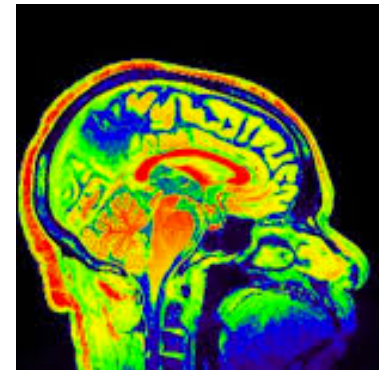
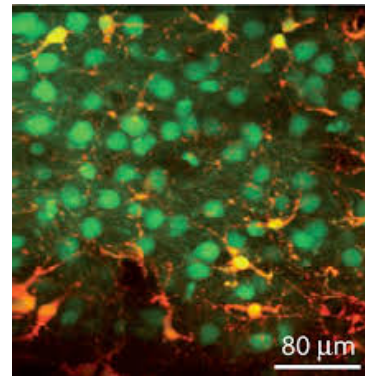
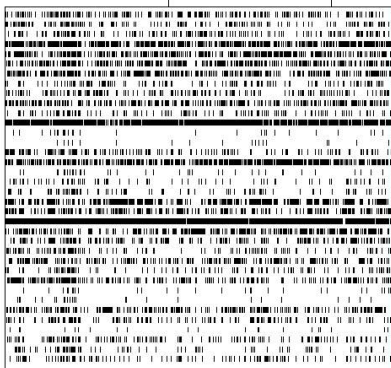


# Statistical Models for Neural Data: from Regression / GLMs to Latent Variables

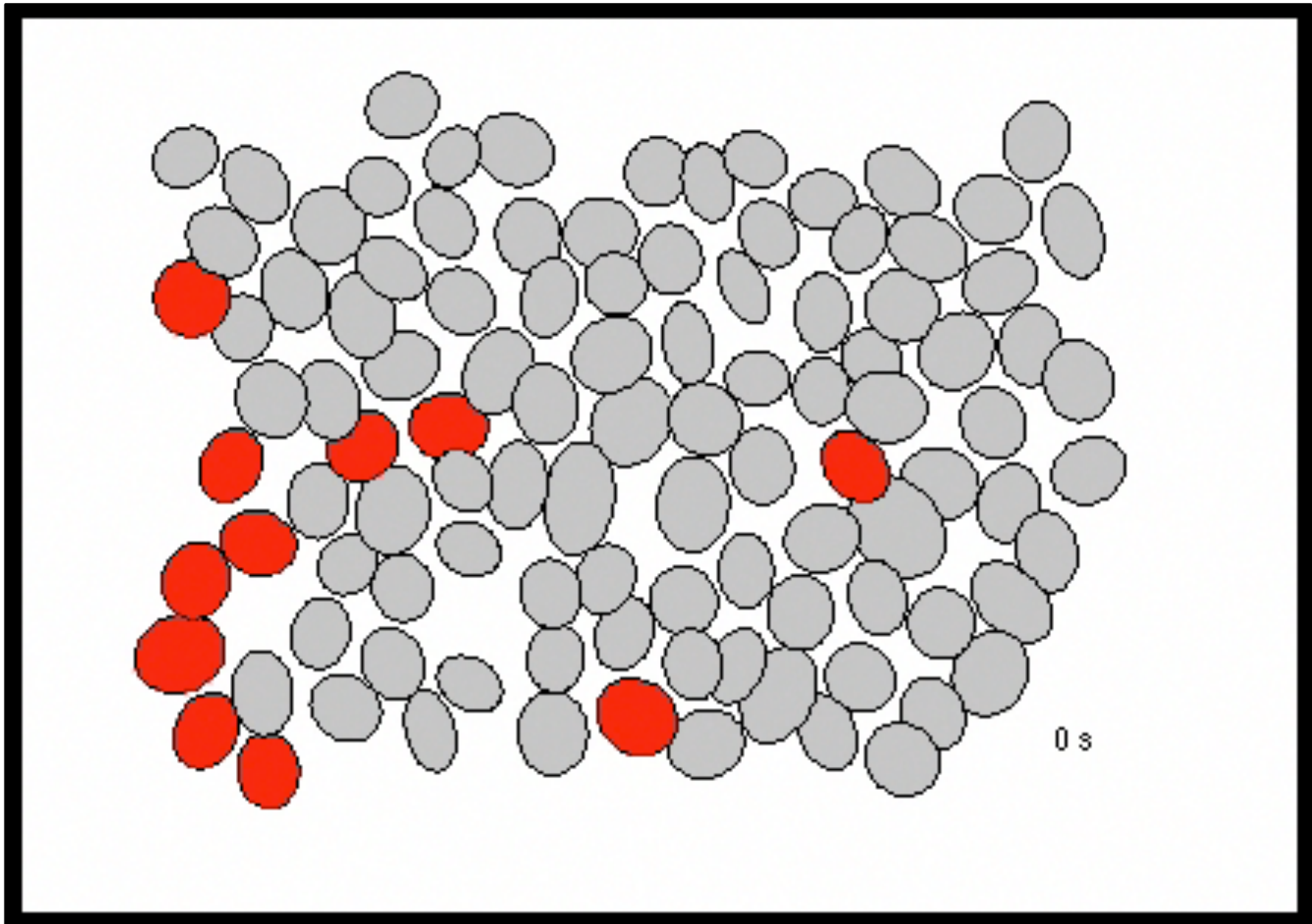
Jonathan Pillow  
*Princeton Neuroscience Institute*



Tutorial  
Cosyne 2018

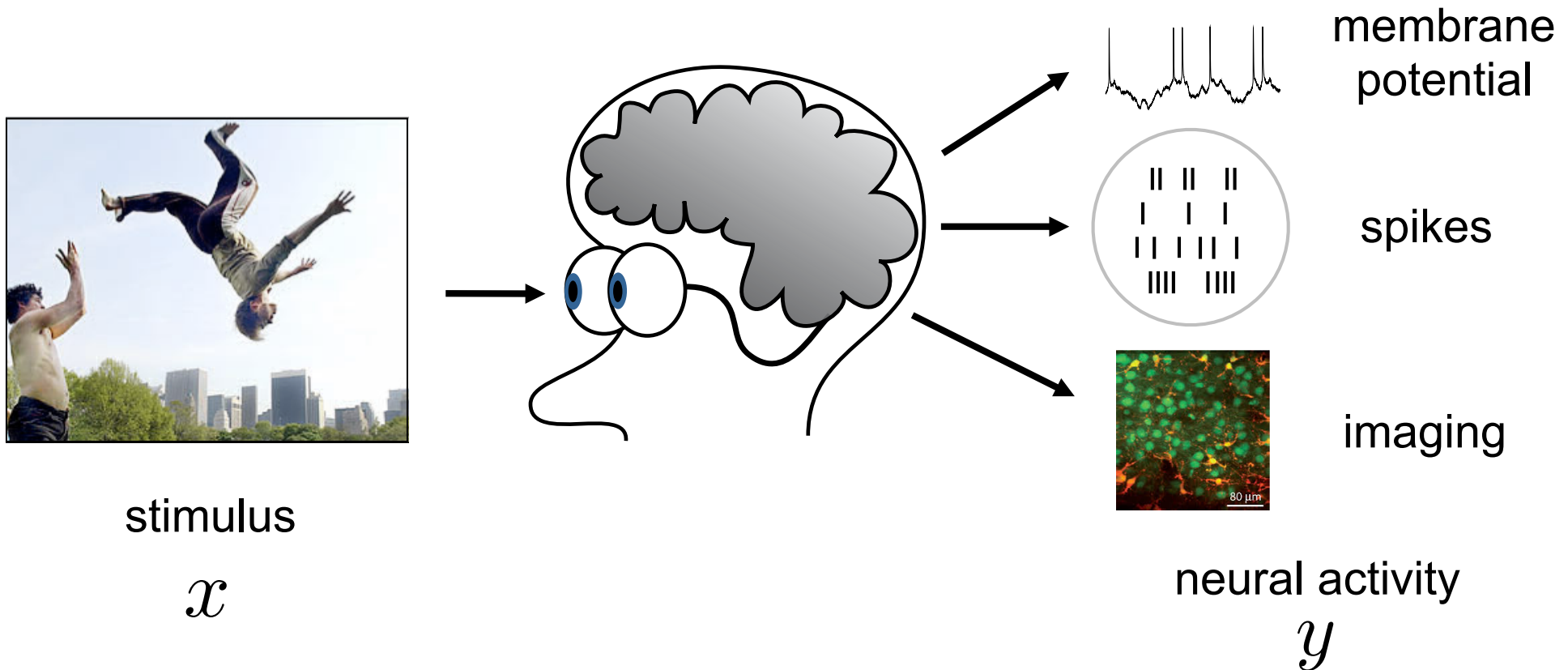
# Retinal responses to white noise

(ON parasol cells)



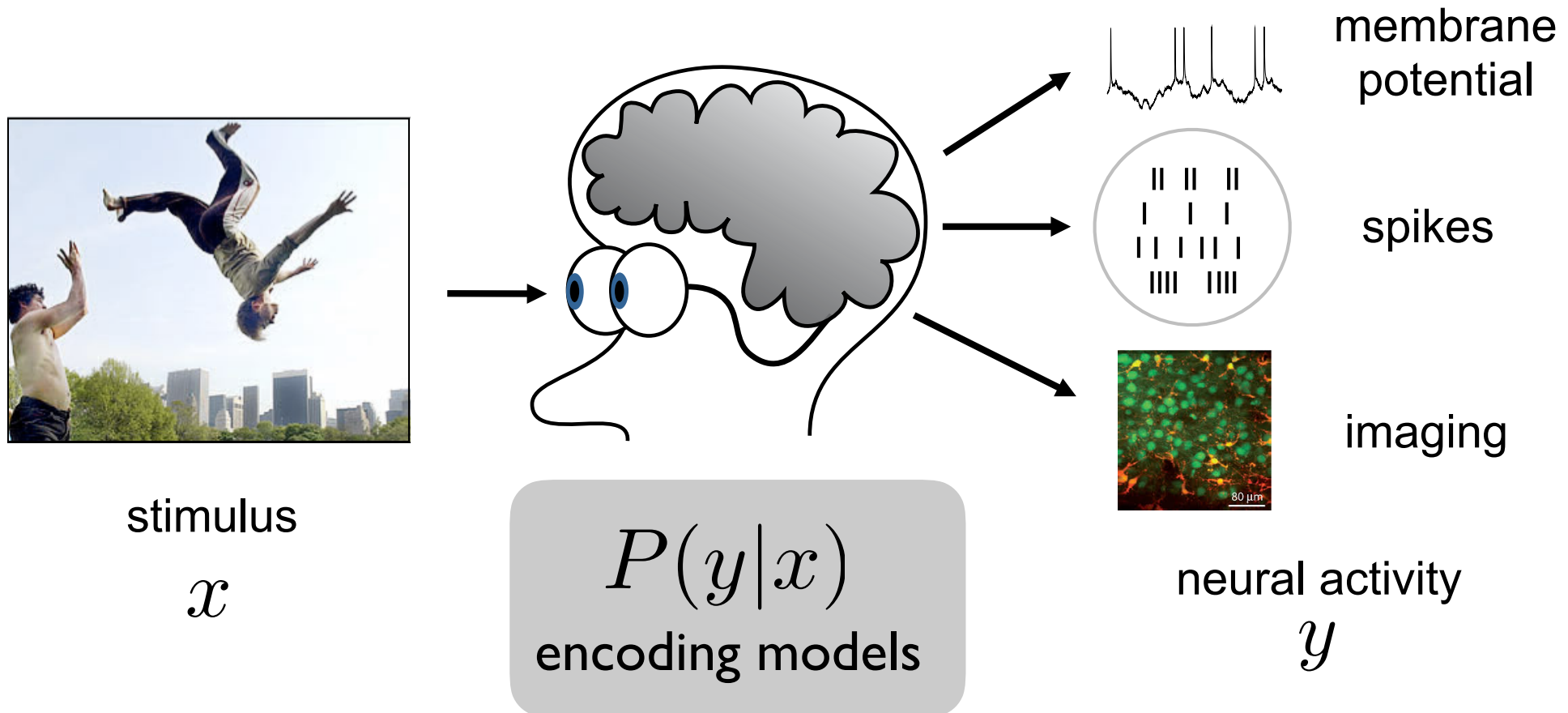
Shlens, Field, Gauthier, Greschner, Sher, Litke & Chichilnisky (2009).

# neural coding problem



- How are stimuli and actions encoded in neural activity?
- What aspects of neural activity carry information?

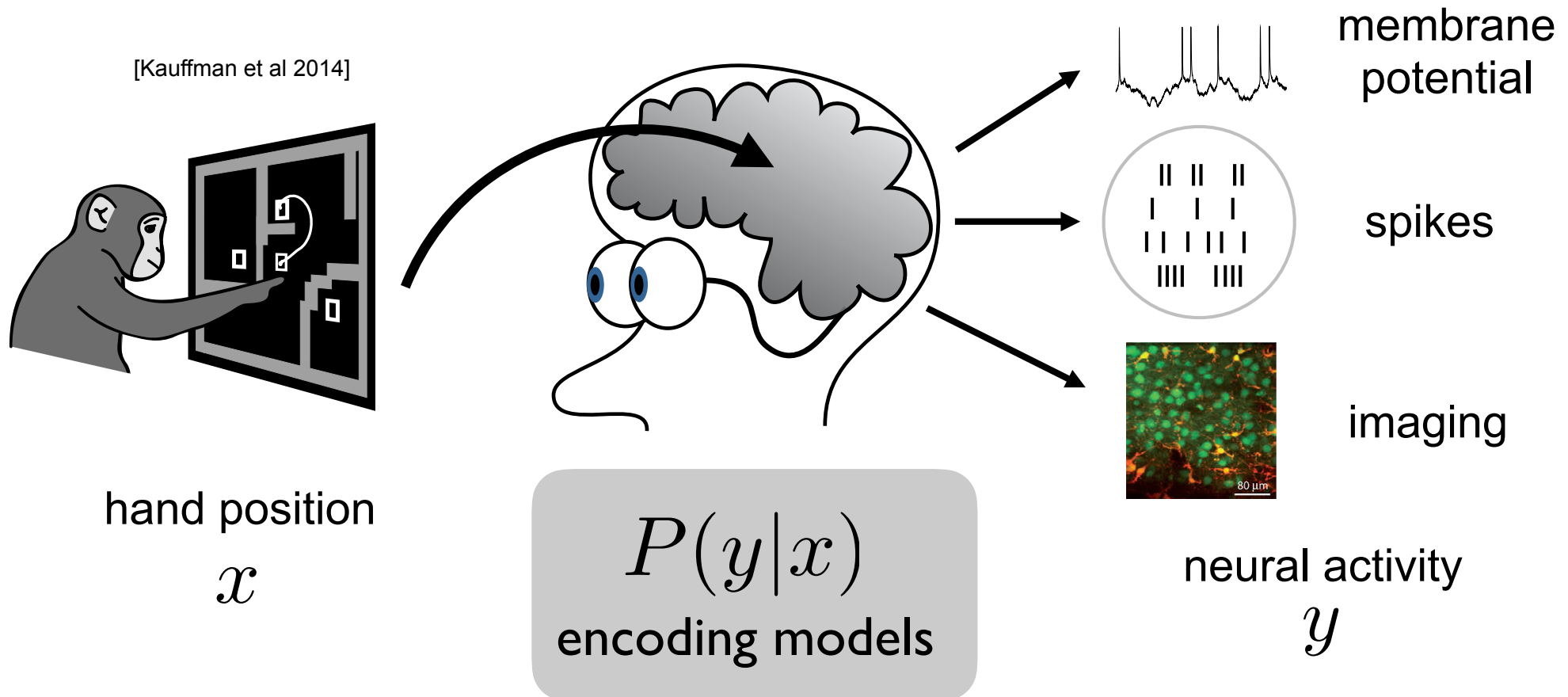
# neural coding problem



- Approach:**
- develop flexible statistical models of  $P(y|x)$
  - quantify information carried in neural responses



# neural coding problem

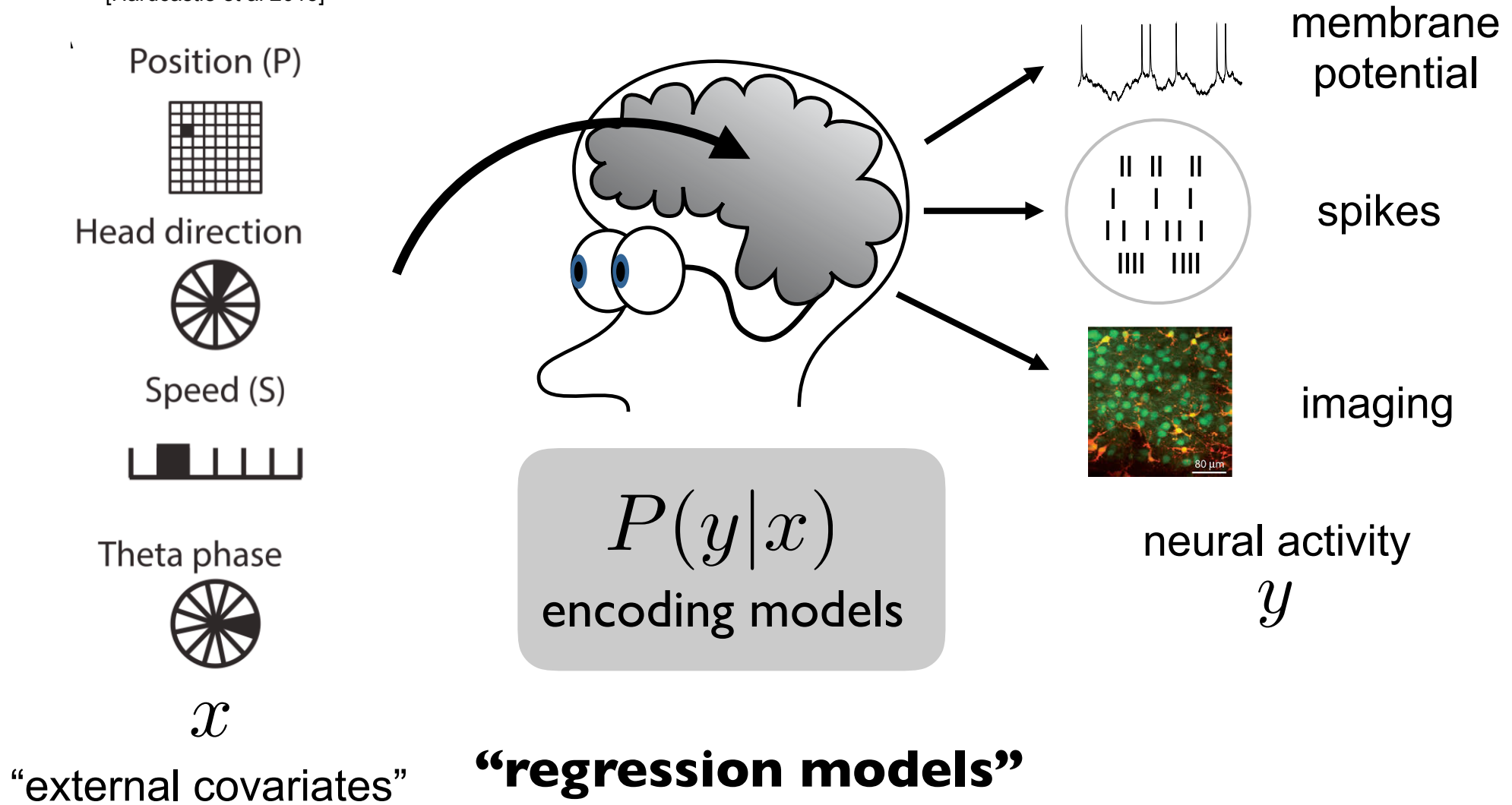


## “regression models”

- not restricted to sensory variables

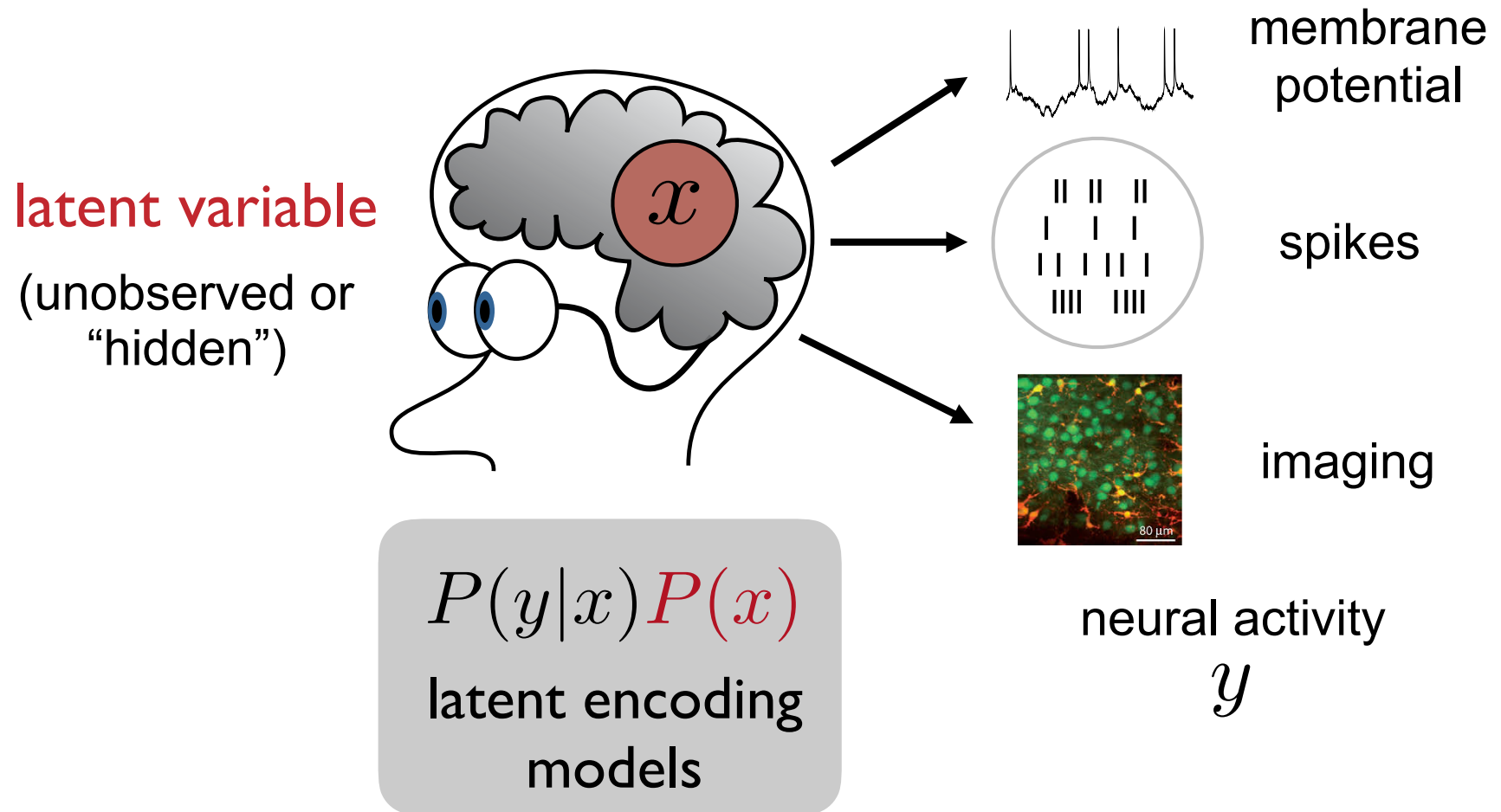
# neural coding problem

[Hardcastle et al 2015]



