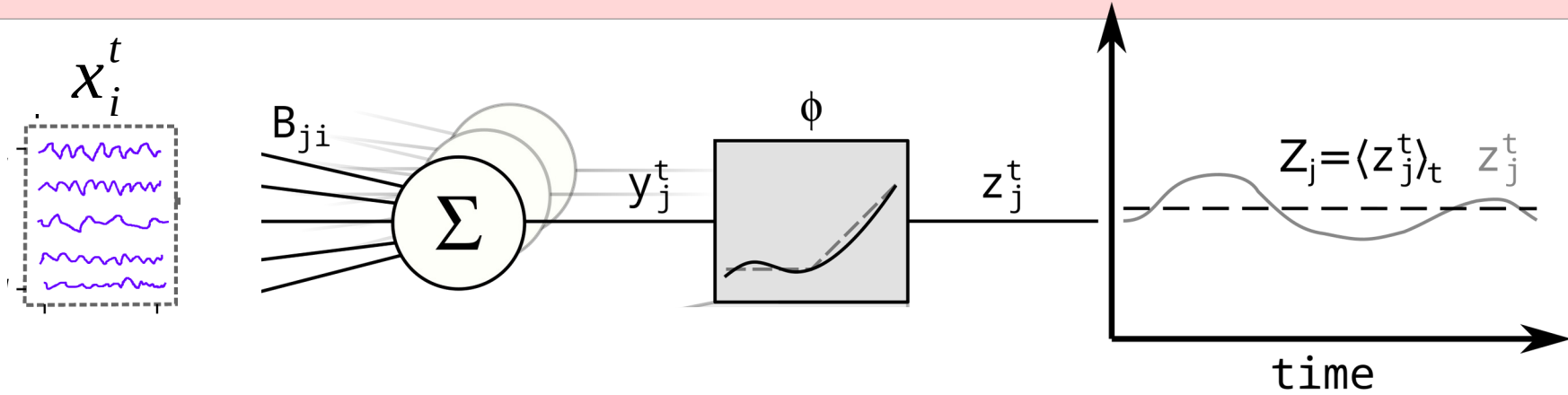
- Automated selection of relevant statistical orders
- Mapping to first-order in output

# Outline

- Statistical learning for time series
  - mean versus covariance decoding
  - processing by neuronal reservoir
  - decoding by biological architecture
- **Theory**

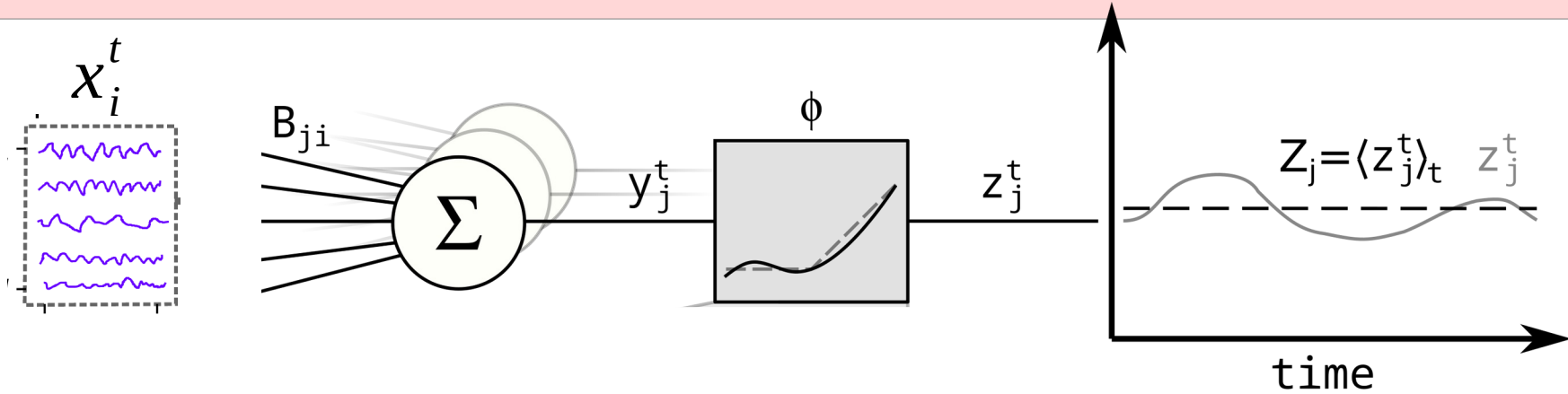


# Order-selective perceptron (OSP)



- Decision based on mean  $Z$
- Training with target  $\bar{Z}$

# Order-selective perceptron (OSP)



Equations for activation dynamics

$$y^t = \sum_i B_i x_i^t$$

$$z^t = \Phi(y^t) = \alpha_1 y^t + \alpha_2 (y^t - \langle y^t \rangle)^2 + \alpha_3 (y^t - \langle y^t \rangle)^3$$

## Training the nonlinearity

$$z^t = \alpha_1 y^t + \alpha_2 (y^t - \langle y^t \rangle)^2 + \alpha_3 (y^t - \langle y^t \rangle)^3$$

$$Z = \alpha_1 Y'^{1'} + \alpha_2 Y'^{2'} + \alpha_3 Y'^{3'} \quad \text{after averaging over observation window}$$

cumulants of orders 1, 2 and 3

# Training the nonlinearity

$$z^t = \alpha_1 y^t + \alpha_2 (y^t - \langle y^t \rangle)^2 + \alpha_3 (y^t - \langle y^t \rangle)^3$$

$$\mathbf{Z} = \alpha_1 \mathbf{Y}^{\prime 1} + \alpha_2 \mathbf{Y}^{\prime 2} + \alpha_3 \mathbf{Y}^{\prime 3} \quad \text{after averaging over observation window}$$

$$\text{Error } \epsilon = \frac{1}{2} \|\mathbf{Z} - \bar{\mathbf{Z}}\|^2 \quad \text{for target } \bar{\mathbf{Z}} \quad \text{when presenting activity } \mathbf{X}_i^t$$

# Training the nonlinearity

$$z^t = \alpha_1 y^t + \alpha_2 (y^t - \langle y^t \rangle)^2 + \alpha_3 (y^t - \langle y^t \rangle)^3$$

$$Z = \alpha_1 Y'^{1'} + \alpha_2 Y'^{2'} + \alpha_3 Y'^{3'} \quad \text{after averaging over observation window}$$

Error  $\epsilon = \frac{1}{2} \|Z - \bar{Z}\|^2$  for target  $\bar{Z}$  when presenting activity  $X_i^t$

$$\Delta \alpha_k = - \frac{\partial \epsilon}{\partial Z} \frac{\partial Z}{\partial \alpha_k} = (\bar{Z} - Z) Y'^{k'} \quad \text{using} \quad \frac{\partial Z}{\partial \alpha_k} = Y'^{k'}$$

chain rule

# Training the afferent weights using back-propagation

$$Z = \alpha_1 Y'^{1'} + \alpha_2 Y'^{2'} + \alpha_3 Y'^{3'} \quad \text{and} \quad y^t = \sum_i B_i x_i^t$$

so

$$\frac{\partial Z}{\partial B} = \sum_k \alpha_k \frac{\partial Y'^{k'}}$$

# Training the afferent weights using back-propagation

$$Z = \alpha_1 Y'^{1'} + \alpha_2 Y'^{2'} + \alpha_3 Y'^{3'} \quad \text{and} \quad y^t = \sum_i B_i x_i^t$$

so

$$\frac{\partial Z}{\partial B} = \sum_k \alpha_k \frac{\partial Y'^{k'}}$$

we obtain

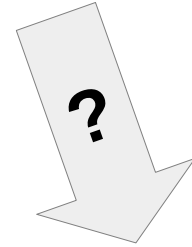
$$\Delta B = - \frac{\partial \epsilon}{\partial Z} \frac{\partial Z}{\partial B} = (\bar{Z} - Z) \sum_k \alpha_k \frac{\partial Y'^{k'}}$$