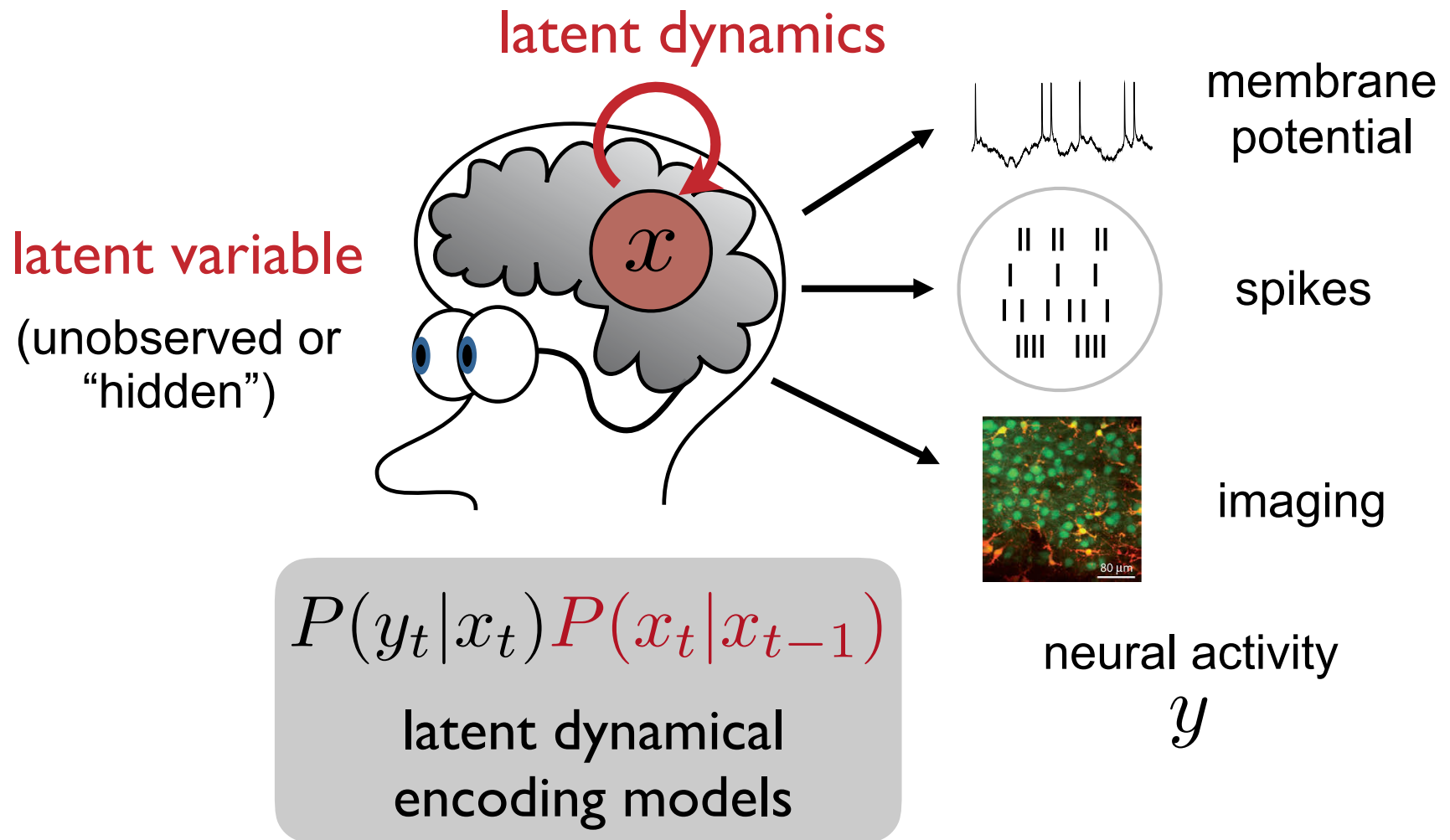
- not restricted to sensory variables

# latent variable models



- capture hidden structure underlying neural activity  
(eg. low-dimensional or discrete states)

# latent variable models



- capture hidden dynamics underlying neural activity

# model desiderata

linear,  
Gaussian

GLM

**sweet  
spot**

multi-  
compartment  
Hodgkin-Huxley



**fittability /  
tractability**

(can be fit to data)

**richness /  
flexibility**

(capture realistic neural  
properties)

normative theories  
(e.g. “efficient coding”)

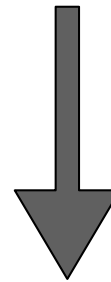
*Why does the code  
take this form?*



**descriptive  
statistical models**

$$P(y|x)$$

*What is the code?*



anatomy,  
biophysics

*How is it implemented?*



# Outline

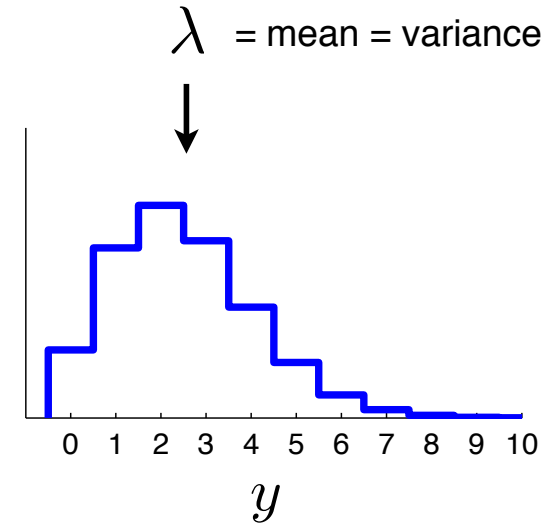
1. Spike count models & Maximum Likelihood
2. Spike train models (GLMs with spike history)
3. Multiple Spike Train Models (GLMs with coupling)
4. Regularization
5. Beyond GLM
6. Latent variable models

# simple example #1: linear Poisson neuron

parameter stimulus

spike rate  $\lambda = \theta x$

spike count  $y \sim \text{Pois}(\lambda)$

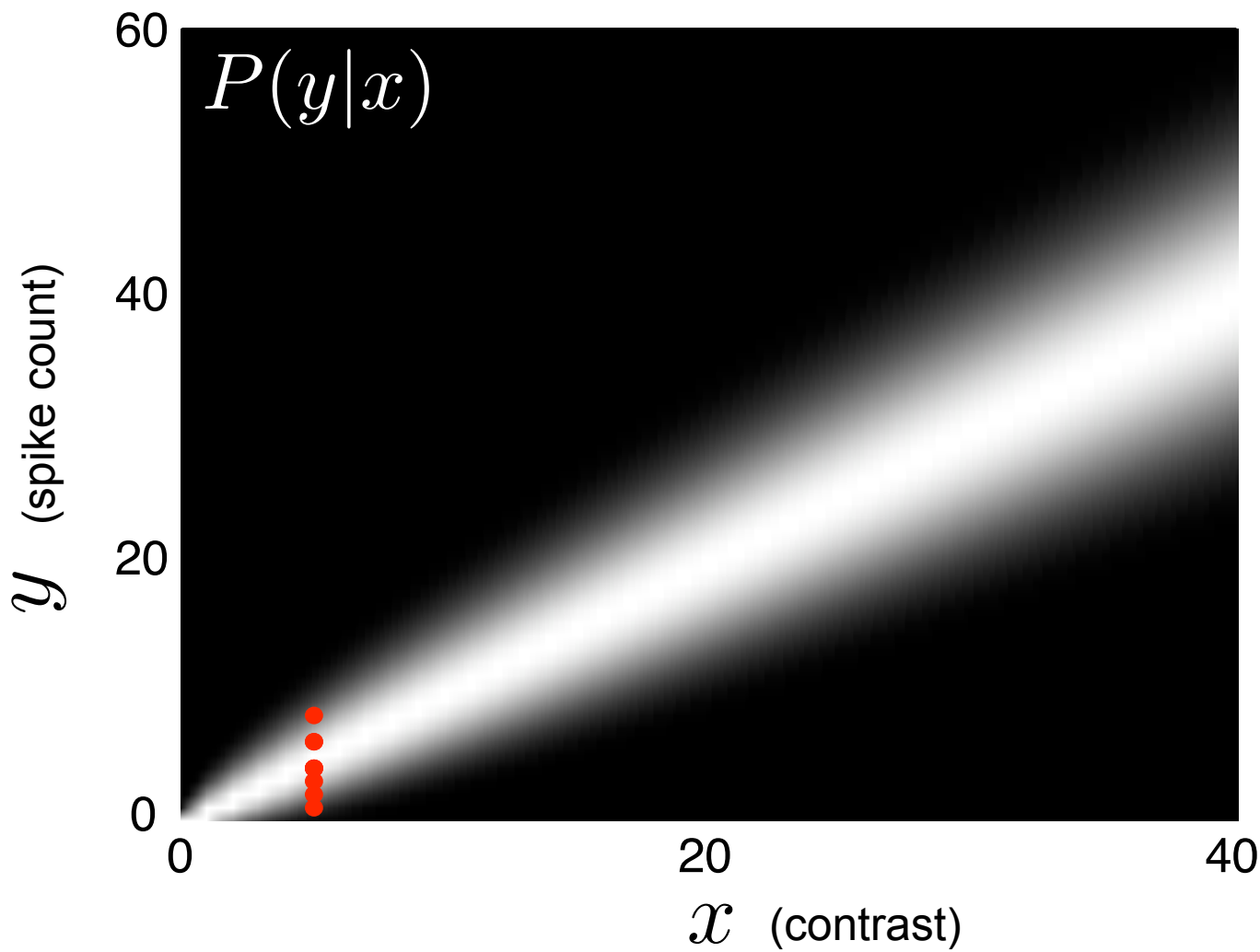


encoding model:

$$P(y|x, \theta) = \frac{1}{y!} \lambda^y e^{-\lambda}$$
$$= \frac{1}{y!} (\theta x)^y e^{-(\theta x)}$$

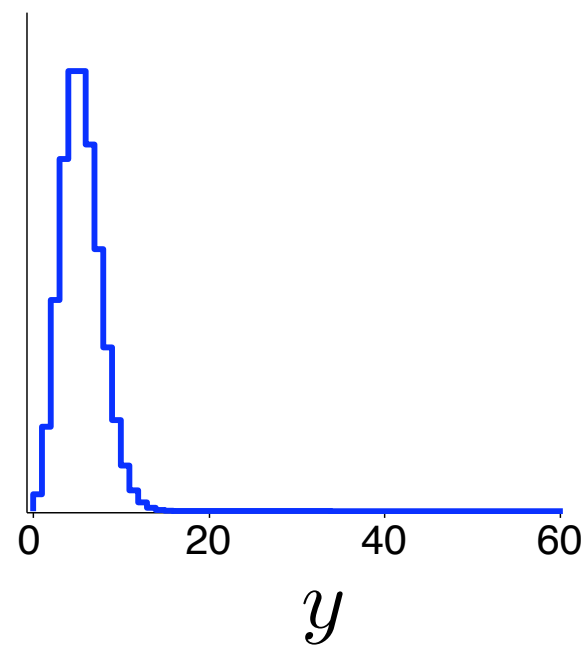
$$\text{mean}(y) = \theta x$$

$$\text{var}(y) = \theta x$$



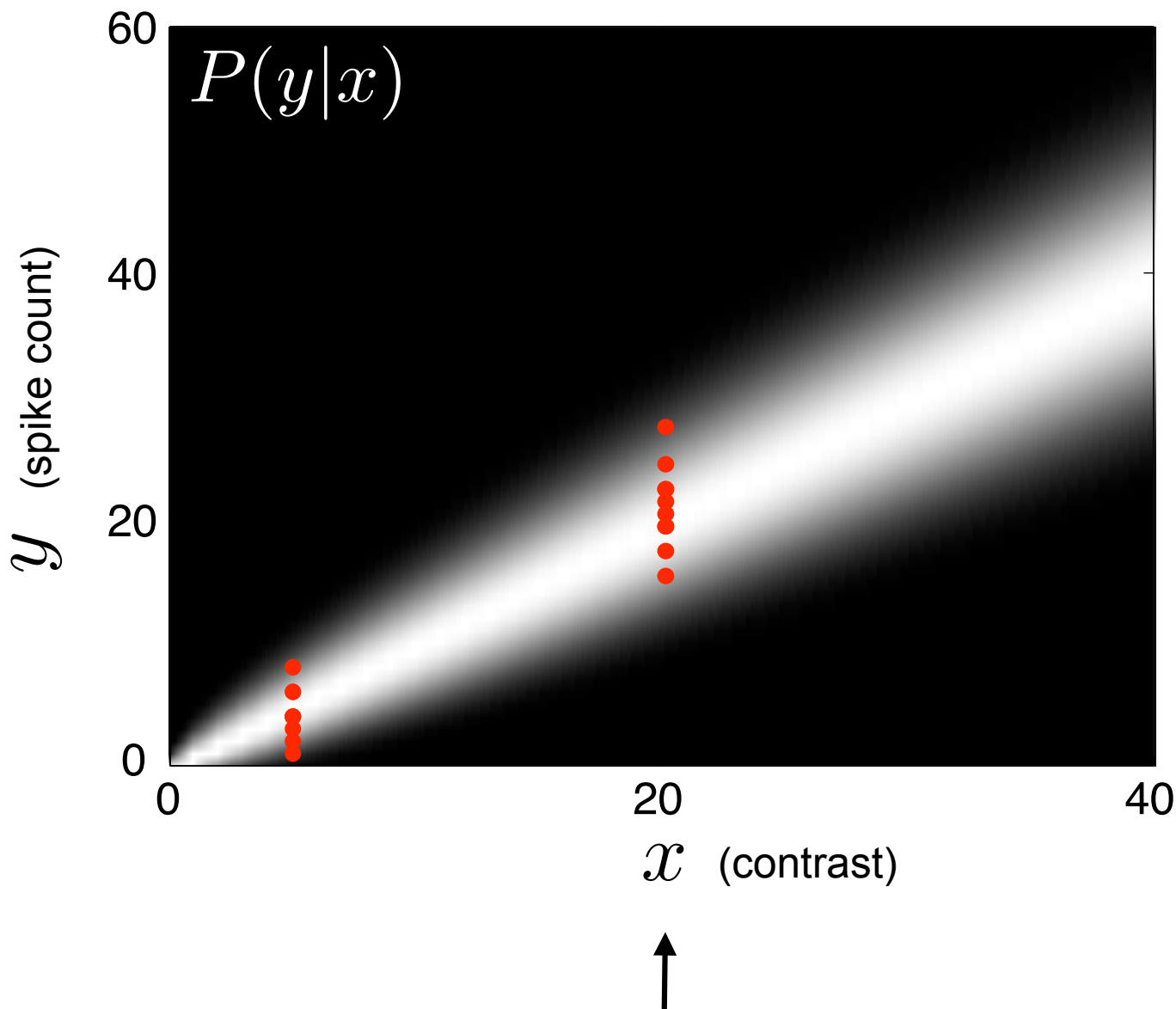
conditional distribution

$$p(y|x = 5)$$



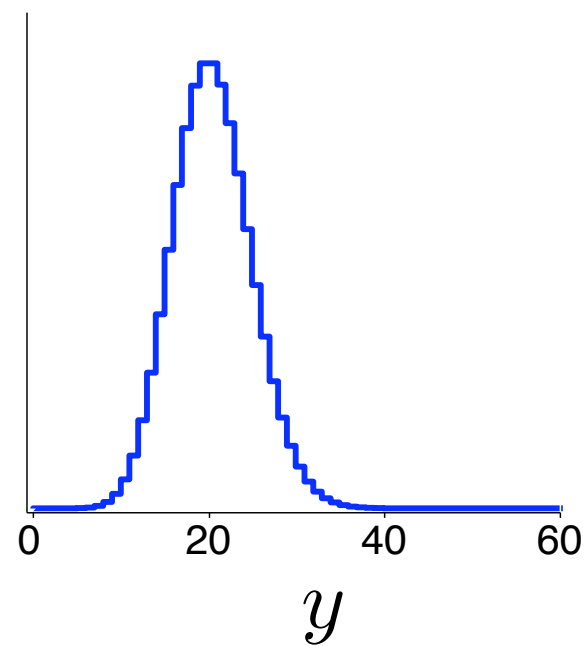
$$\text{mean}(y) = \theta x$$

$$\text{var}(y) = \theta x$$



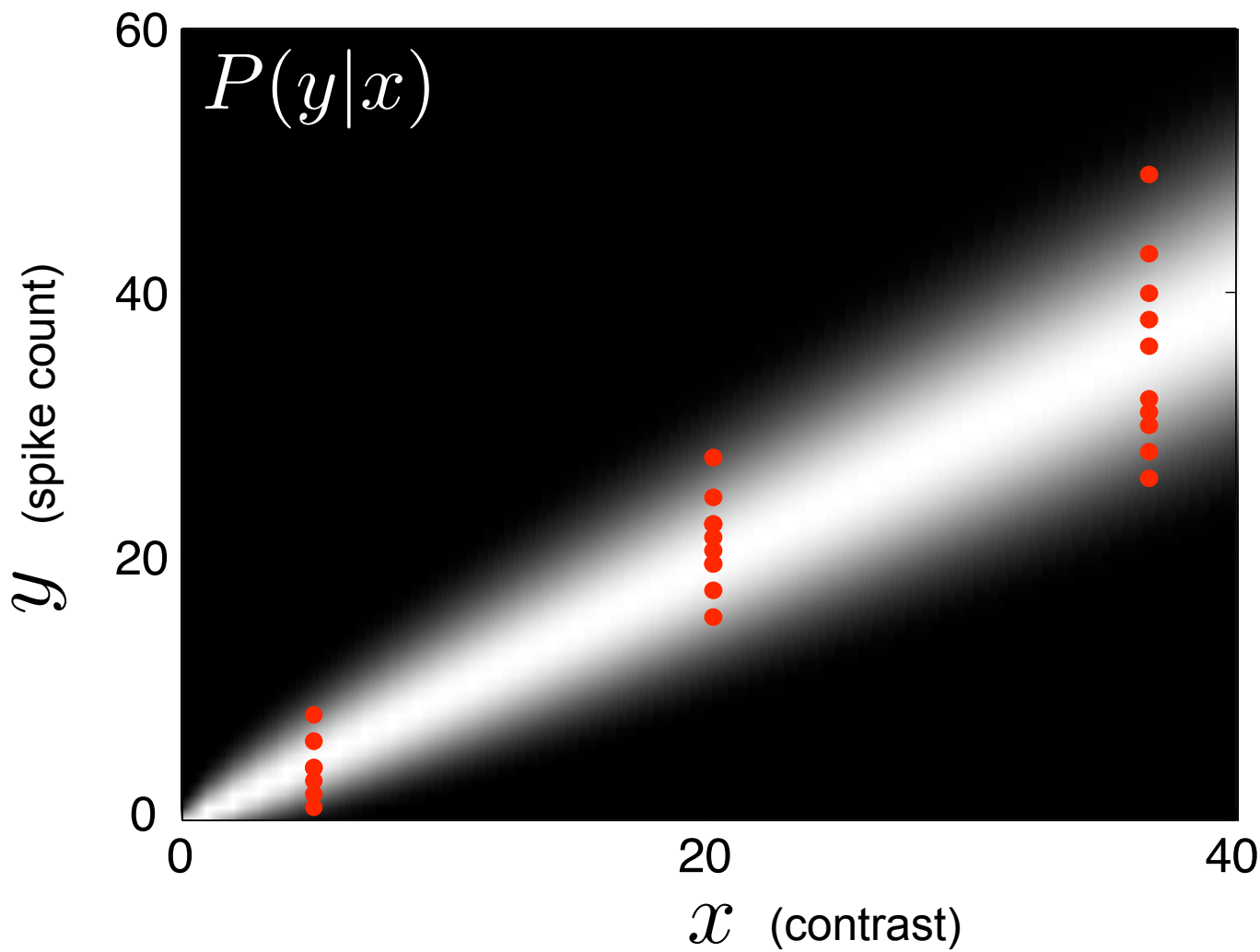
conditional distribution

$$p(y|x = 20)$$



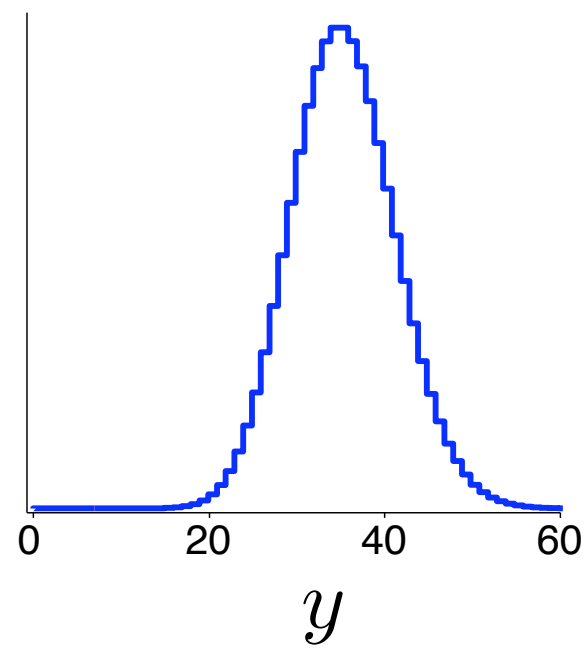
$$\text{mean}(y) = \theta x$$

$$\text{var}(y) = \theta x$$




conditional distribution

$$p(y|x = 35)$$



## Maximum Likelihood Estimation:

- given observed data  $(Y, X)$ , find  $\theta$  that maximizes  $P(Y|X, \theta)$

  
all spike counts   all stimuli   parameters

$$P(Y|X, \theta) = \prod_{i=1}^N \underbrace{P(y_i|x_i, \theta)}_{\text{single-trial probability}}$$


Q: what assumption are we making about the responses?

A: conditional independence across trials!



## Maximum Likelihood Estimation:

- given observed data  $(Y, X)$ , find  $\theta$  that maximizes  $P(Y|X, \theta)$

  
all spike counts    all stimuli    parameters

$$P(Y|X, \theta) = \prod_{i=1}^N \underbrace{P(y_i|x_i, \theta)}_{\text{single-trial probability}}$$

Q: what assumption are we making about the responses?