# Training the afferent weights using back-propagation

$$Z = \alpha_1 Y'^{1'} + \alpha_2 Y'^{2'} + \alpha_3 Y'^{3'} \quad \text{and} \quad y^t = \sum_i B_i x_i^t$$

so

$$\frac{\partial Z}{\partial B} = \sum_k \alpha_k \frac{\partial Y'^{k'}}$$

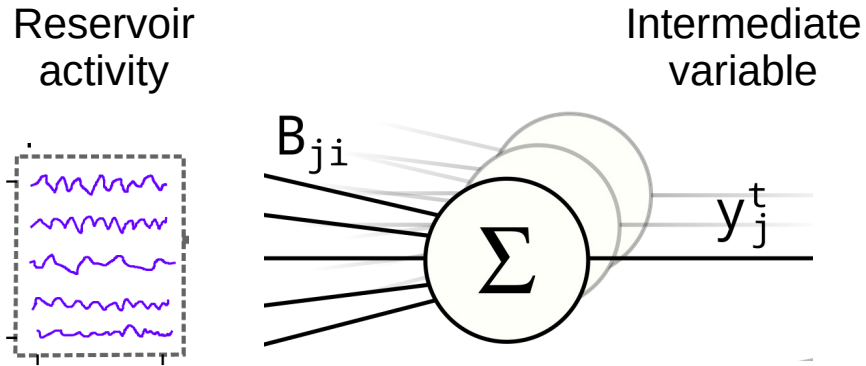


we obtain

$$\Delta B = - \frac{\partial \epsilon}{\partial Z} \frac{\partial Z}{\partial B} = (\bar{Z} - Z) \sum_k \alpha_k \frac{\partial Y'^{k'}}$$



# Derivatives of cumulants with respect to afferent weights



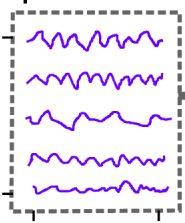
$$x_i^t$$

$$y_j^t = \sum_i B_{ji} x_i^t$$

Linear perceptron

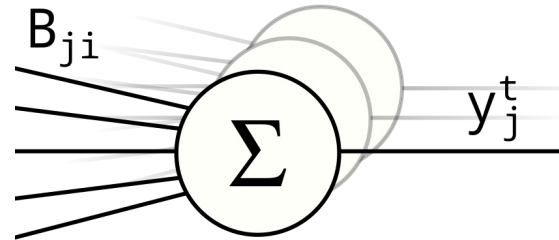
# Derivatives of cumulants with respect to afferent weights

Reservoir activity



$$x_i^t$$

Intermediate variable



$$y_j^t = \sum_i B_{ji} x_i^t$$

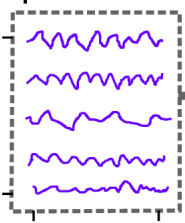
Linear perceptron

Temporal average

$$Y'^{1'} = \langle y^t \rangle = B X'^{1'}$$

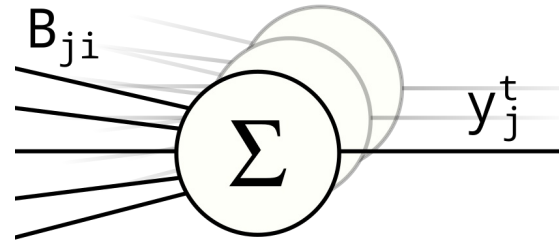
# Derivatives of cumulants with respect to afferent weights

Reservoir activity



$$x_i^t$$

Intermediate variable



$$y_j^t = \sum_i B_{ji} x_i^t$$

Linear perceptron

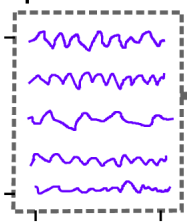
Temporal average

$$Y'^{1'} = \langle y^t \rangle = B X'^{1'}$$

$$\frac{\partial Y'^{1'}}{\partial B} = X'^{1'}$$

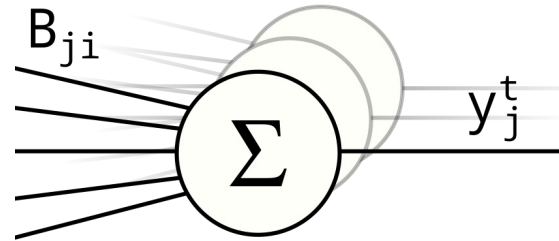
# Derivatives of cumulants with respect to afferent weights