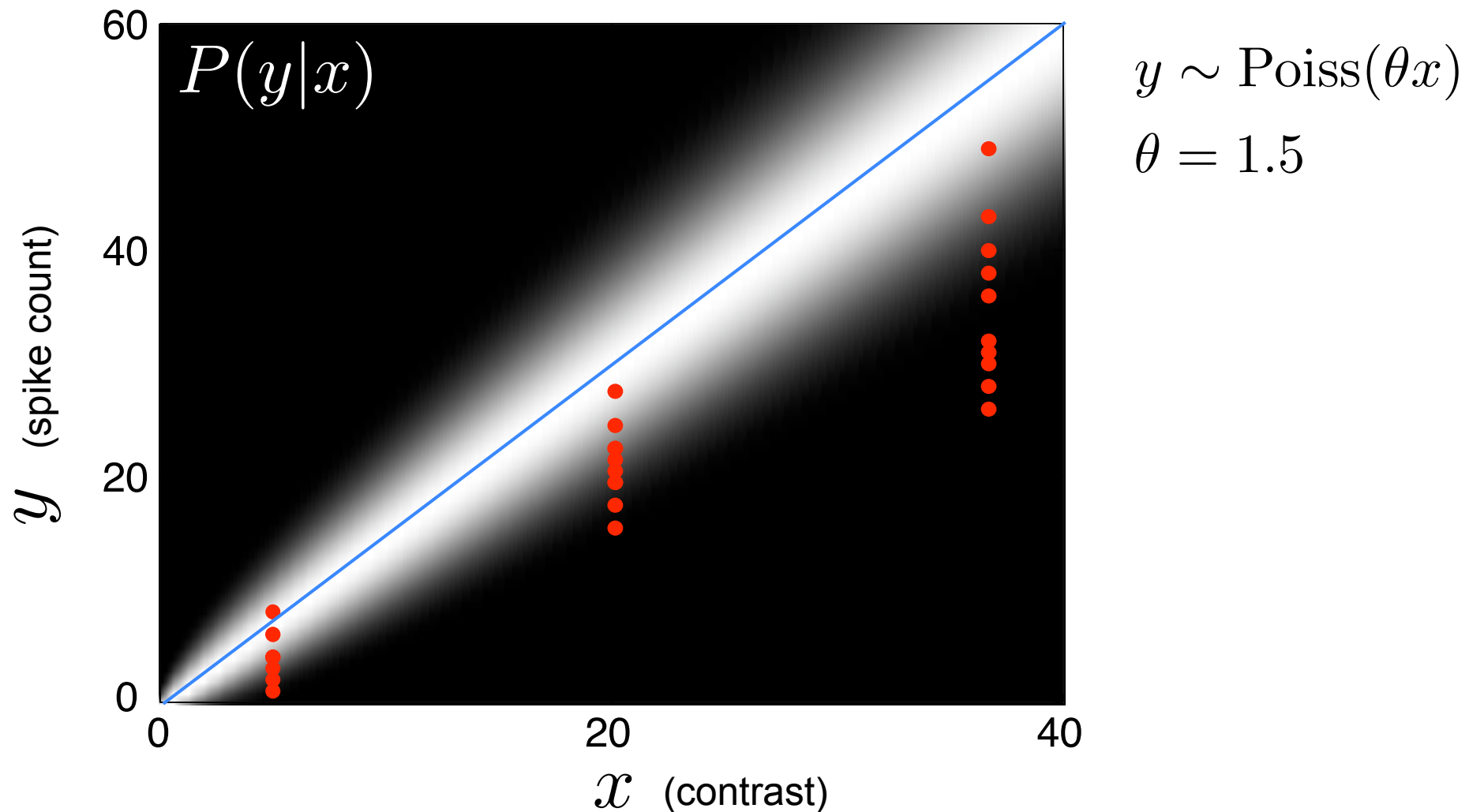
A: conditional independence across trials!

Q: when do we call  $P(Y|X, \theta)$  a *likelihood*?

A: when considering it as a function of  $\theta$  !

## Maximum Likelihood Estimation:

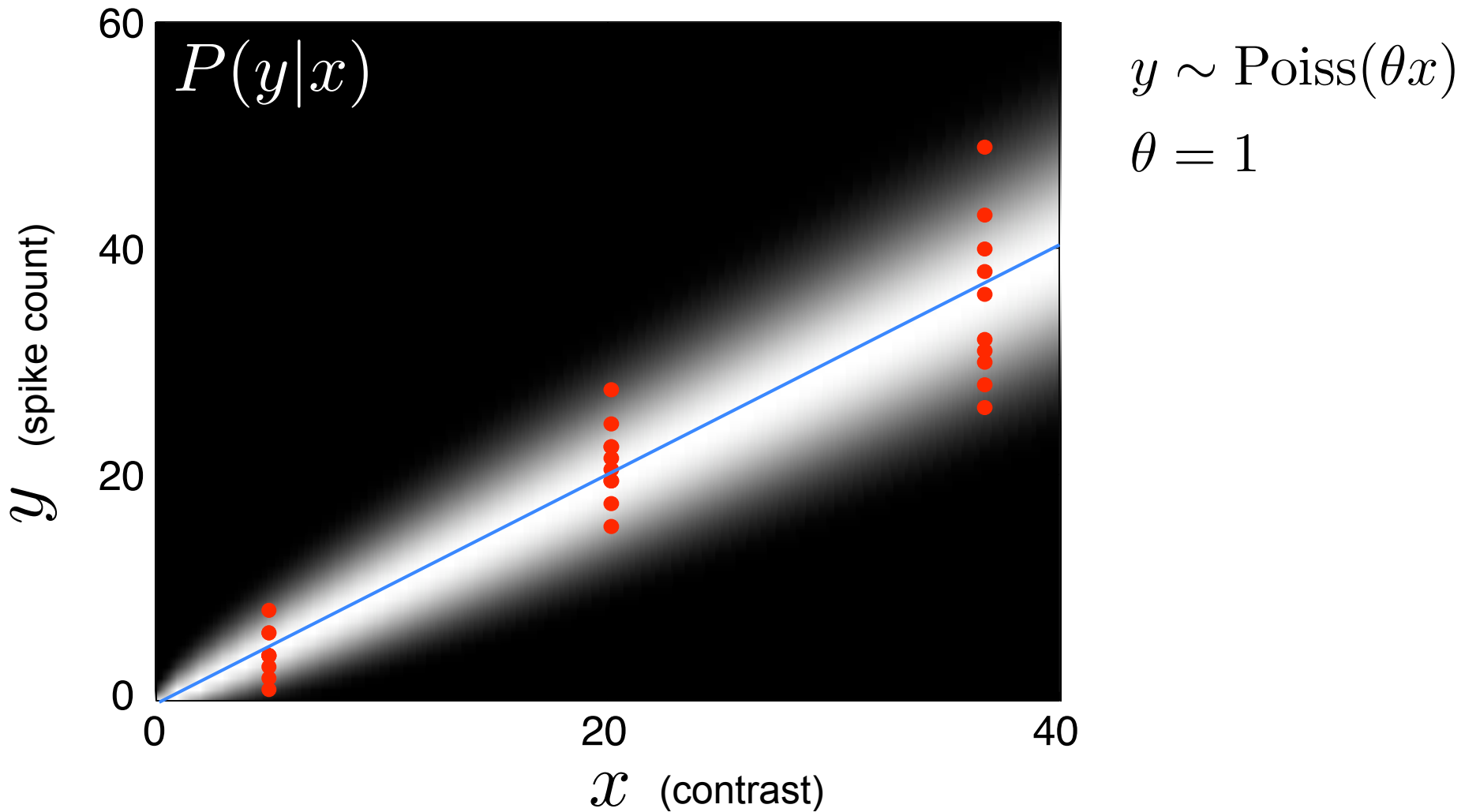
- given observed data  $(Y, X)$ , find  $\theta$  that maximizes  $P(Y|X, \theta)$



- could in theory do this by turning a knob

## Maximum Likelihood Estimation:

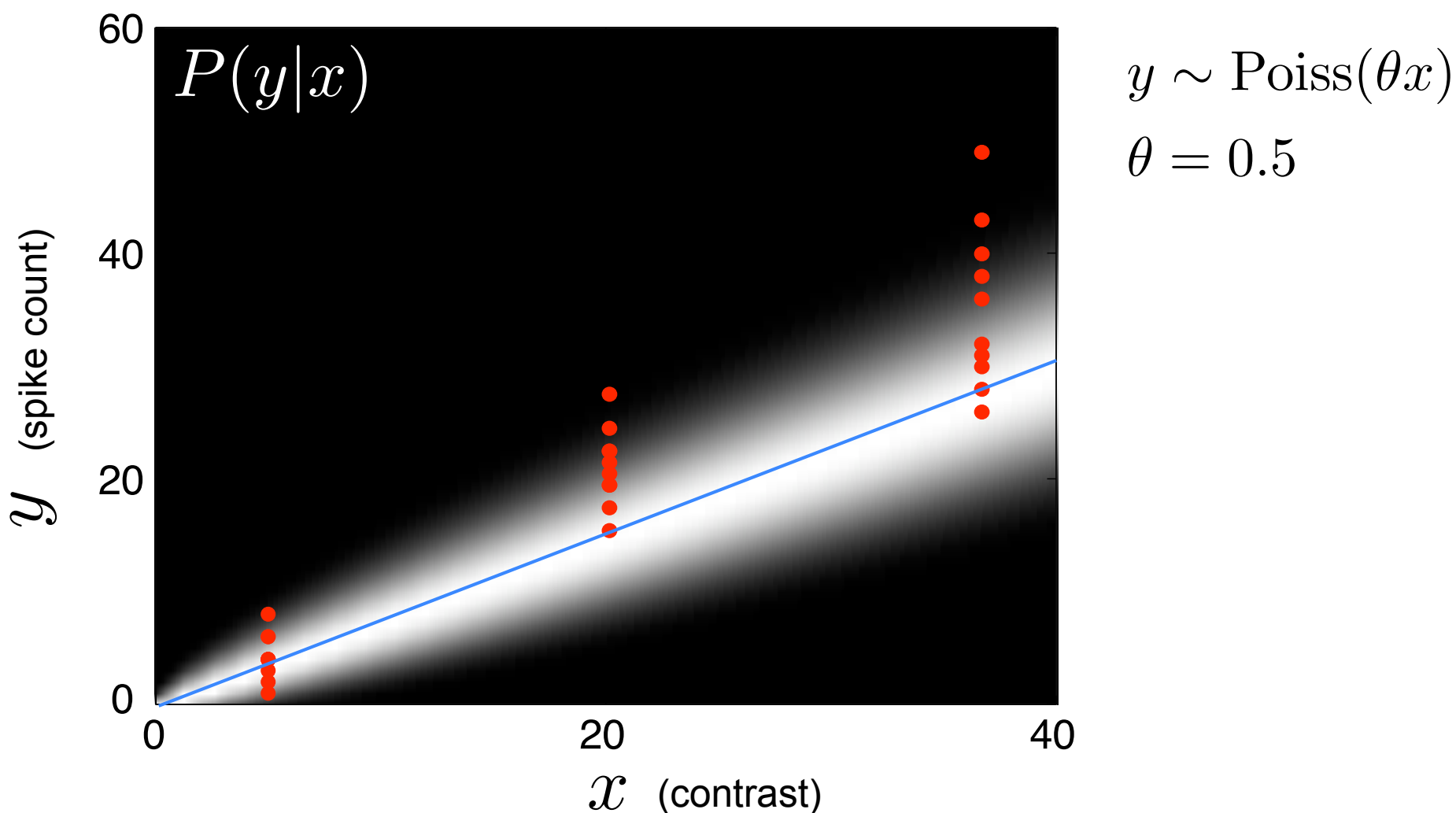
- given observed data  $(Y, X)$ , find  $\theta$  that maximizes  $P(Y|X, \theta)$



- could in theory do this by turning a knob

## Maximum Likelihood Estimation:

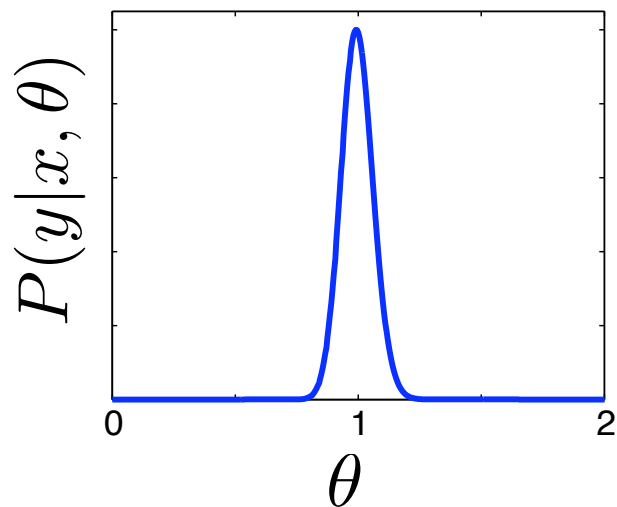
- given observed data  $(Y, X)$ , find  $\theta$  that maximizes  $P(Y|X, \theta)$



- could in theory do this by turning a knob

Likelihood function:  $P(Y|X, \theta)$  as a function of  $\theta$

likelihood

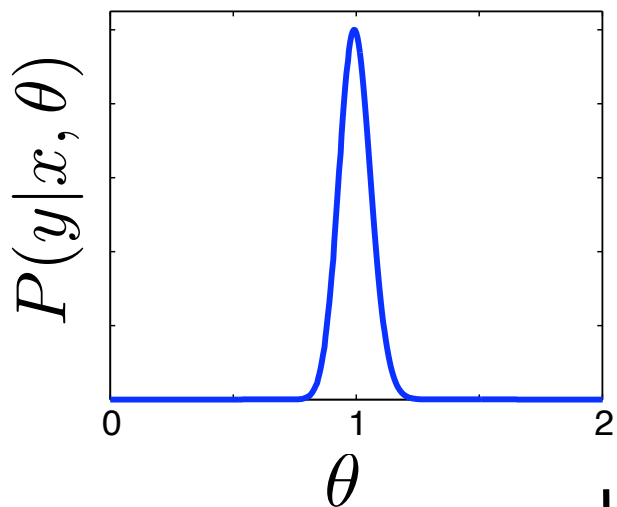


Because data are independent:

$$\begin{aligned} P(Y|X, \theta) &= \prod_i P(y_i|x_i, \theta) \\ &= \prod \frac{1}{y_i!} (\theta x_i)^{y_i} e^{-(\theta x_i)} \end{aligned}$$

Likelihood function:  $P(Y|X, \theta)$  as a function of  $\theta$

likelihood

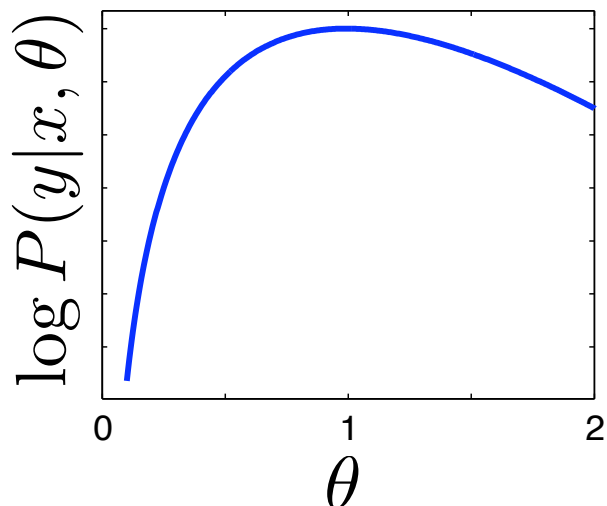


Because data are independent:

$$\begin{aligned} P(Y|X, \theta) &= \prod_i P(y_i|x_i, \theta) \\ &= \prod \frac{1}{y_i!} (\theta x_i)^{y_i} e^{-(\theta x_i)} \end{aligned}$$

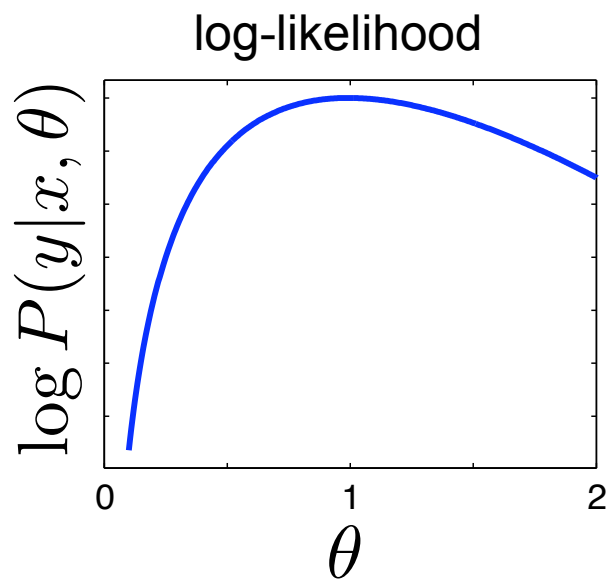
log

log-likelihood



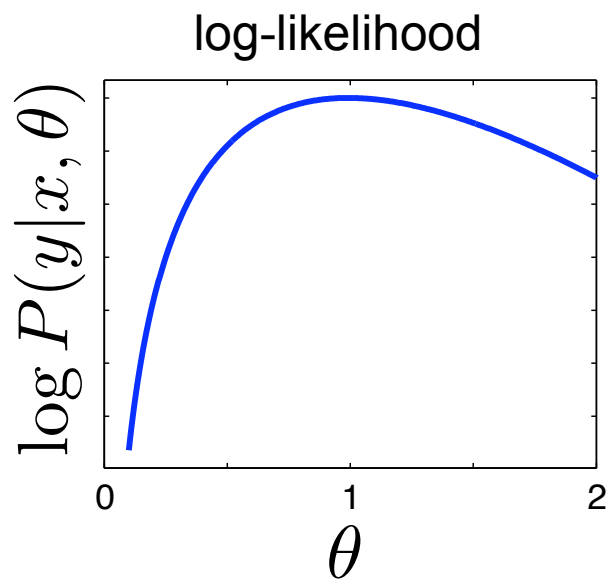
$$\begin{aligned} \log P(Y|X, \theta) &= \sum_i \log P(y_i|x_i, \theta) \\ &= \sum y_i \log \theta - \theta x_i + c \end{aligned}$$





$$\begin{aligned}\log P(Y|X, \theta) &= \sum_i \log P(y_i|x_i, \theta) \\ &= \sum y_i \log \theta - \theta x_i + c \\ &= \log \theta (\sum y_i) - \theta (\sum x_i)\end{aligned}$$

Do it: solve for  $\theta$



$$\begin{aligned}\log P(Y|X, \theta) &= \sum_i \log P(y_i|x_i, \theta) \\ &= \sum y_i \log \theta - \theta x_i + c \\ &= \log \theta (\sum y_i) - \theta (\sum x_i)\end{aligned}$$

- Closed-form solution when model in “exponential family”

$$\begin{aligned}\frac{d}{d\theta} \log P(Y|X, \theta) &= \frac{1}{\theta} \sum y_i - \sum x_i = 0 \\ \implies \hat{\theta}_{ML} &= \frac{\sum y_i}{\sum x_i}\end{aligned}$$

# Properties of the MLE (maximum likelihood estimator)

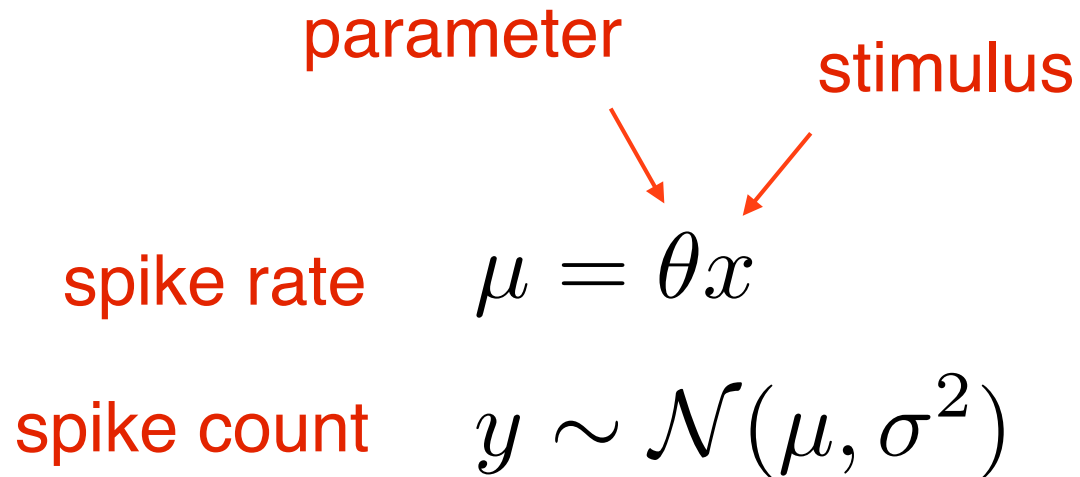
- **consistent**  
(converges to true  $\theta$  in limit of infinite data)
- **efficient**  
(converges as quickly as possible,  
i.e., achieves minimum possible asymptotic error)

# simple example #2: linear Gaussian neuron

parameter stimulus

spike rate  $\mu = \theta x$

spike count  $y \sim \mathcal{N}(\mu, \sigma^2)$

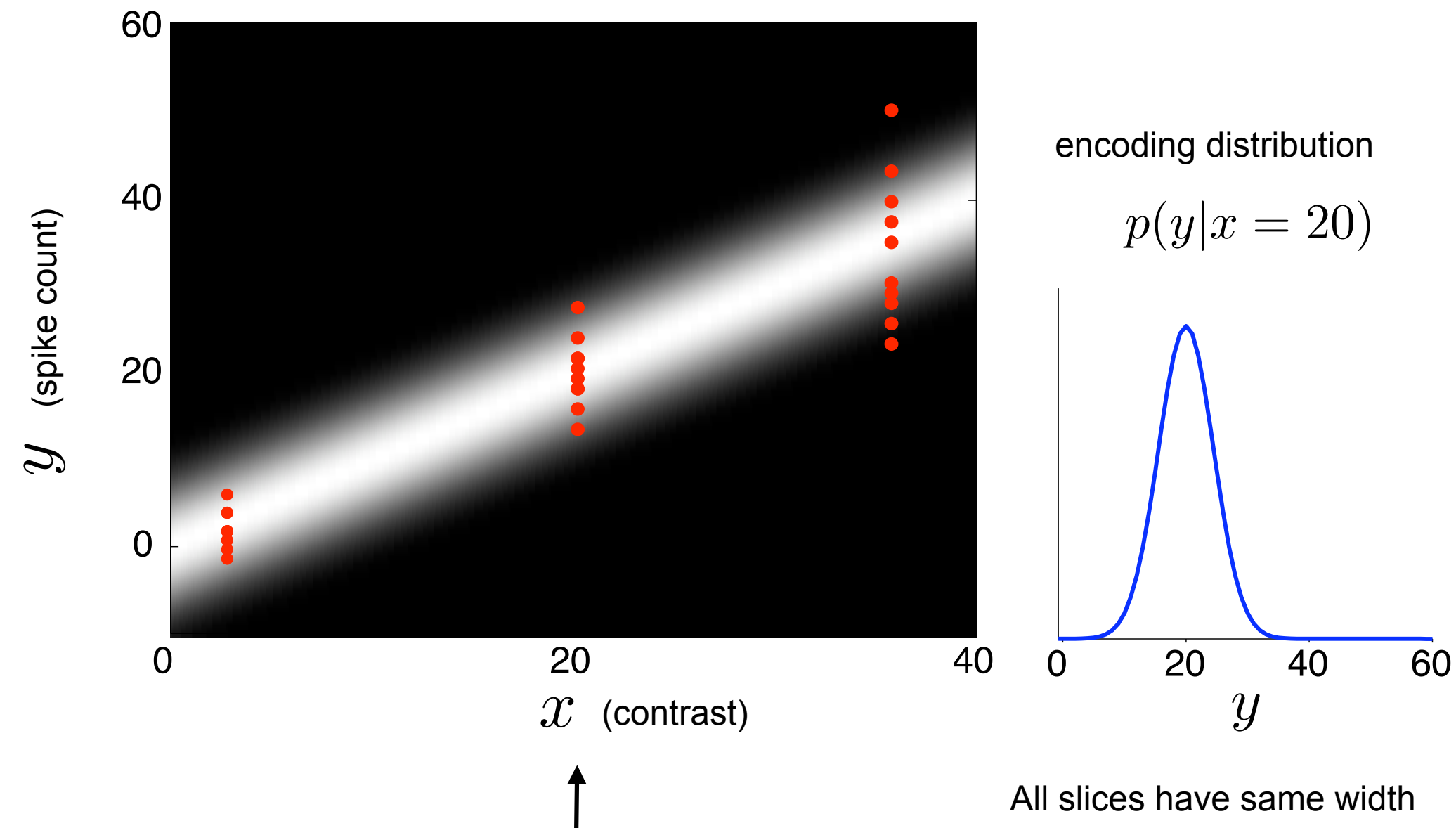


encoding model:

$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(y - \theta x)^2}{2\sigma^2}}$$

$$\text{mean}(y) = \theta x$$

$$\text{var}(y) = \sigma^2$$



$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(y-\theta x)^2}{2\sigma^2}}$$

Log-Likelihood

$$\log P(Y|X, \theta) = - \sum \frac{(y_i - \theta x_i)^2}{2\sigma^2} + c$$

Differentiate, set to zero, and solve for  $\theta$



$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(y-\theta x)^2}{2\sigma^2}}$$

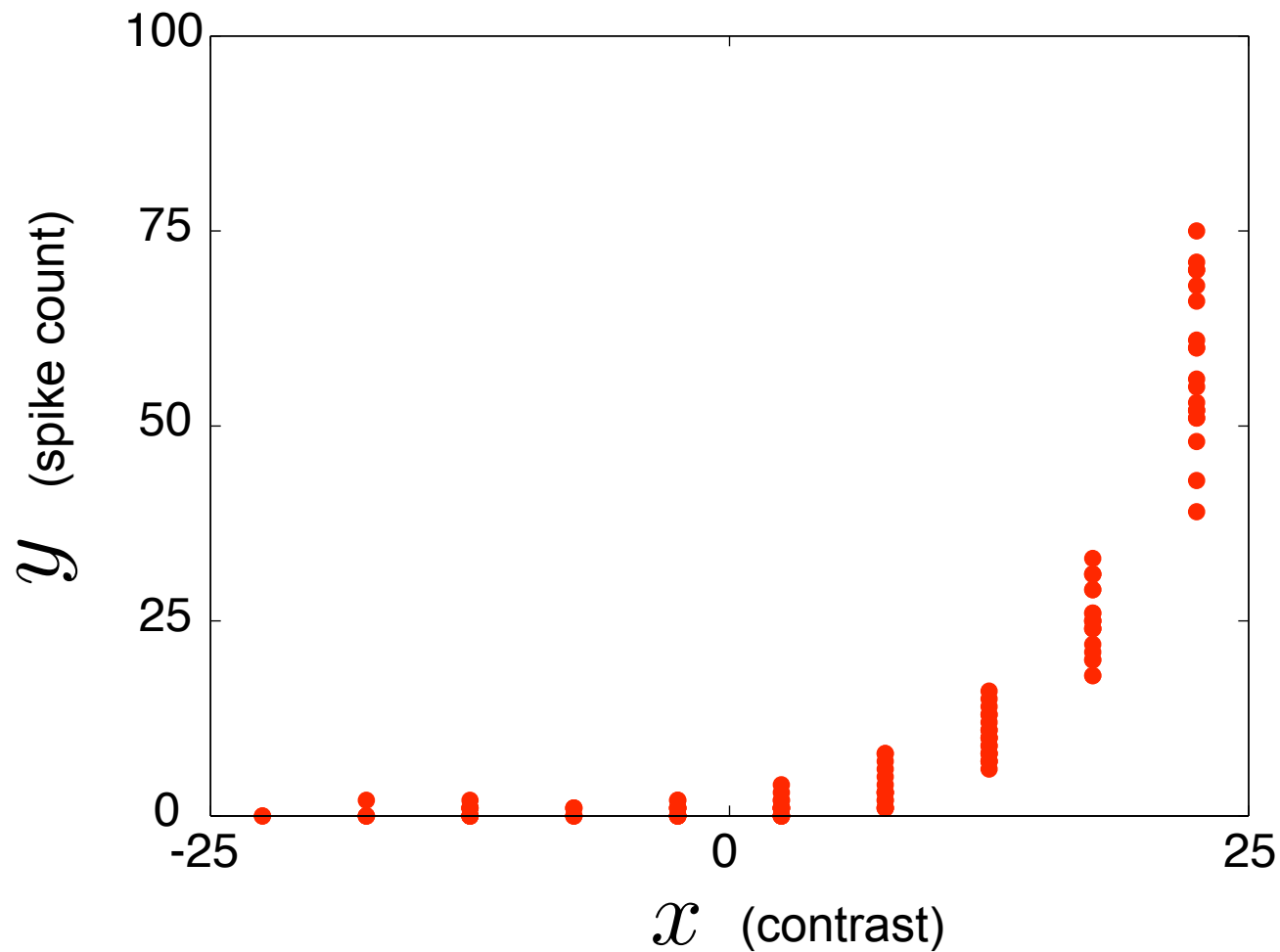
Log-Likelihood  $\log P(Y|X, \theta) = - \sum \frac{(y_i - \theta x_i)^2}{2\sigma^2} + c$

$$\frac{d}{d\theta} \log P(Y|X, \theta) = - \sum \frac{(y_i - \theta x_i)x_i}{\sigma^2} = 0$$

Maximum-Likelihood Estimator:  $\hat{\theta}_{ML} = \frac{\sum y_i x_i}{\sum x_i^2}$   
 (“Least squares regression” solution)

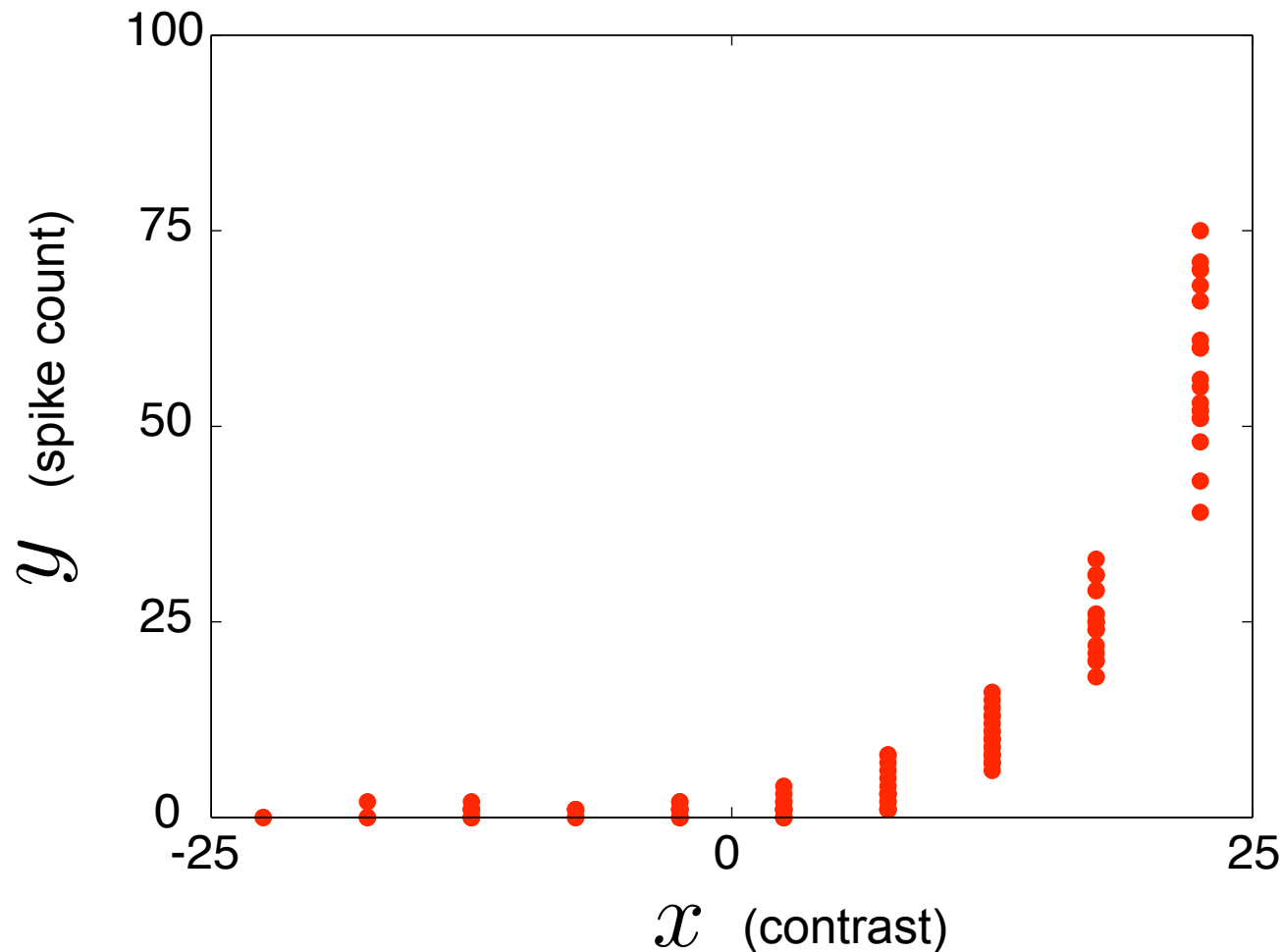
(Recall that for Poisson,  $\hat{\theta}_{ML} = \frac{\sum y_i}{\sum x_i}$  )

# example #3: unknown neuron



Be the computational neuroscientist: what model would you use?

# Example 3: unknown neuron



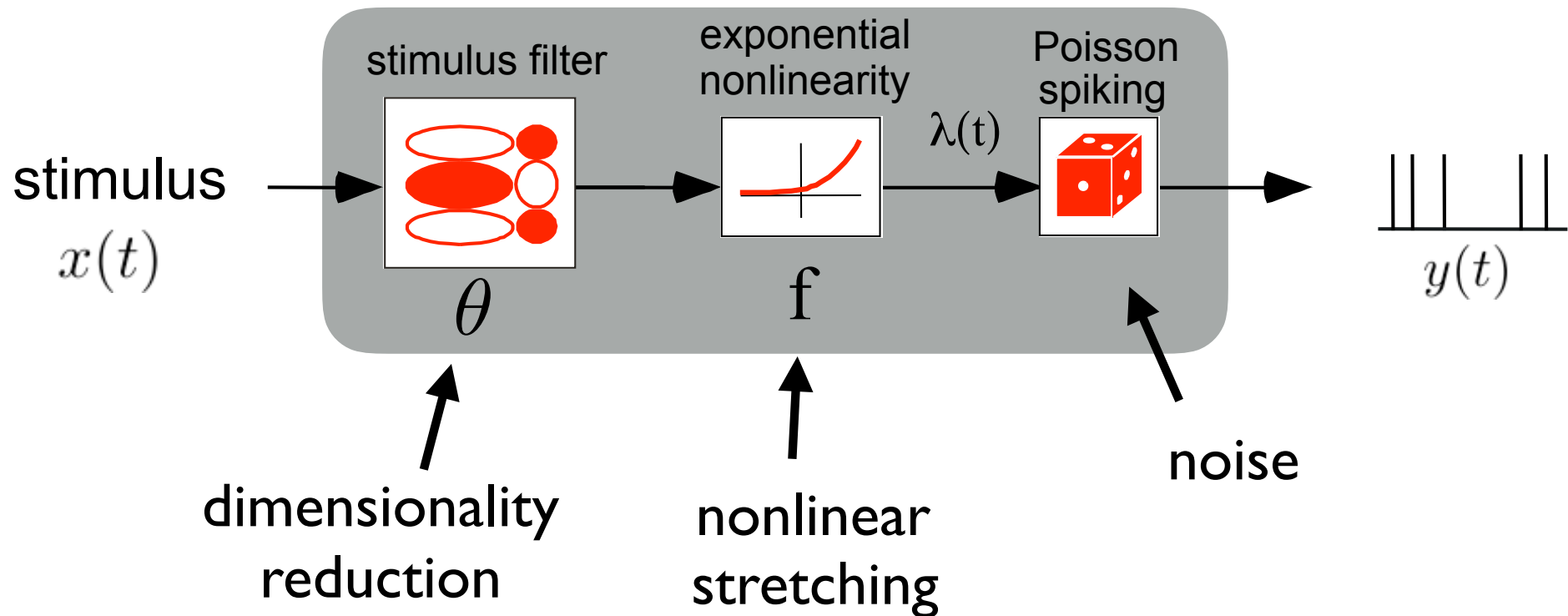
More general setup:  $\lambda = f(\theta x)$  firing rate is nonlinear

**This is a GLM!**

$y \sim \text{Poisson}(\lambda)$  Poisson firing

# “basic” Poisson generalized linear model (GLM)

## Linear-Nonlinear-Poisson (LNP) model



$$\begin{array}{ll} \text{spike rate} & \lambda = f(\vec{k} \cdot \vec{x}) \\ \text{spike count} & y \sim \text{Pois}(\lambda) \end{array}$$

- also known as a “cascade” model

# What is a GLM?

Be careful about terminology:

GLM

≠

GLM

General Linear Model

Generalized Linear Model

(Nelder 1972)

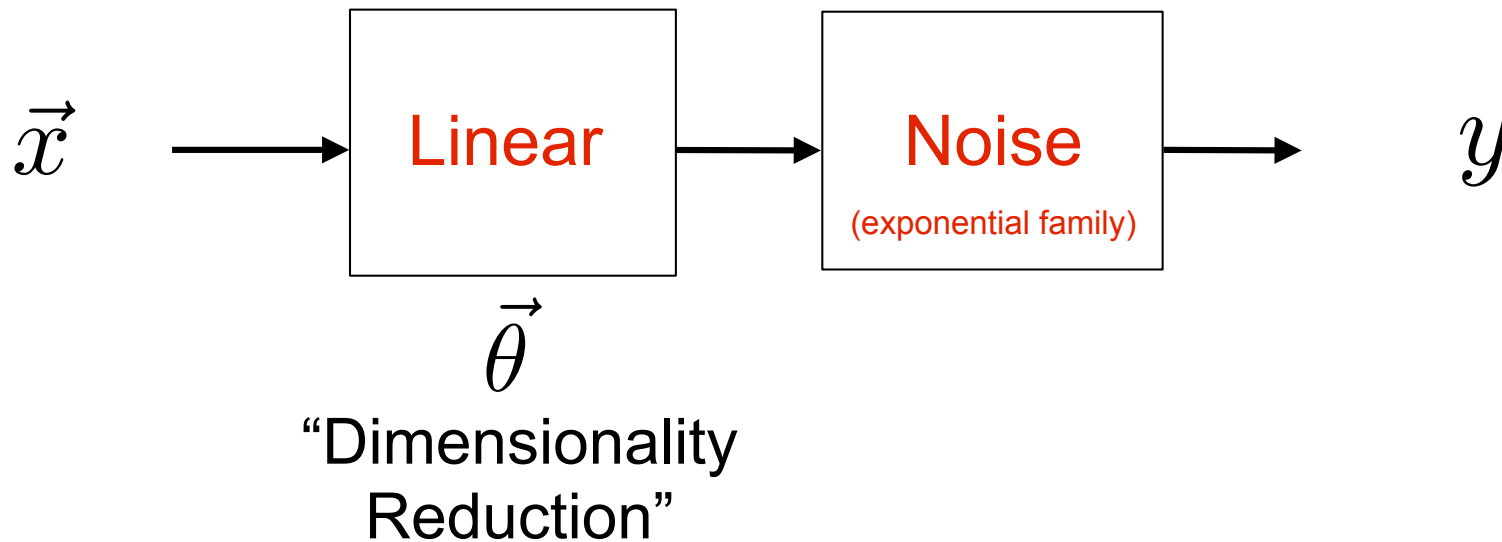
Linear

~~Linear~~

## **Moral:**

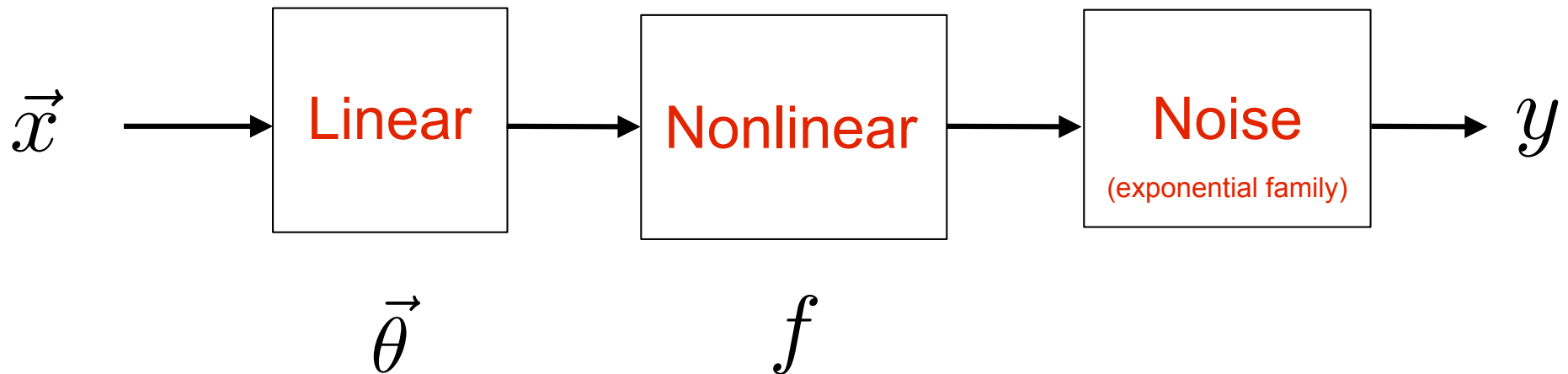
Be careful when naming your model!

# 1. General Linear Model



- Examples:
1. Gaussian  $y = \vec{\theta} \cdot \vec{x} + \epsilon$
  2. Poisson  $y \sim \text{Pois}(\vec{\theta} \cdot \vec{x})$

## 2. Generalized Linear Model

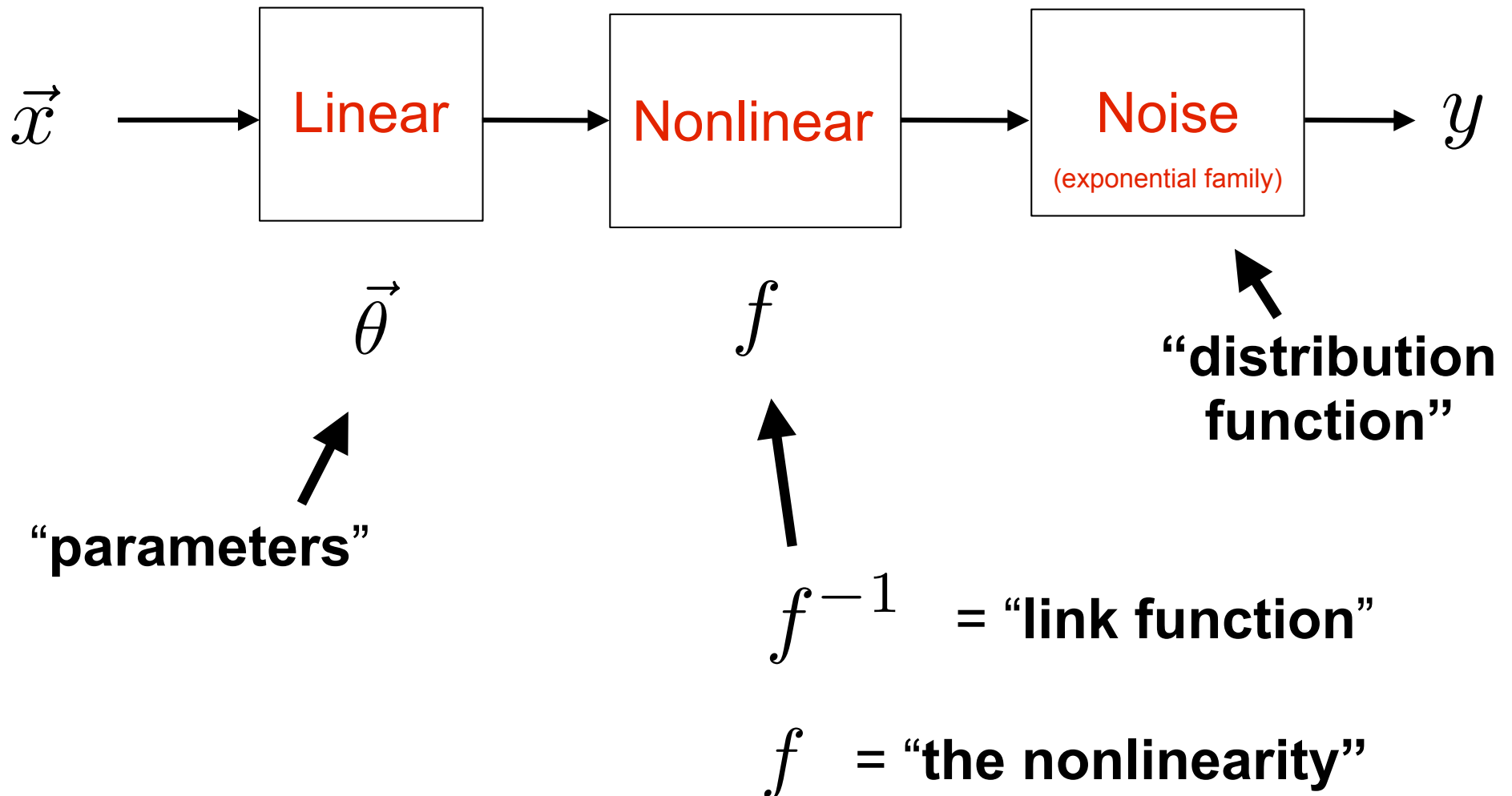


- Examples:
1. Gaussian  $y = f(\vec{\theta} \cdot \vec{x}) + \epsilon$
  2. Poisson  $y \sim \text{Pois}(f(\vec{\theta} \cdot \vec{x}))$



## 2. Generalized Linear Model

Terminology:



# Applying it to data

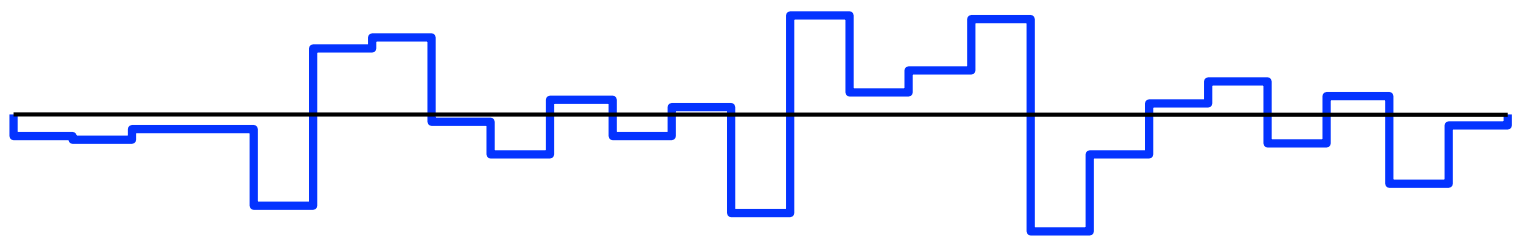
response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