Reservoir activity



$$x_i^t$$

Intermediate variable



$$y_j^t = \sum_i B_{ji} x_i^t$$

Linear perceptron

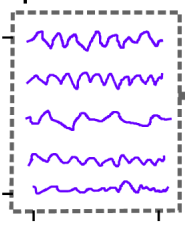
Temporal covariance

$$Y'^{2'} = \langle y^t y^t \rangle = B X'^{2'} B^T$$

after centering

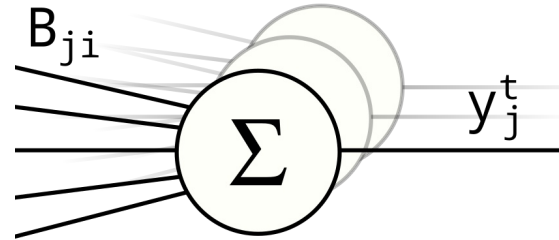
# Derivatives of cumulants with respect to afferent weights

Reservoir activity



$$X_i^t$$

Intermediate variable



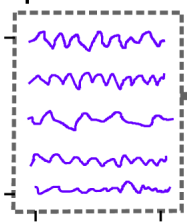
Temporal covariance

$$Y'^{2'} = \langle y^t y^t \rangle = B X'^{2'} B^T$$

$$\frac{\partial Y'^{2'}}{\partial B} = B X'^{2'} + X'^{2'} B^T$$

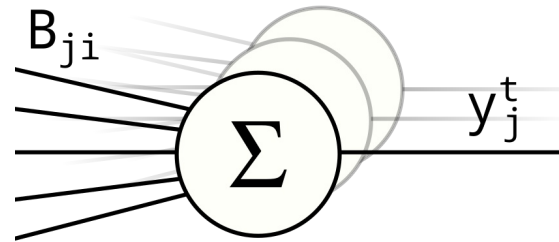
# Derivatives of cumulants with respect to afferent weights

Reservoir activity



$$X_i^t$$

Intermediate variable

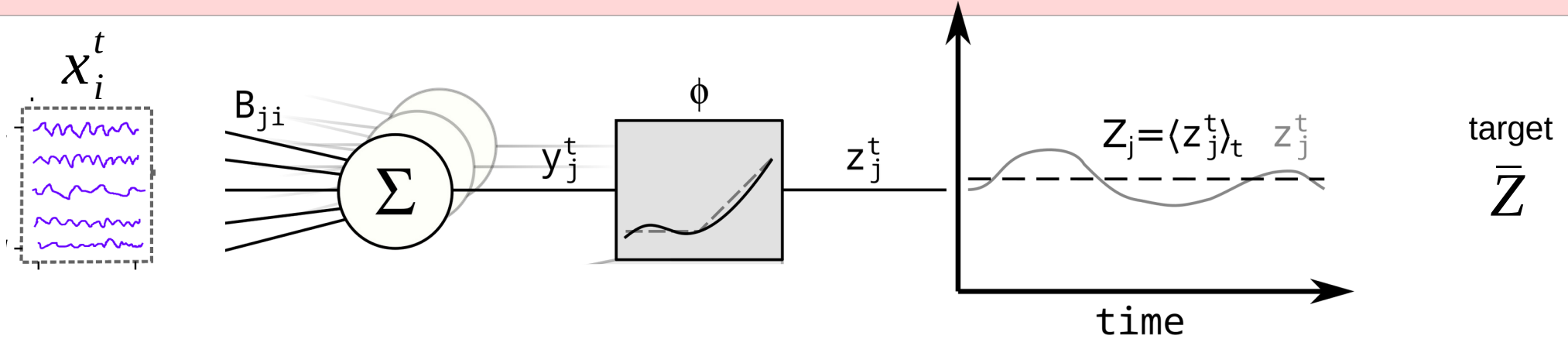


Third-order covariance

$$Y'^{3'} = \langle y^t y^t y^t \rangle \\ = B \odot^1 B \odot^2 B \odot^3 X'^{3'}$$