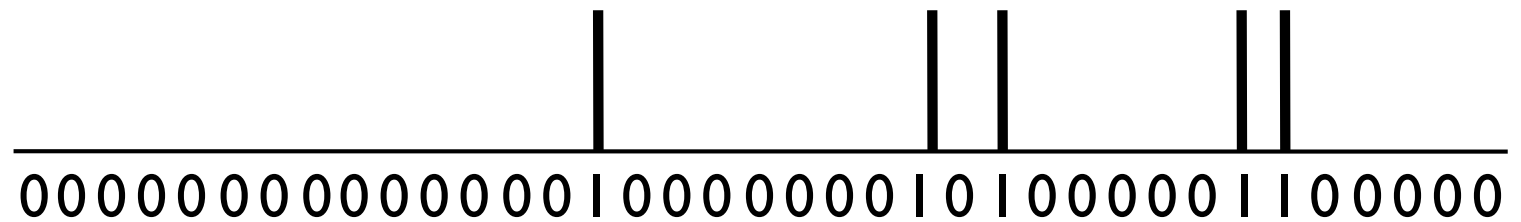
linear  
filter

vector stimulus  
at time  $t$

stimulus



response



time  $\longrightarrow$

response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

walk through the data  
one time bin at a time

$t = 1$

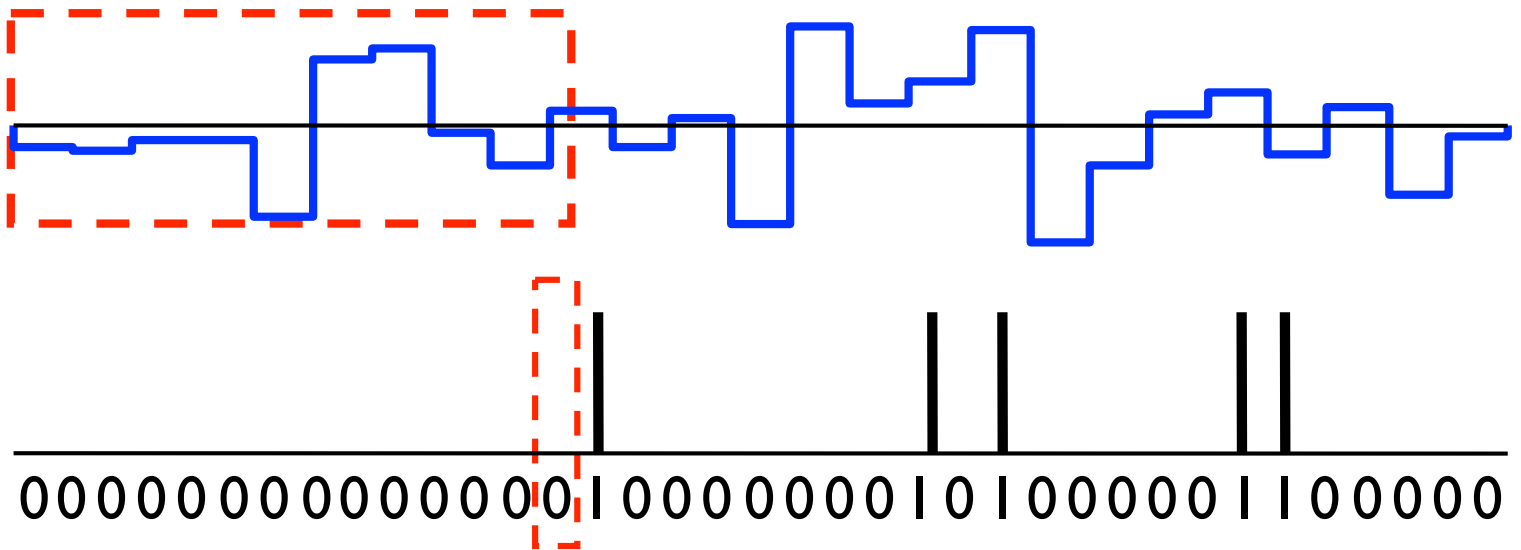
stimulus

$\vec{x}_t$

response

time  $\longrightarrow$

$y_t$



response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

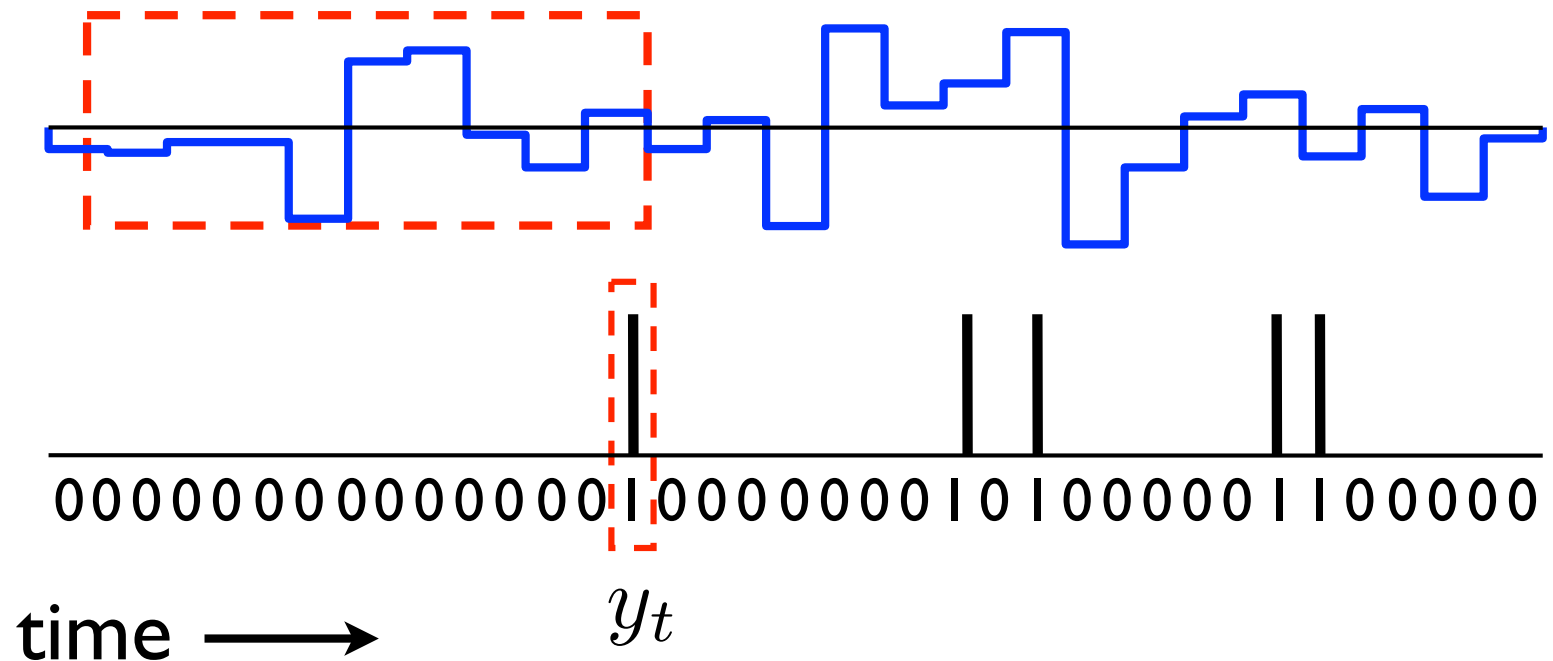
walk through the data  
one time bin at a time

$t = 2$

stimulus

$\vec{x}_t$

response



response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

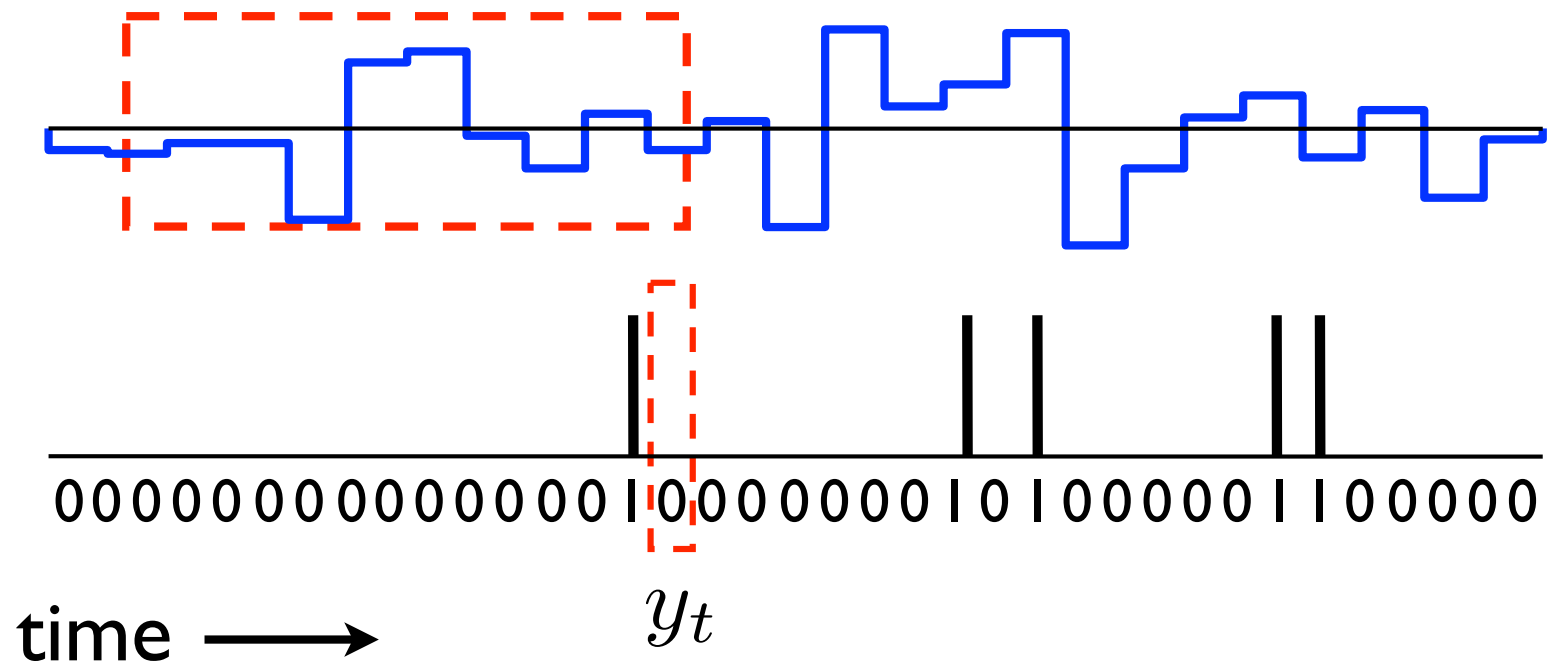
walk through the data  
one time bin at a time

$t = 3$

stimulus

$\vec{x}_t$

response



response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

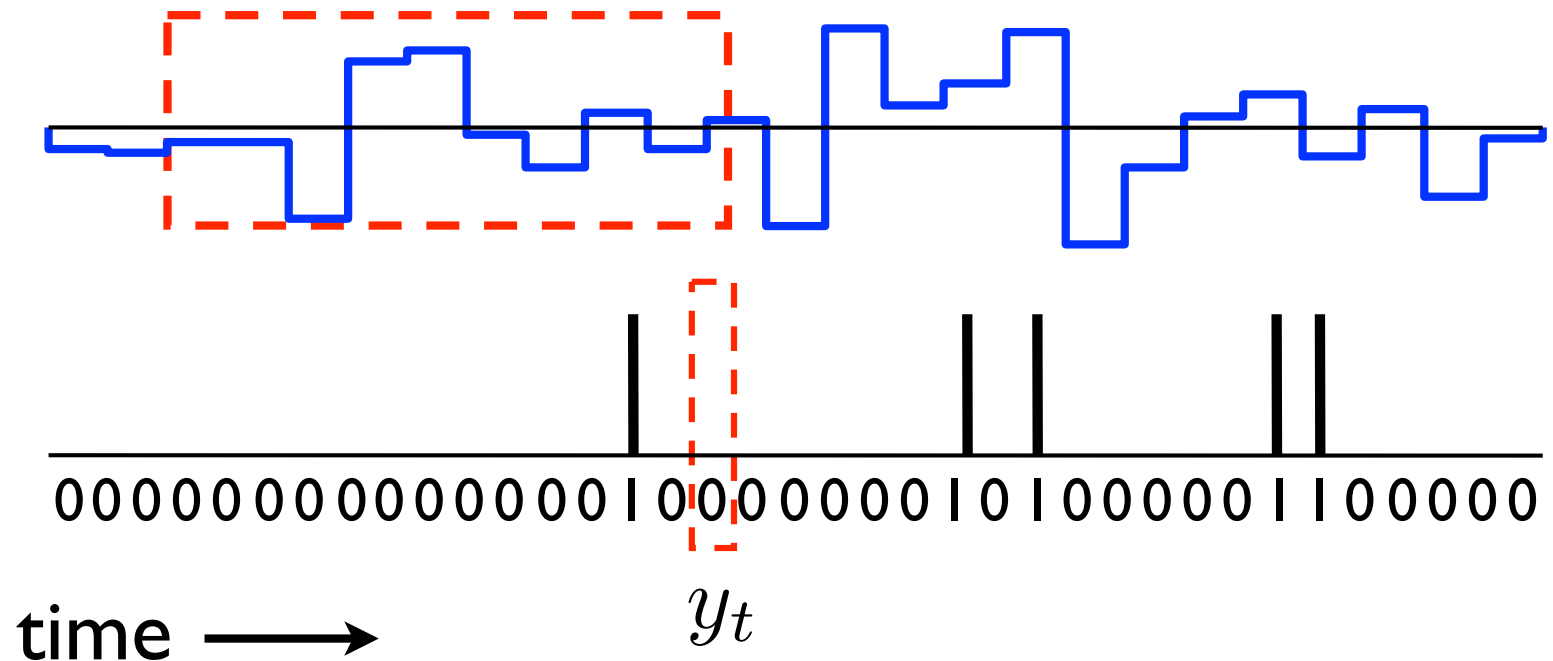
walk through the data  
one time bin at a time

$t = 4$

stimulus

$\vec{x}_t$

response



response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

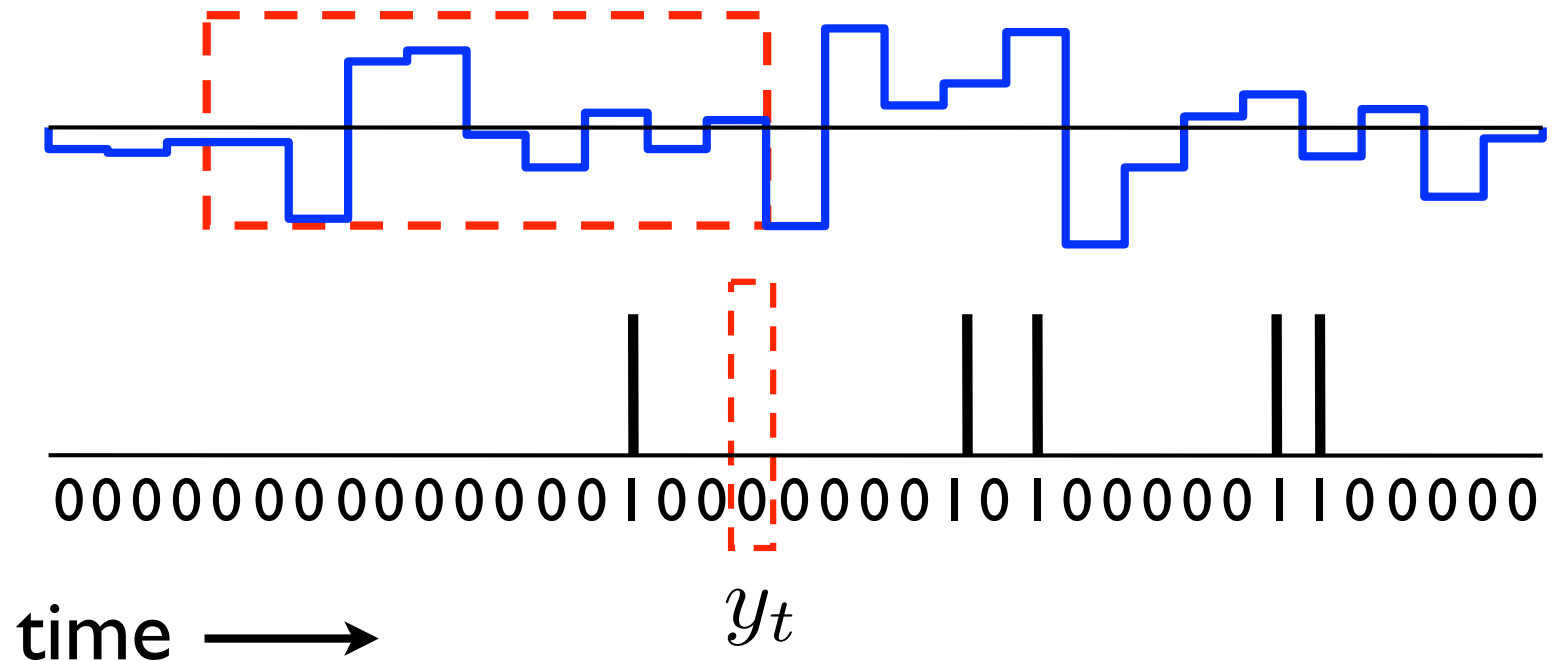
walk through the data  
one time bin at a time

$t = 5$

stimulus

$\vec{x}_t$

response



response  
at time  $t$

$$y_t = \vec{k} \cdot \vec{x}_t + \text{noise}$$

linear  
filter

vector stimulus  
at time  $t$

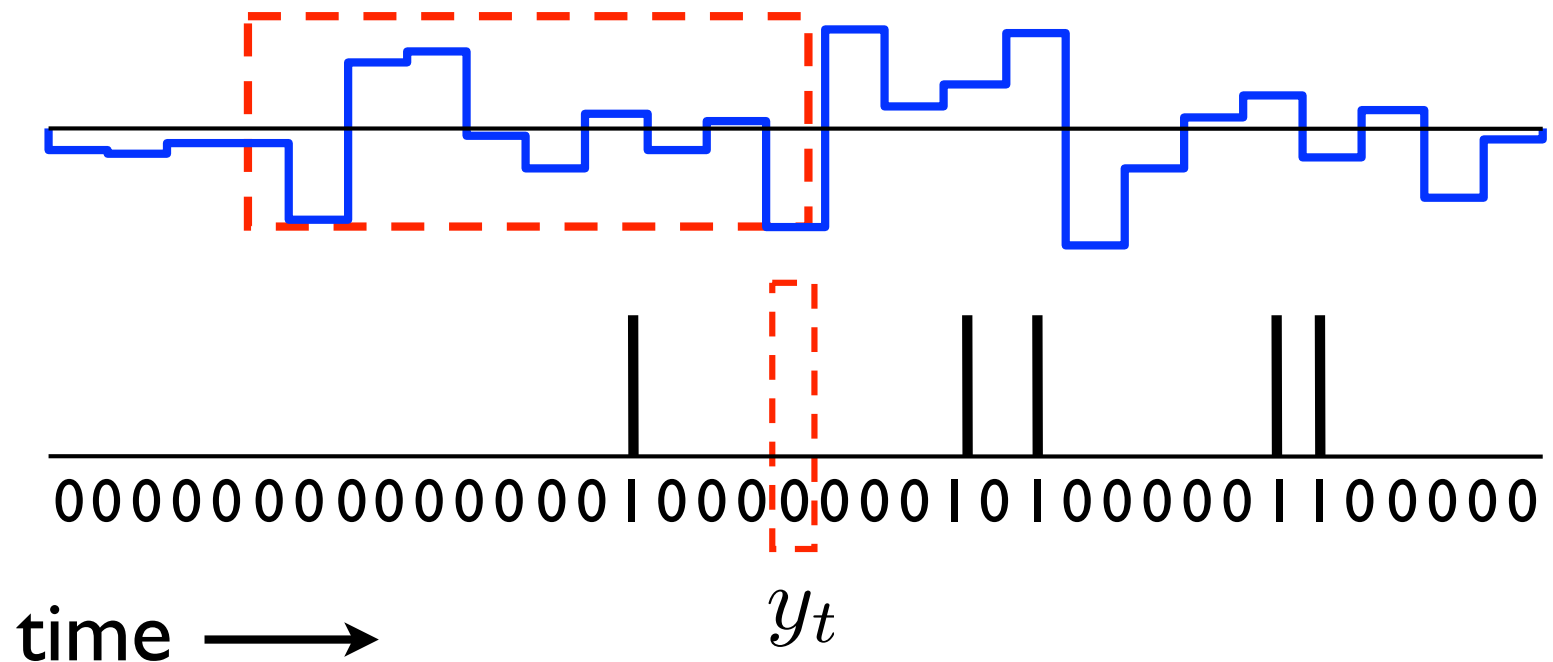
walk through the data  
one time bin at a time

$t = 6$

stimulus

$\vec{x}_t$

response






Build up to following matrix version:

$$\begin{array}{c} \text{time} \\ \downarrow \end{array} \quad Y = X \vec{k} + \text{noise}$$

$\nearrow$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix} = \begin{bmatrix} \text{blue step function} \\ \text{blue step function} \\ \text{blue step function} \\ \vdots \end{bmatrix} \begin{bmatrix} \vec{k} \end{bmatrix}$$

$\nwarrow$

 **design matrix**

# Computing maximum likelihood estimate

$$Y = X \vec{k} + \text{noise}$$

time ↓

$\begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix} = \begin{bmatrix} \text{step function} \\ \text{step function} \\ \text{step function} \\ \vdots \end{bmatrix} \begin{bmatrix} \vec{k} \end{bmatrix}$

I. “Linear-Gaussian” GLM:

$$\hat{k} = \underbrace{(X^T X)^{-1}}_{\text{stimulus covariance}} \underbrace{X^T Y}_{\text{spike-triggered avg (STA)}}$$

# Computing maximum likelihood estimate

$$Y = f(X\vec{k}) + \text{noise}$$

time ↓

The diagram illustrates the components of the maximum likelihood estimation equation. On the left, a vertical vector  $Y$  is shown with elements 0, 0, 1, and a vertical ellipsis. To its left, a vertical arrow labeled "time" points downwards. In the center is an equals sign. To the right of the equals sign is a large square matrix  $X$  containing three blue step functions and a vertical ellipsis. To the right of the matrix is a vertical vector  $\vec{k}$ . An arrow points from the vector  $\vec{k}$  to the function  $f$  in the equation above.

2. Poisson GLM: `k = glmfit(X,Y,'Poisson');`

maximum likelihood fit  
(assumes exponential nonlinearity by default)

# Computing maximum likelihood estimate

$$Y = f(X\vec{k}) + \textit{noise}$$

time ↓

The diagram illustrates the relationship between the observed data  $Y$ , the feature matrix  $X$ , and the parameter vector  $\vec{k}$ . The vector  $Y$  is shown as a column of values  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$ . The matrix  $X$  is shown as a column of blue step functions. The vector  $\vec{k}$  is shown as a column of values  $\begin{bmatrix} \vec{k} \end{bmatrix}$ . The equation  $Y = f(X\vec{k}) + \textit{noise}$  is shown above the matrix  $X$ . An arrow points from the word "time" to the vector  $Y$ , indicating that the values in  $Y$  are observed over time.

3. Bernoulli GLM: `k = glmfit(X,Y,'binomial');`  
outputs 0 and 1 (assumes **logistic** nonlinearity by default)

**“logistic regression”**

# GLM summary

1. Linear-Gaussian GLM:  $Y|X, \vec{k} \sim \mathcal{N}(X\vec{k}, \sigma^2 I)$  **continuous**

log-likelihood:  $-\frac{1}{2\sigma^2} (Y - X\vec{k})^\top (Y - X\vec{k}) + \text{const}$

MLE:  $\hat{k} = (X^\top X)^{-1} X^\top Y$

2. Poisson GLM:  $y|\vec{x}, \vec{k} \sim \text{Pois}(f(\vec{x}_t \cdot \vec{k}))$  **integer counts**

log-likelihood:  $\mathcal{L} = Y^\top \log f(X\vec{k}) - \mathbf{1}^\top f(X\vec{k})$

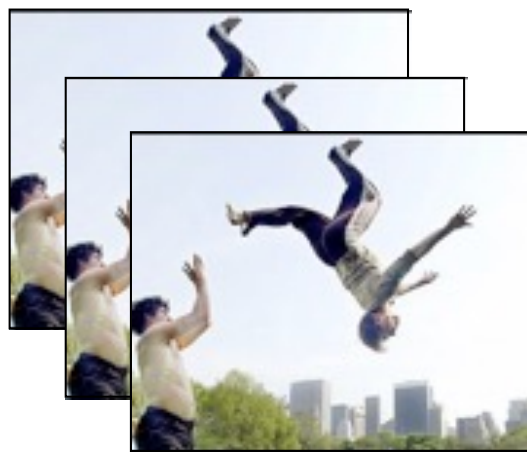
3. Bernoulli GLM:  $y_t|\vec{x}_t, \vec{k} \sim \text{Ber}(f(\vec{x}_t \cdot \vec{k}))$  **binary counts**

log-likelihood:  $\mathcal{L} = Y^\top \log f(X\vec{k}) - (1 - Y)^\top \log(1 - f(X\vec{k}))$

“logistic regression” if  $f(x) = \frac{1}{1 + e^{-x}}$

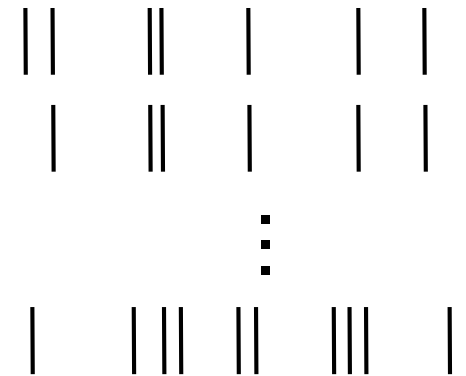
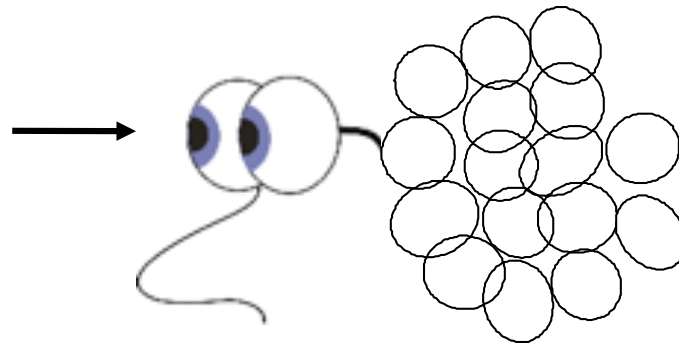
**NEXT:**

## GLMs with spike-history and coupling



$X$

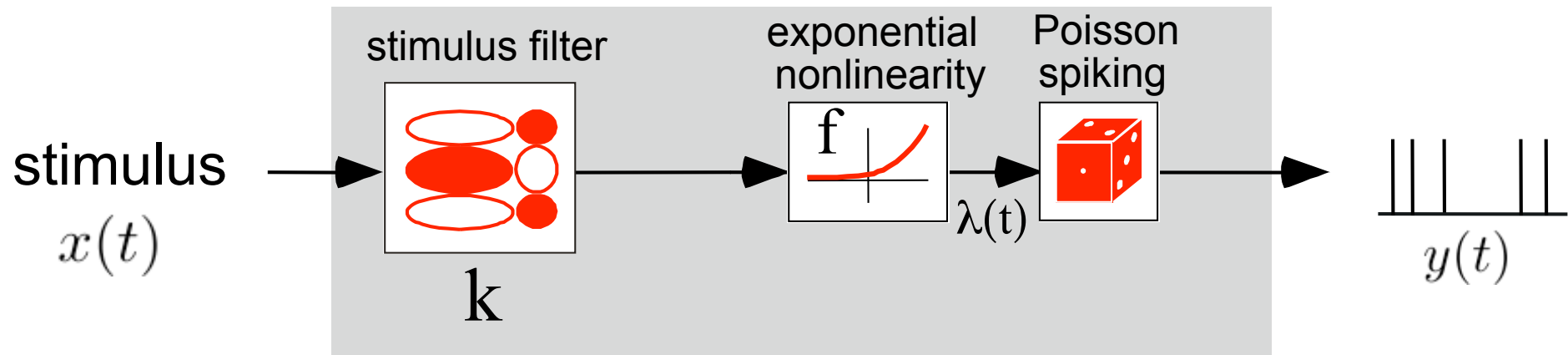
stimuli



$y$

spike trains

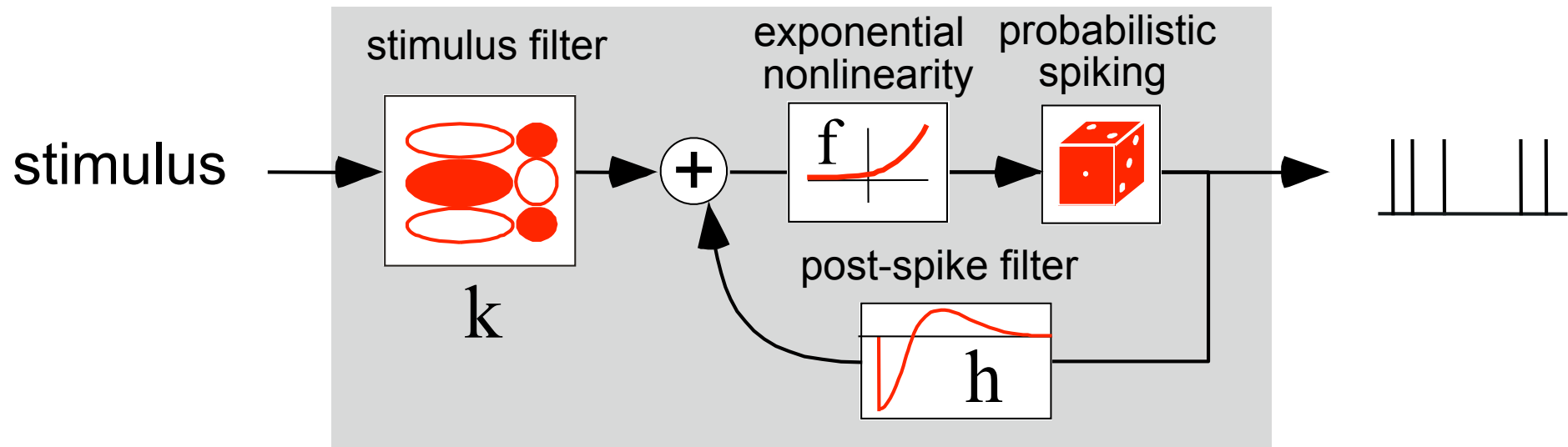
# Poisson GLM



spike rate  $\lambda(t) = f(k \cdot x(t))$

- problem: assumes spiking depends only on stimulus!

# Poisson GLM with spike-history dependence



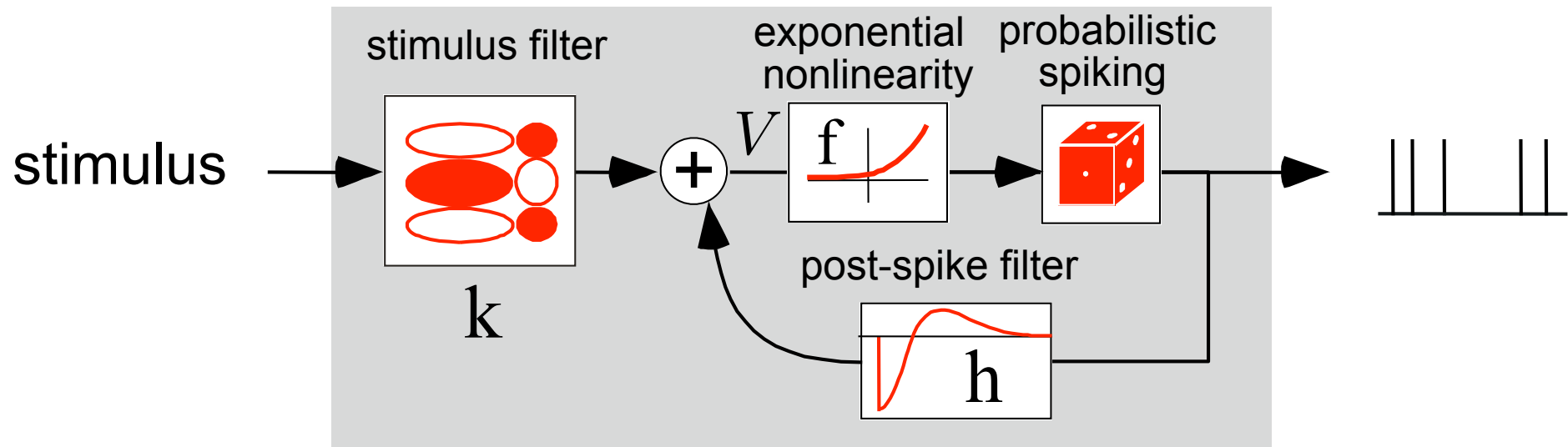
spike rate: 
$$\lambda(t) = f( \vec{k} \cdot \vec{x}(t) + \vec{h} \cdot \vec{y}_{hst}(t) )$$

$$= e^{\vec{k} \cdot \vec{x}(t)} \cdot e^{\vec{h} \cdot \vec{y}_{hst}(t)}$$

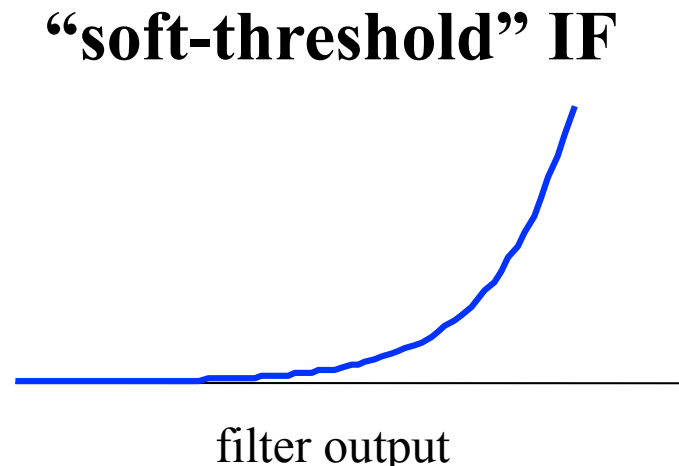
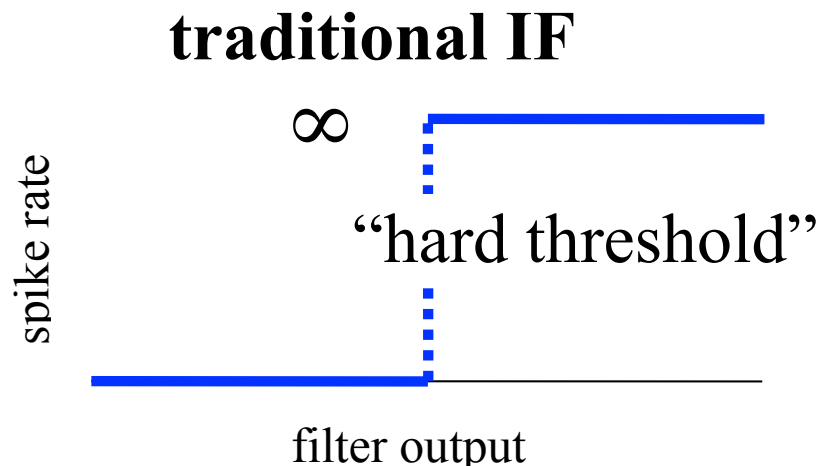
- output: no longer a Poisson process
- interpretation: “soft-threshold” integrate-and-fire model