$$\frac{\partial Y'^{3'}}{\partial B} = B \odot^1 B \odot^2 X'^{3'} + \dots$$

# Putting everything together: Automated selection of relevant input cumulant for classification



$$\Delta \alpha_k = (\bar{Z} - Z) Y'^{k'}$$

$$\Delta B = (\bar{Z} - Z) \sum_k \alpha_k [B, \dots, B] \odot X'^{k'}$$

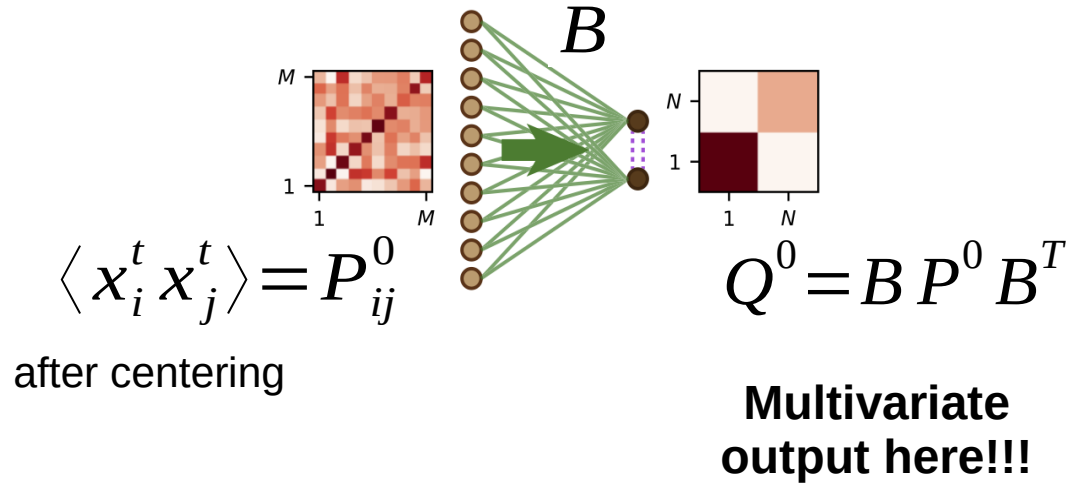
depends on k

# Weight update in linear network to tune output covariance

## Covariance mapping

Target

Error (matrix norm)

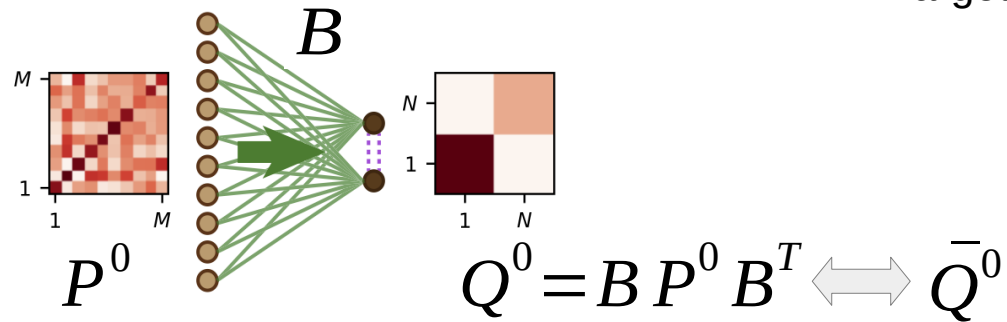


Weight update



# Weight update in linear network to tune output covariance

## Covariance mapping



Target

Error (matrix norm)

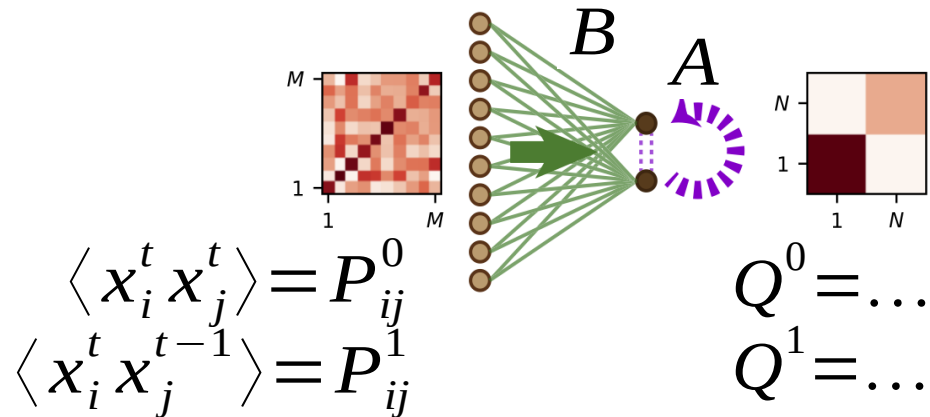
$$\epsilon = \|\bar{Q}^0 - Q^0\|^2$$

Weight update

$$\Delta B = - \frac{\partial \epsilon}{\partial Q^0} \frac{\partial Q^0}{\partial B}$$

# Naturally extends to recurrent networks

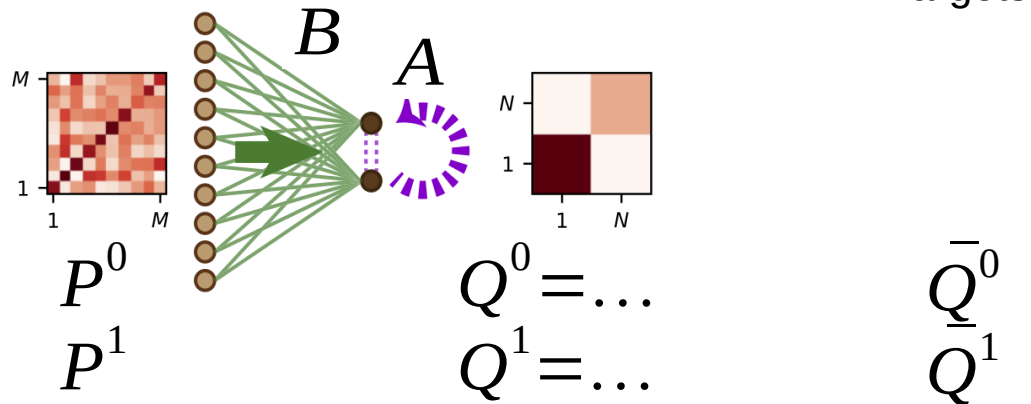
## Covariance mapping



$$y_j^t = \sum_i A_{ji} y_i^{t-1} + \sum_i B_{ji} x_i^t$$

# Naturally extends to recurrent networks

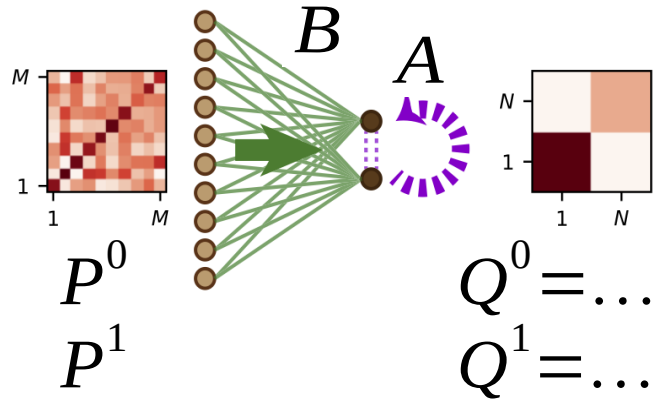
## Covariance mapping



$$y_j^t = \sum_i A_{ji} y_i^{t-1} + \sum_i B_{ji} x_i^t$$

# Naturally extends to recurrent networks

## Covariance mapping



Targets

Error (matrix norm)

$$\begin{matrix} \bar{Q}^0 \\ \bar{Q}^1 \end{matrix}$$

$$\epsilon = \|\bar{Q}^0 - Q^0\|^2 + \|\bar{Q}^1 - Q^1\|^2$$

Weight updates

$$\Delta B = -\frac{\partial \epsilon}{\partial Q^0} \frac{\partial Q^0}{\partial B} - \frac{\partial \epsilon}{\partial Q^1} \frac{\partial Q^1}{\partial B}$$

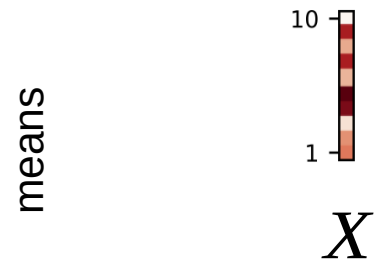
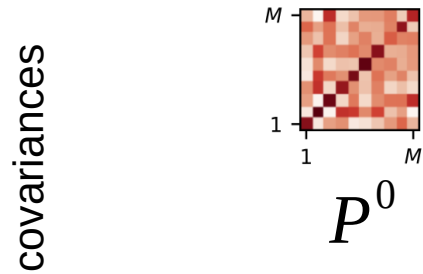
$$\Delta A = -\frac{\partial \epsilon}{\partial Q^0} \frac{\partial Q^0}{\partial A} - \frac{\partial \epsilon}{\partial Q^1} \frac{\partial Q^1}{\partial A}$$

# Summary on theory

- Training nonlinearity and afferent input weights to tune output mean
- Training afferent and recurrent weights in linear network to tune output covariances (also extends to higher orders)
- Now: tune output covariances when including nonlinearity in network

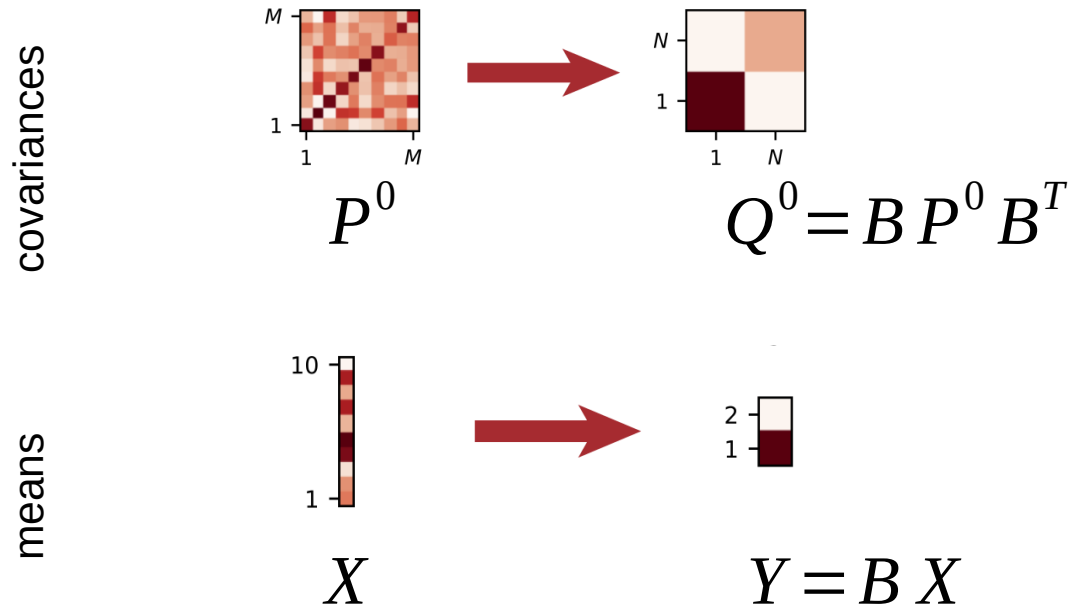
# Capacity of covariance perceptron

- How many binary patterns can be categorized into 2 groups?
- Feedforward linear perceptron



# Capacity of covariance perceptron

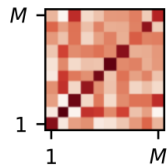
- How many binary patterns can be categorized into 2 groups?
- Feedforward linear perceptron



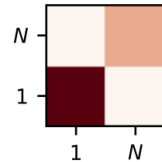
# Capacity of covariance perceptron

- How many binary patterns can be categorized into 2 groups?
- Feedforward linear perceptron

covariances

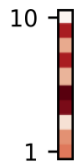


$$P^0$$



$$Q^0 = B P^0 B^T$$

means



$$X$$



$$Y = B X$$

Replica method