# Poisson GLM with spike-history dependence

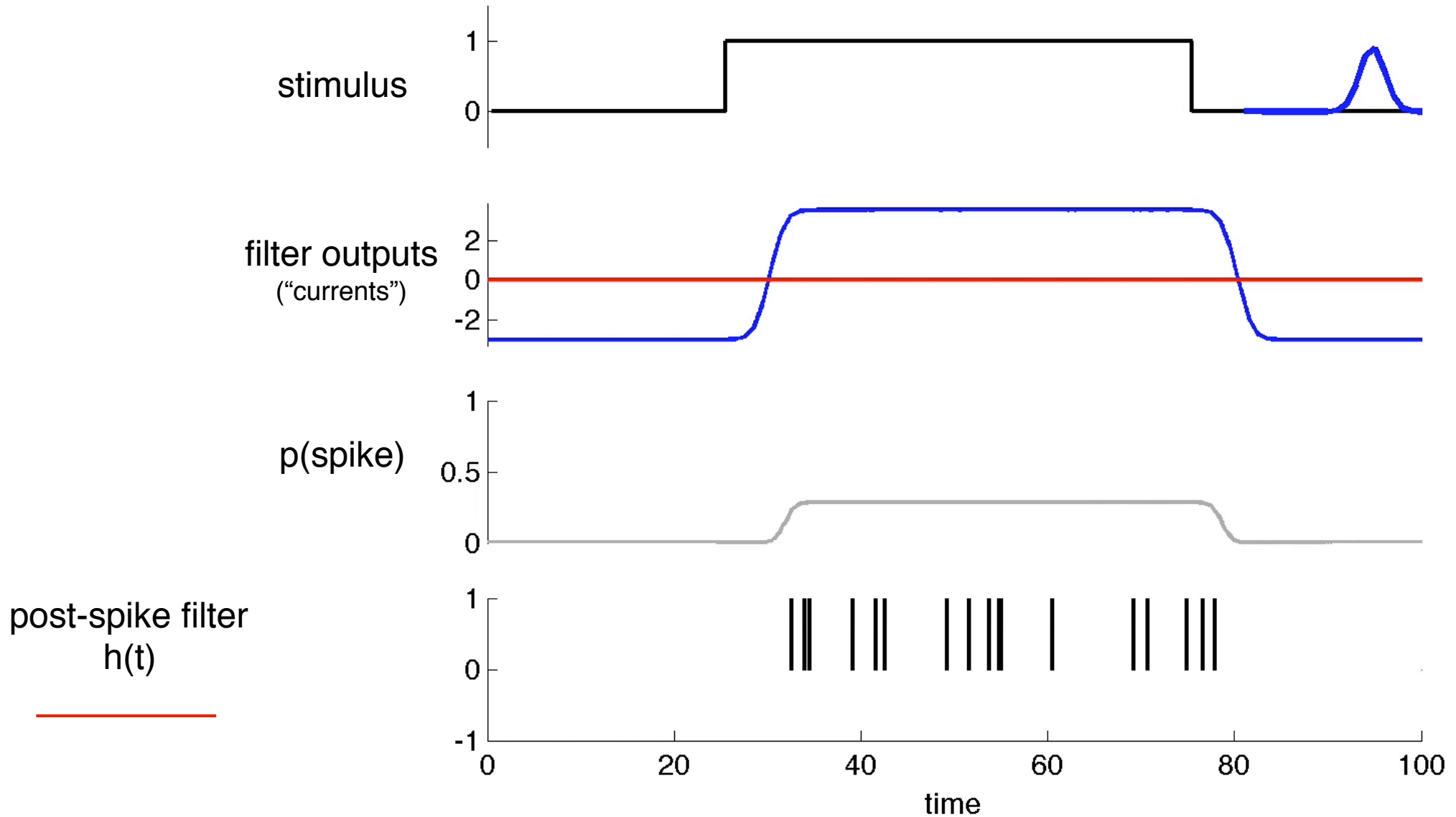


- interpretation: “soft-threshold” integrate-and-fire model



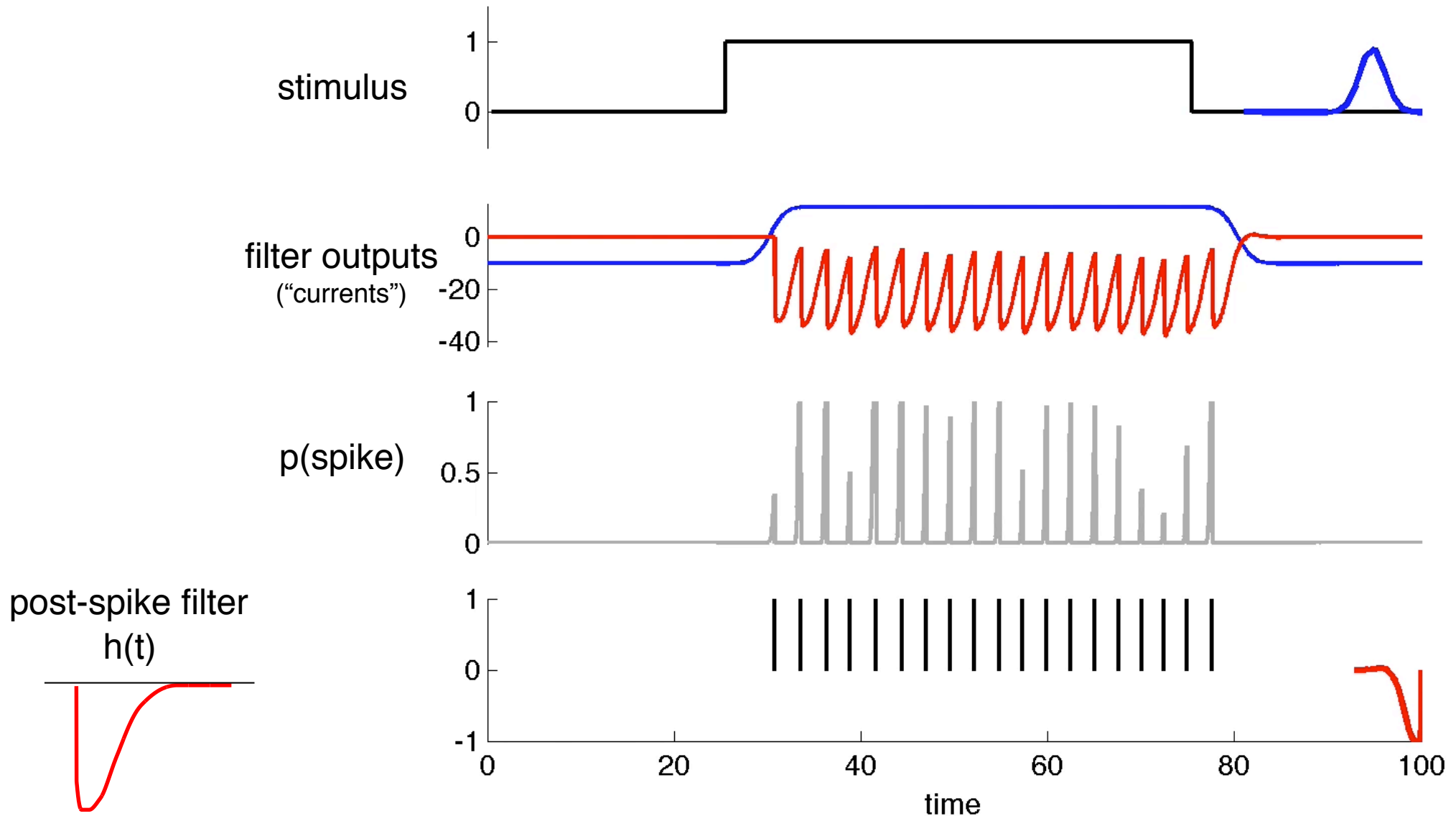
# GLM dynamic behaviors

- irregular spiking



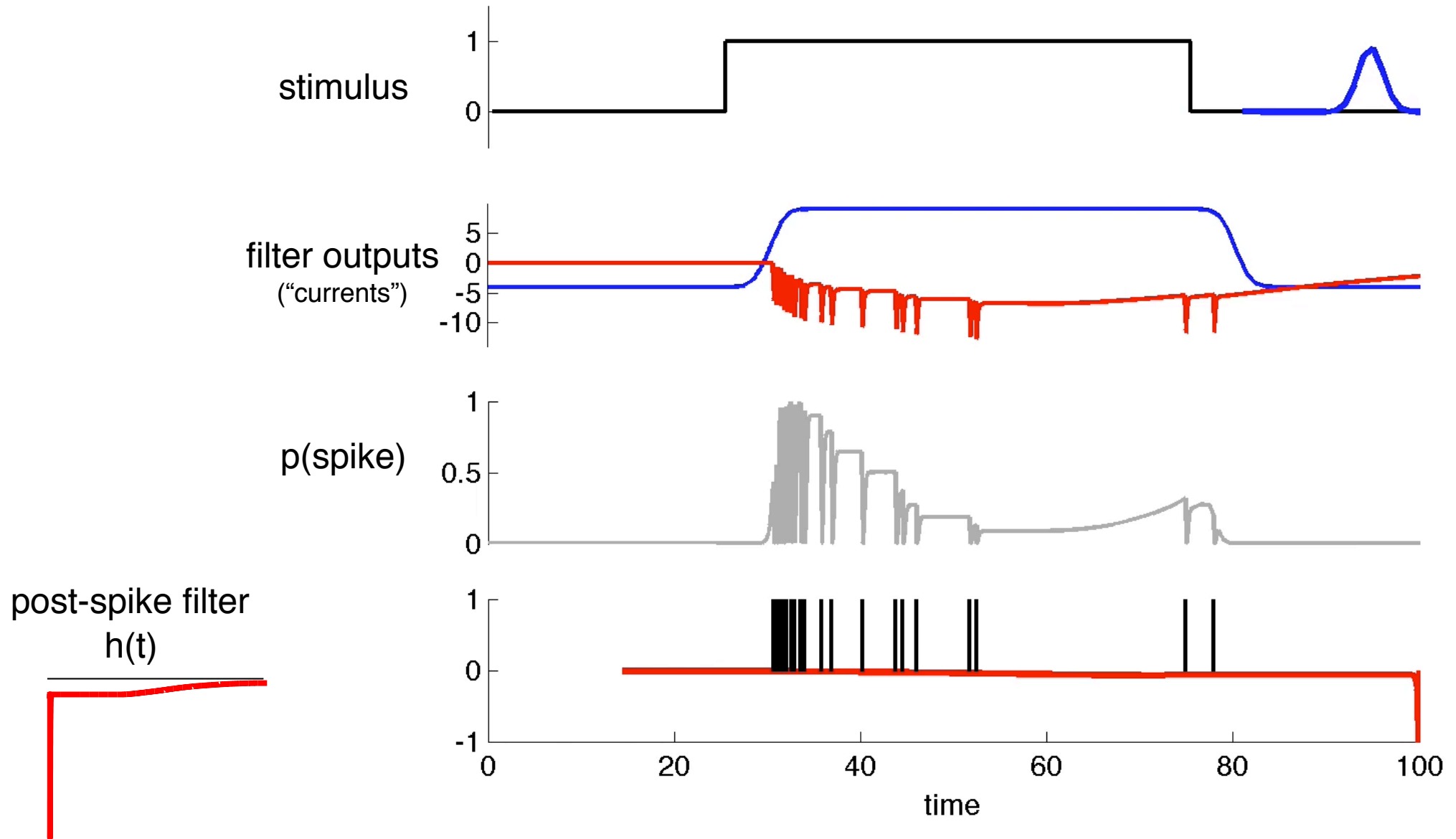
# GLM dynamic behaviors

- regular spiking



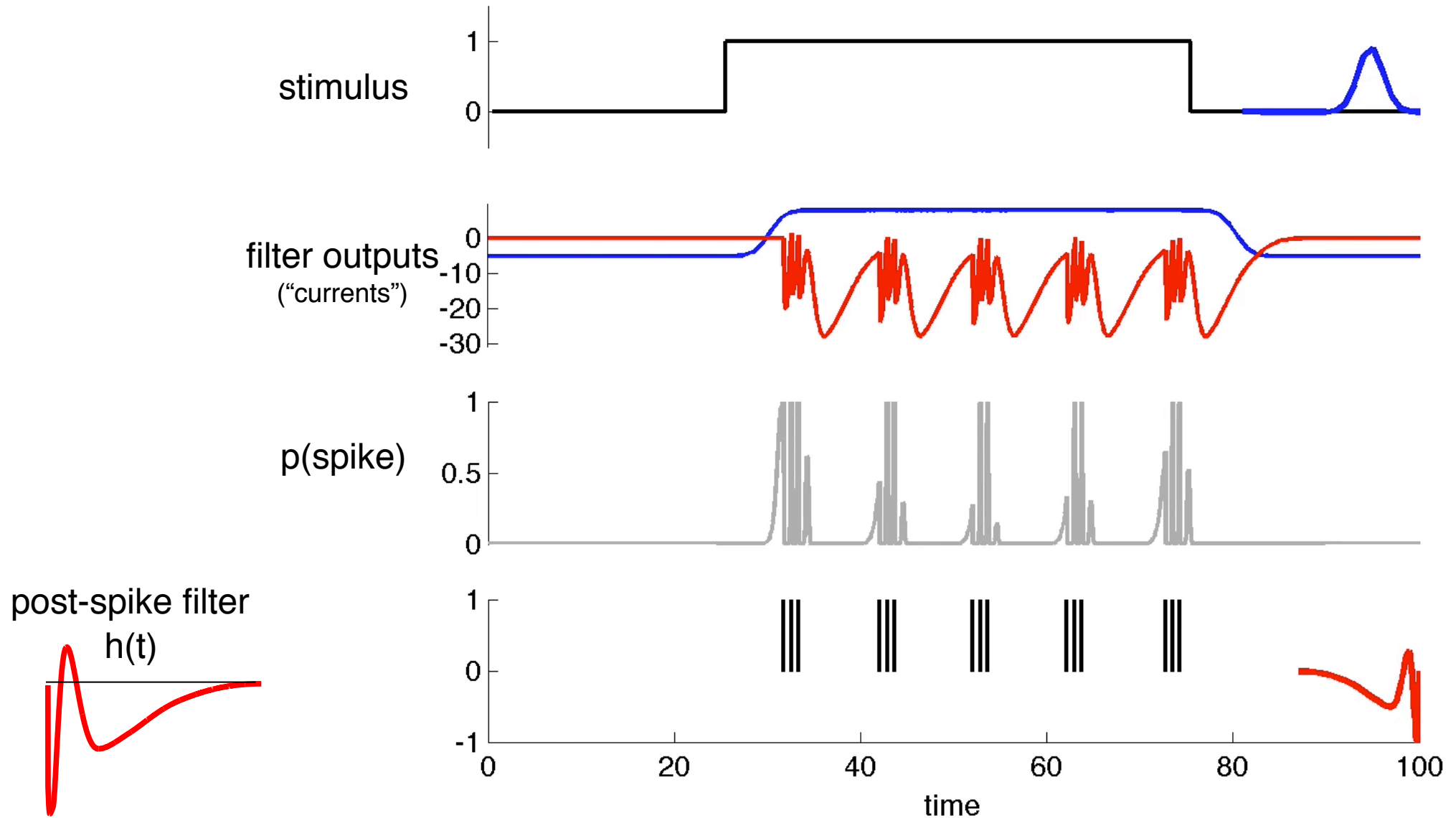
# GLM dynamic behaviors

- adaptation



# GLM dynamic behaviors

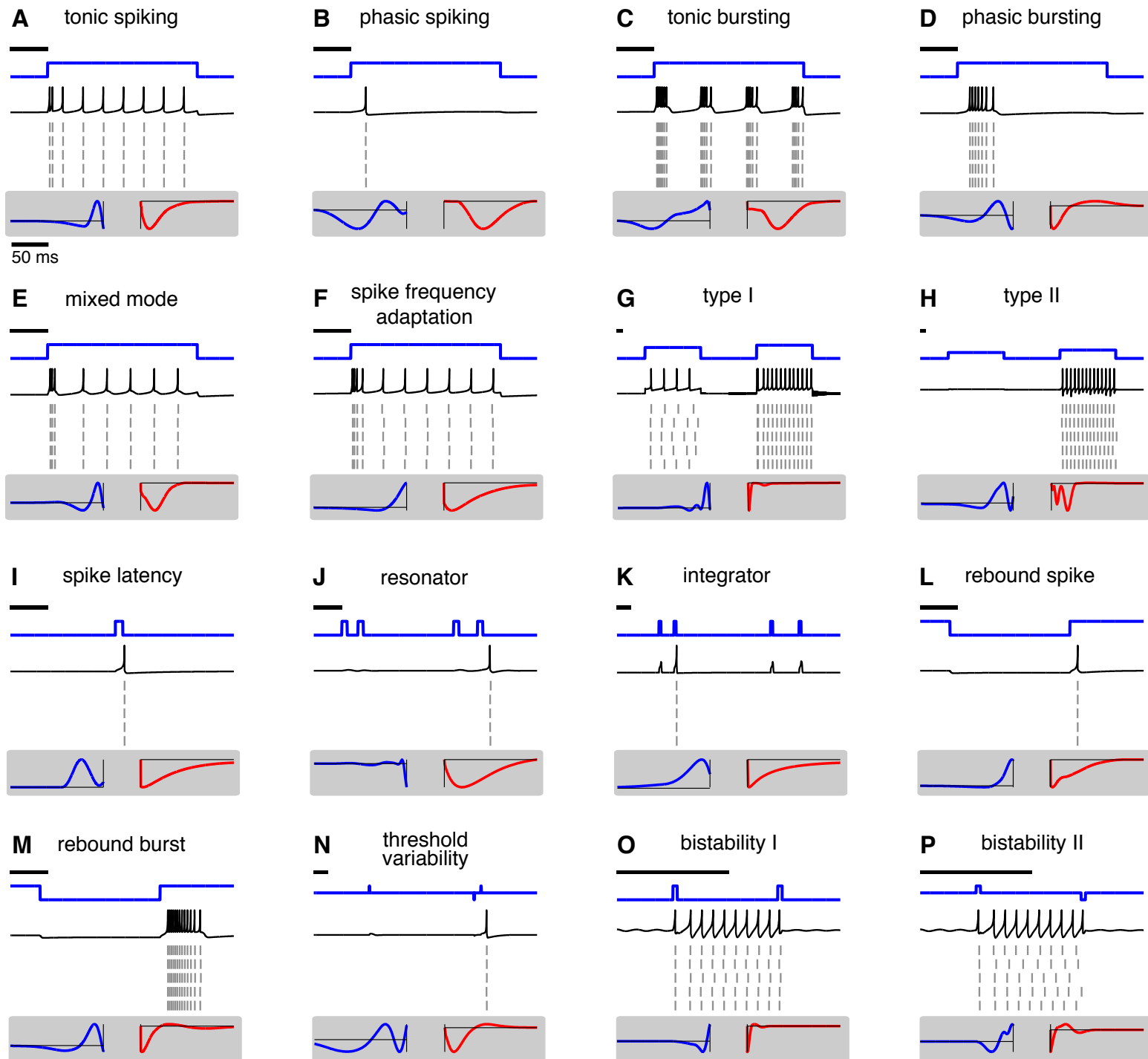
- bursting



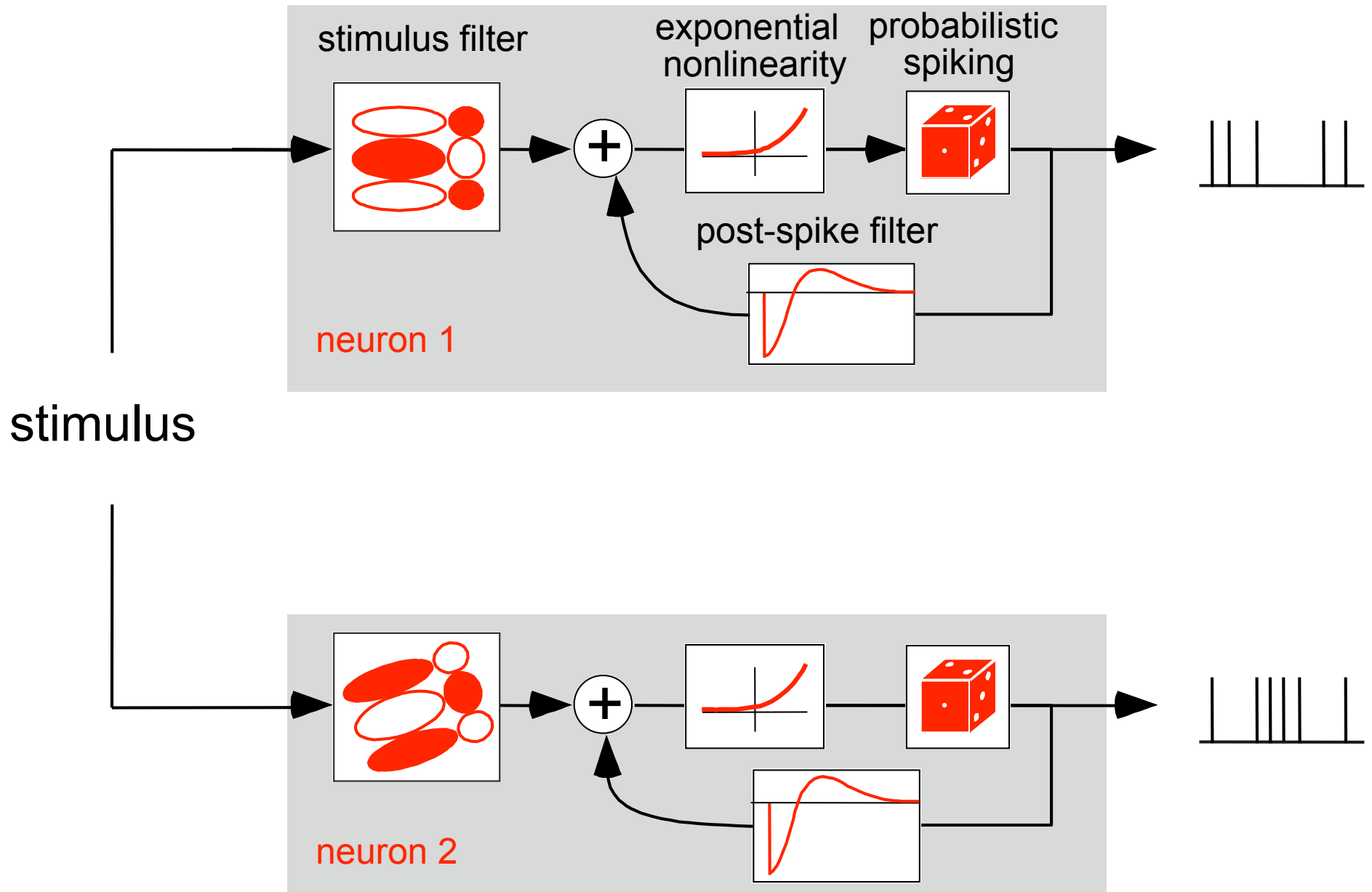
# GLM dynamic behaviors (from Izhikevich)

(Weber & Pillow 2017)

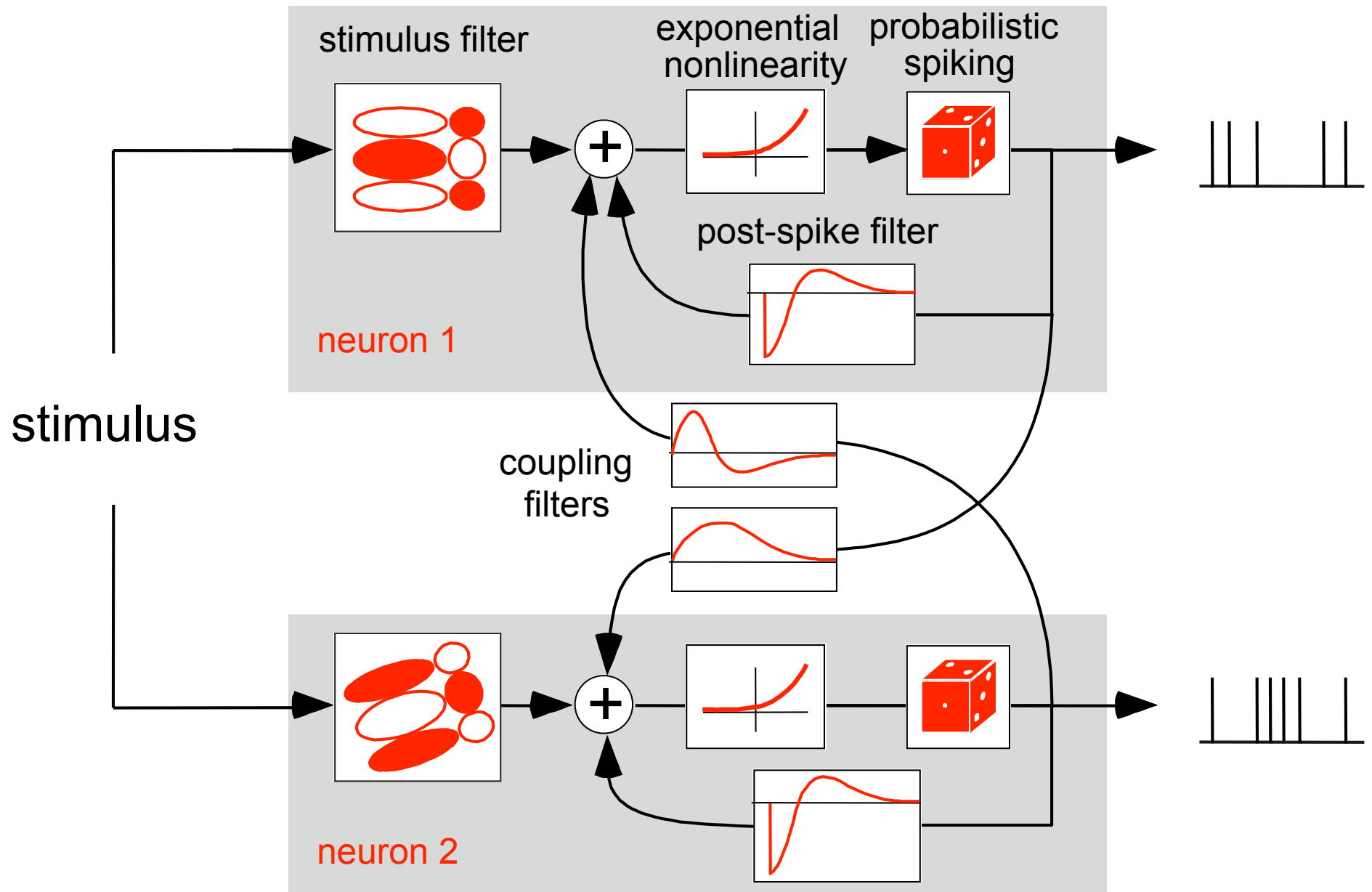
stimulus  
Izhikevich neuron  
GLM spikes  
GLM parameters



# multi-neuron GLM

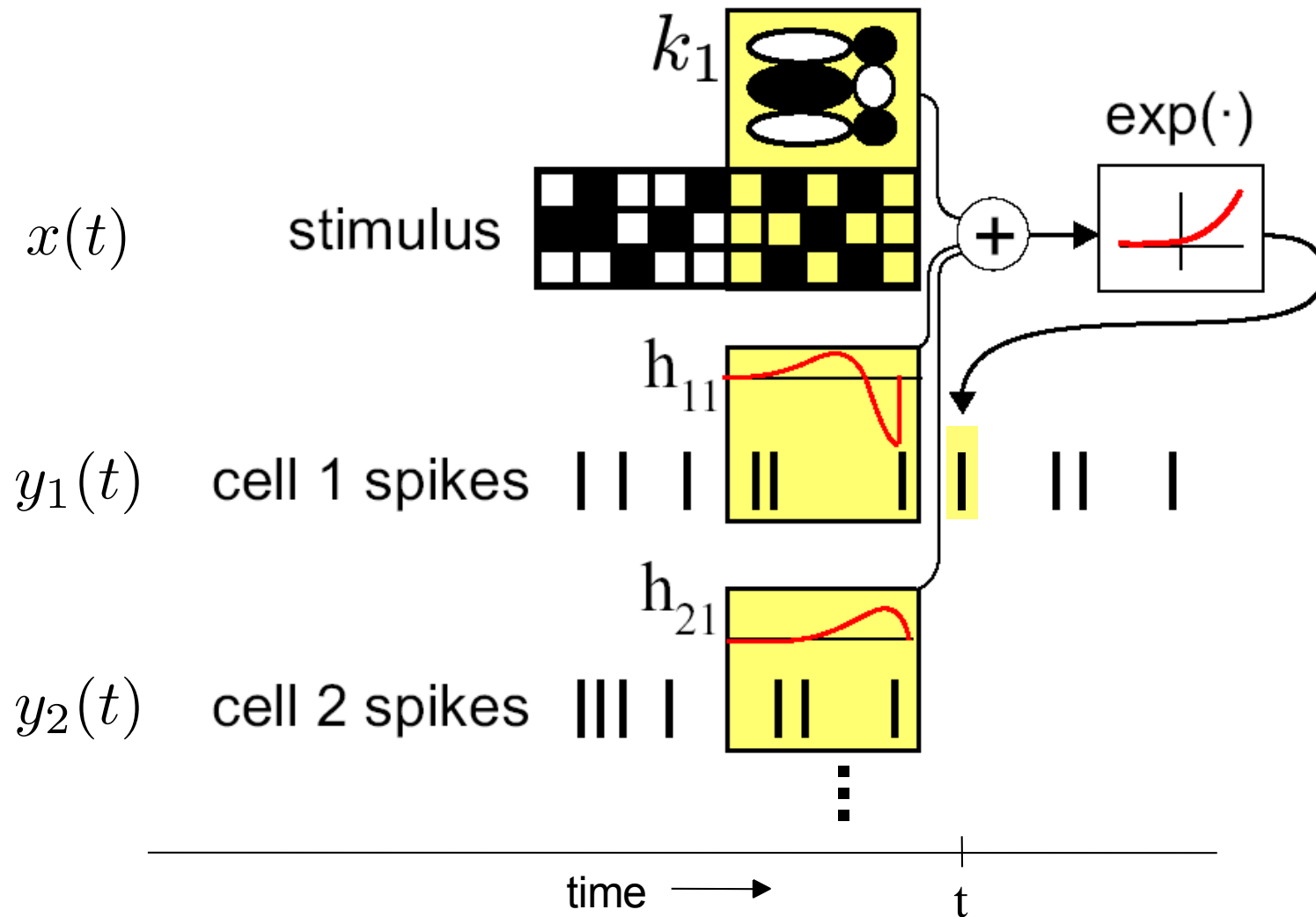


# multi-neuron GLM





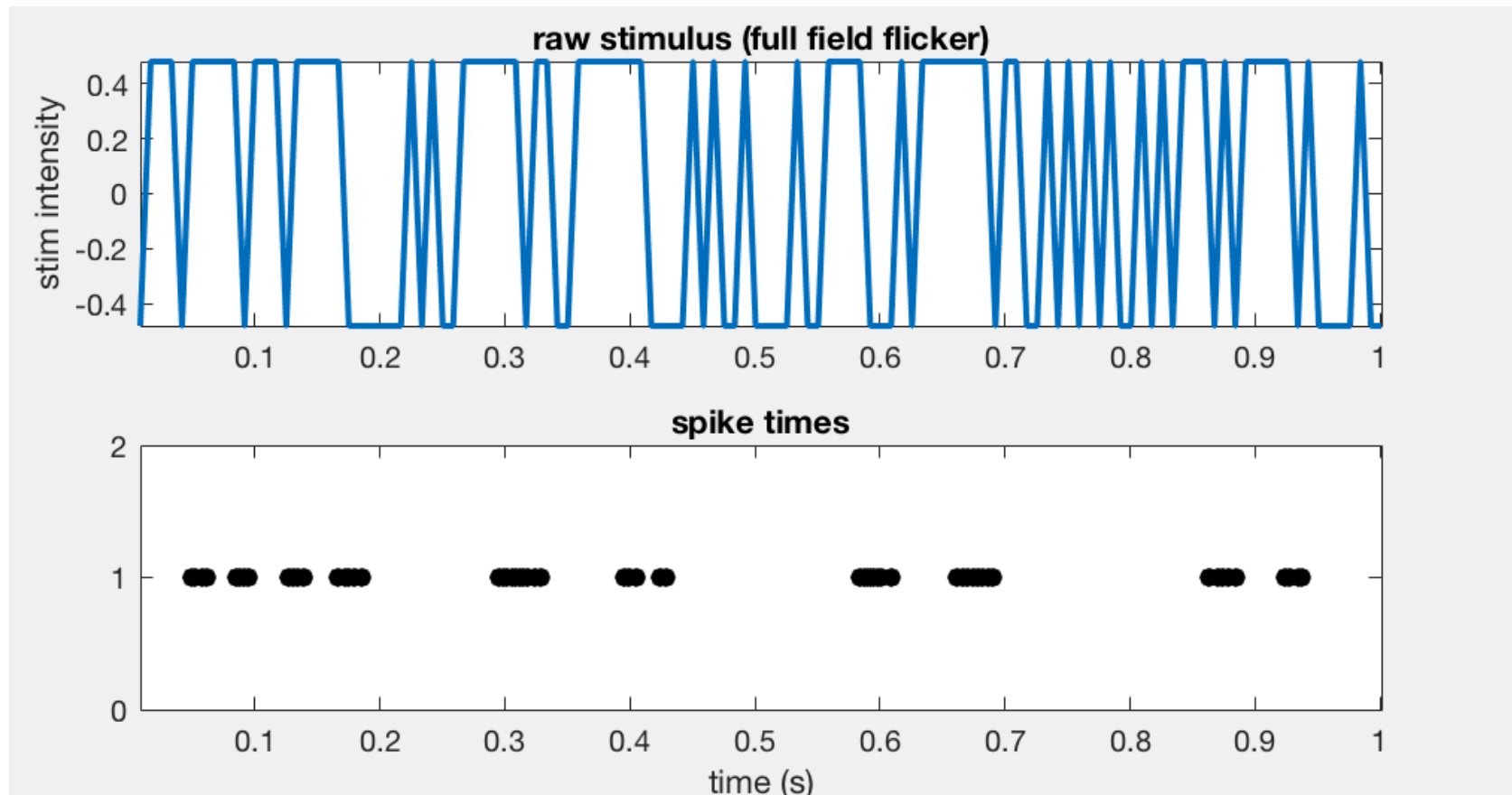
# GLM equivalent diagram:



spike rate  $\lambda_i(t) = \exp(k_i \cdot x(t) + \sum_j h_{ij} \cdot y_j(t))$

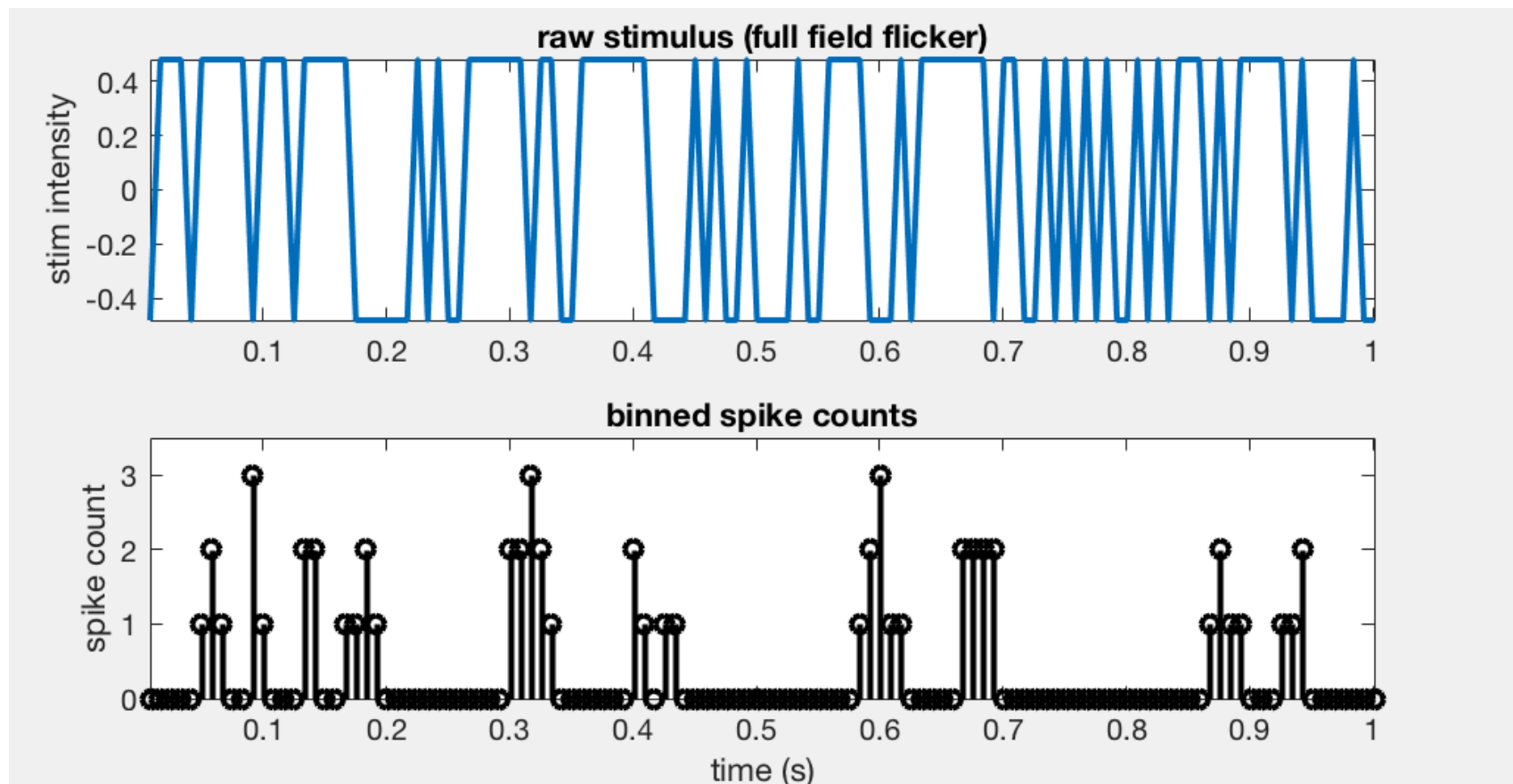
# Example dataset

- stimulus = binary flicker
- parasol retinal ganglion cell spike responses



# Example dataset

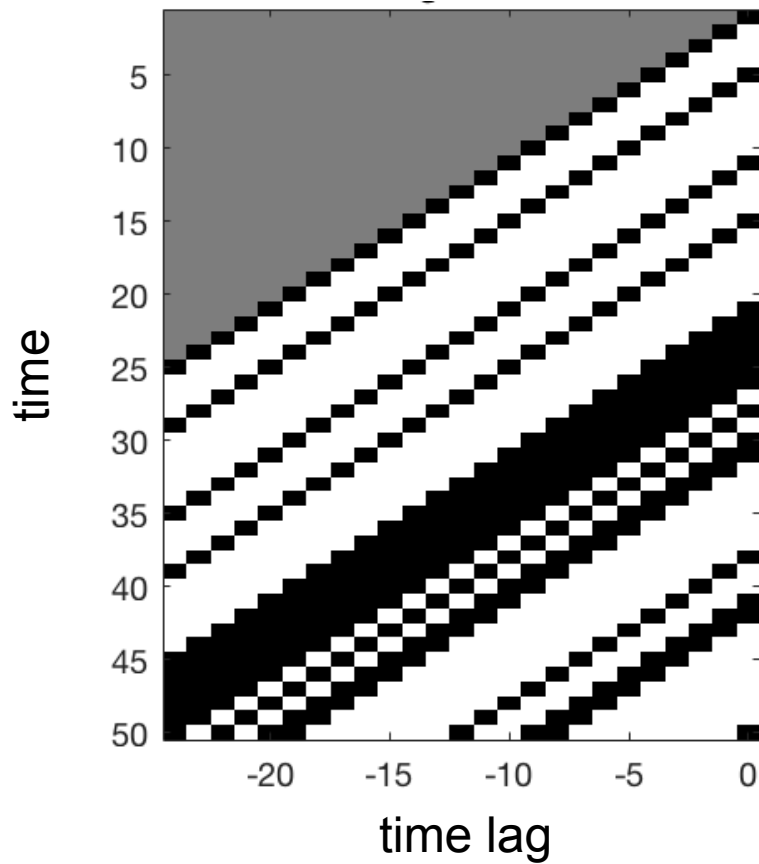
- stimulus = binary flicker
- parasol retinal ganglion cell spike responses



# Stimulus-only GLM

design matrix

$X$



model

$$P(Y|X)$$



spike response

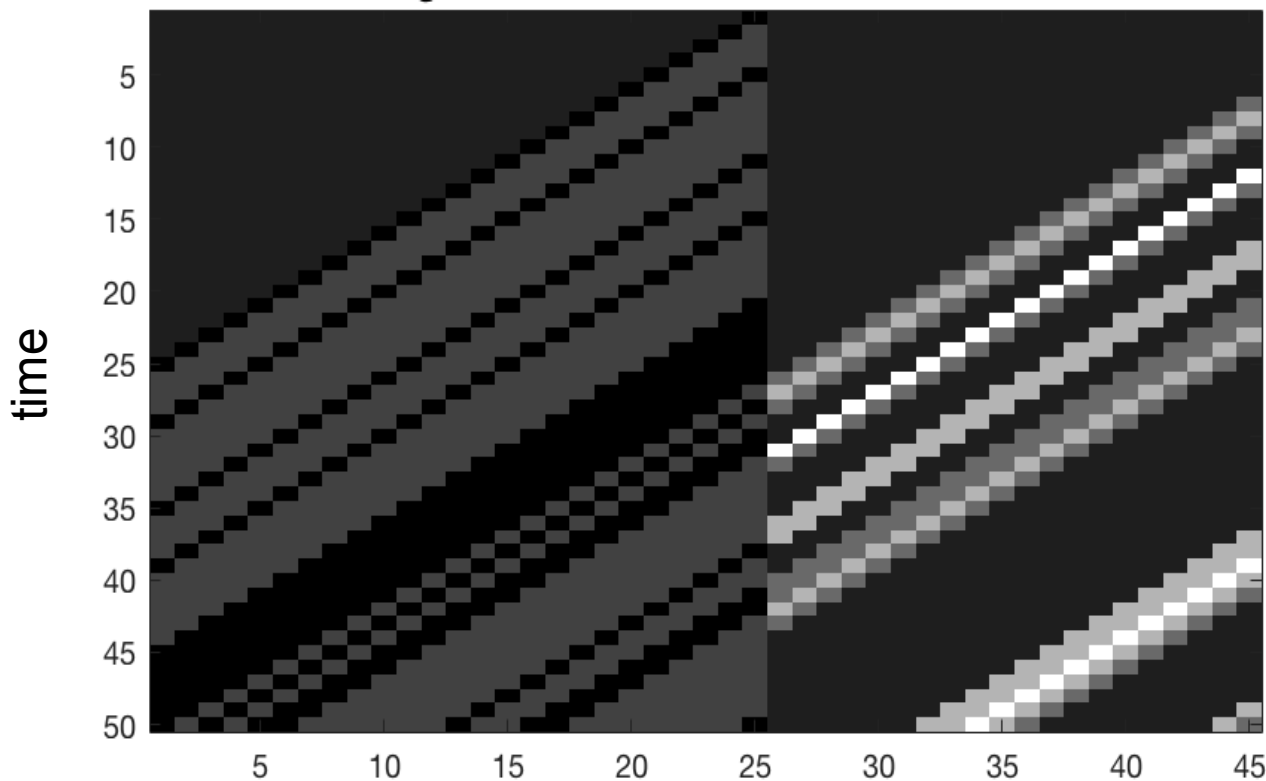
$Y$



# Stimulus + SpikeHistory GLM

design matrix

$X$



stimulus  
portion

spike-history  
portion

spike response

$Y$

model

$$P(Y|X)$$



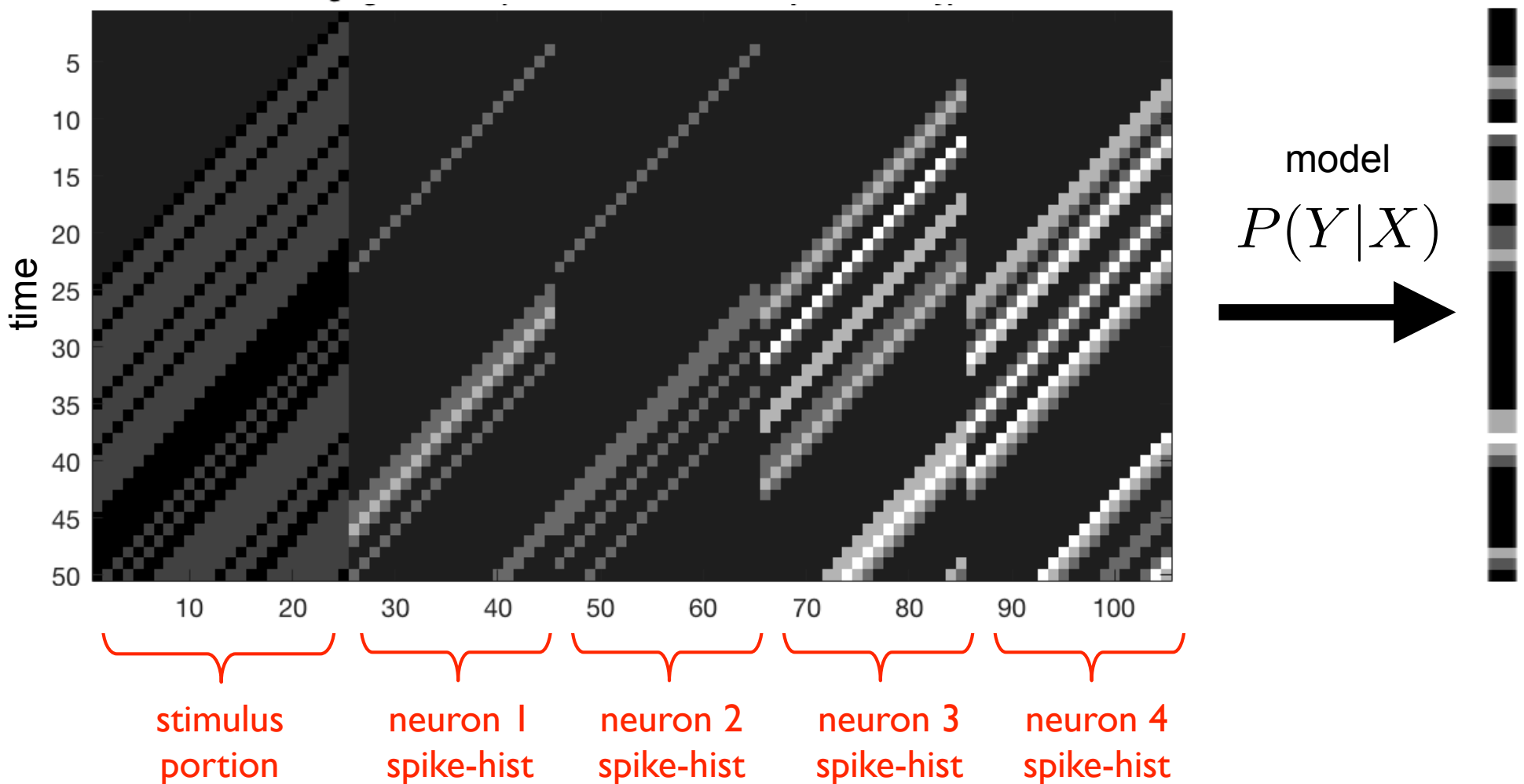
# Stimulus + History + 3 Neuron Coupling GLM

design matrix

spike response

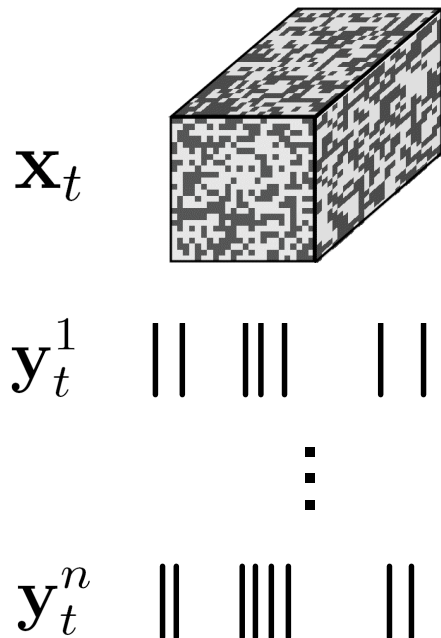
$X$

$Y$

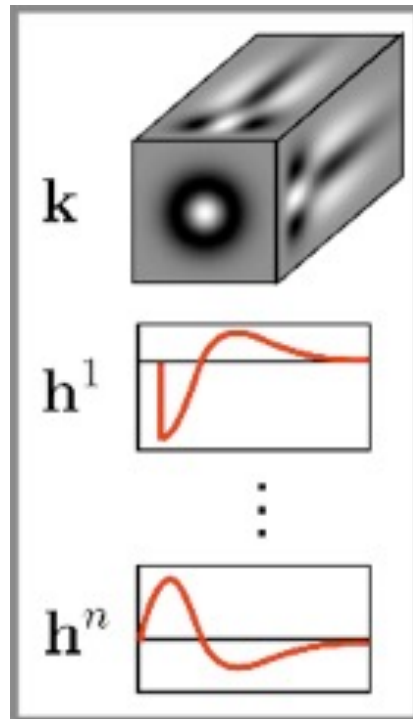


# Fitting: Maximum Likelihood

Data



GLM

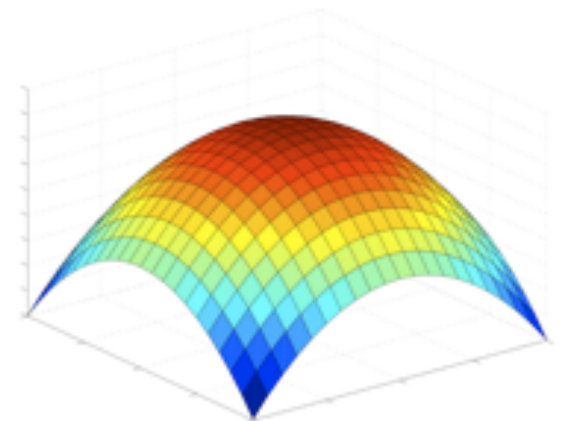


- maximize log-likelihood for filters  $\{\mathbf{k}, h_1, h_2, \dots, h_n\}$

firing rate:  $\lambda_t = f(\vec{x}_t \cdot \vec{k})$

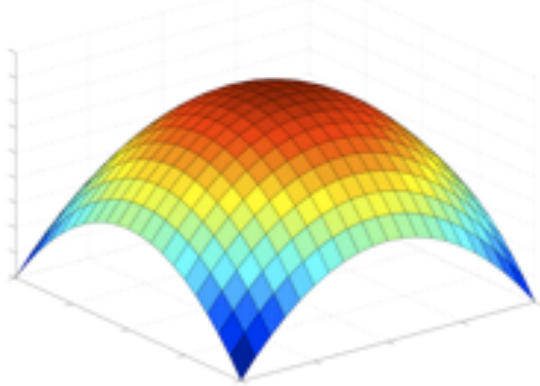
$$\log P(Y|X) = \sum_t y_t \log \lambda_t - \lambda_t$$

- log-likelihood is concave
- no local maxima [Paninski 04]



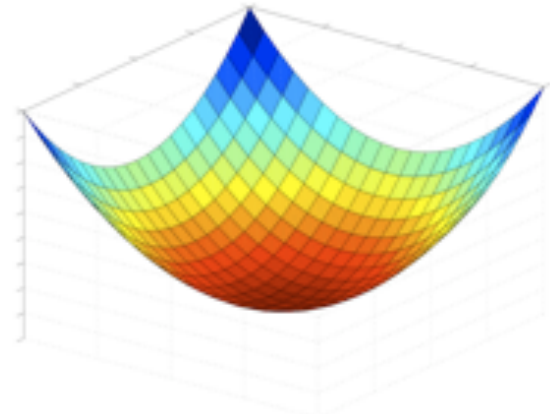
# convexity and concavity

concave



- everywhere downward curvature

convex



- everywhere upward curvature
- 
- maximizing concave function  $\iff$  minimizing a convex function
  - preclude existence of non-global local optima

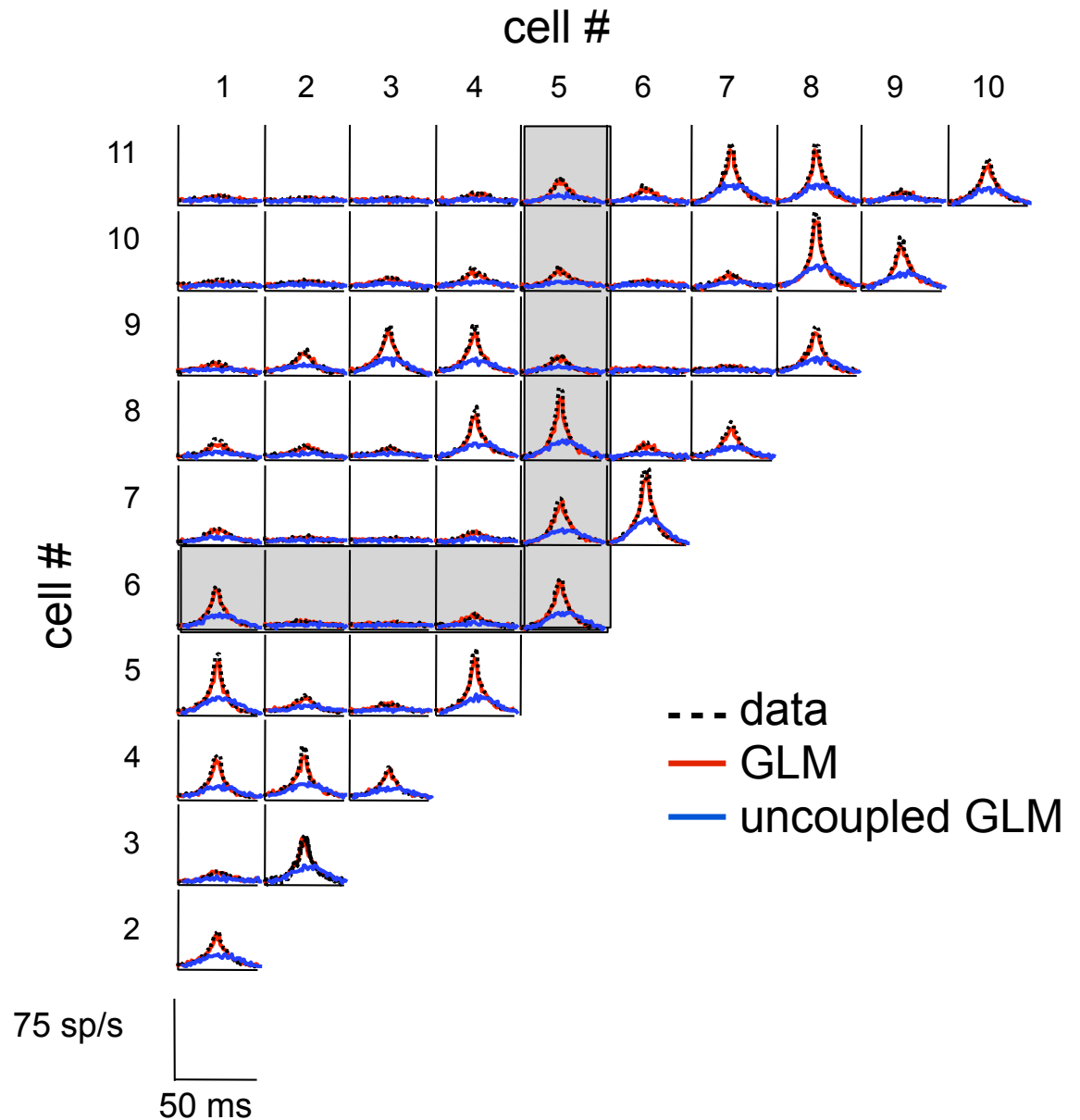
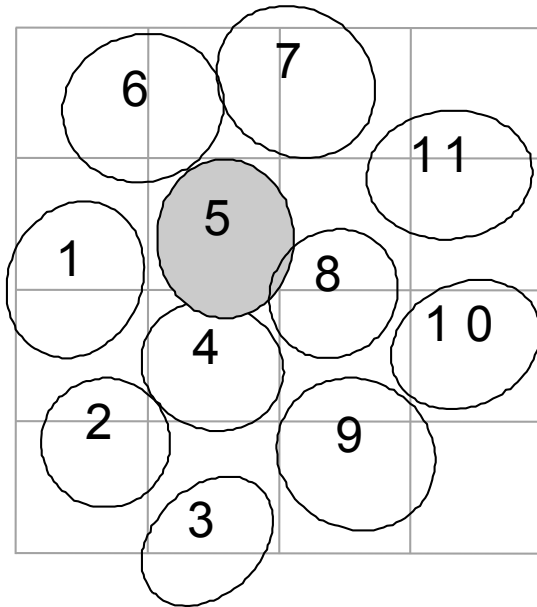


# capturing dependencies in multi-neuron responses

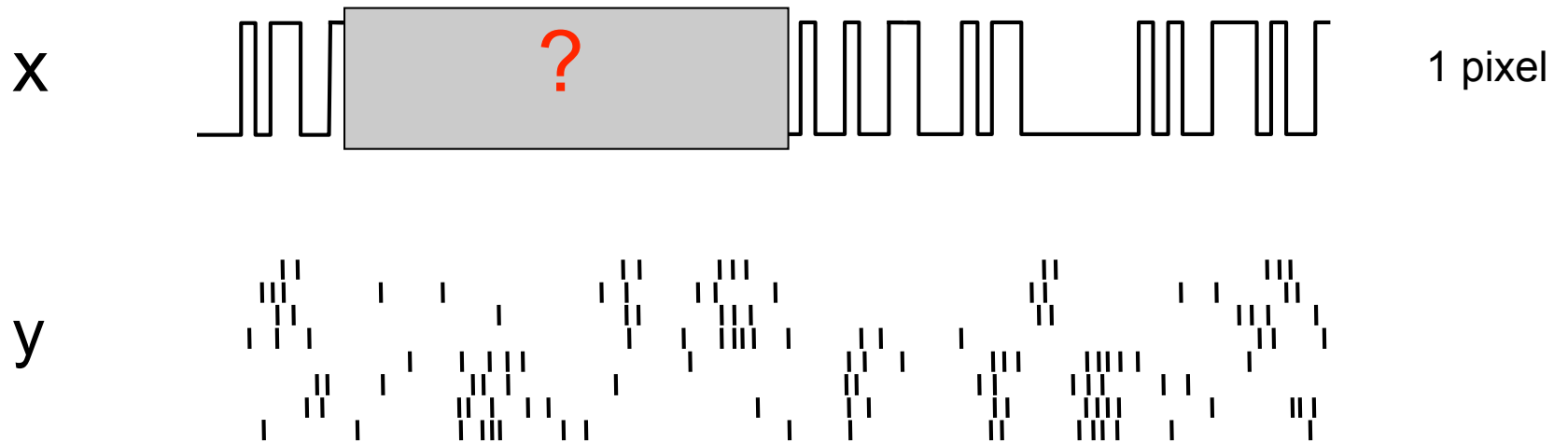
[Pillow et al 2008]

## cross-correlations

## retinal receptive fields

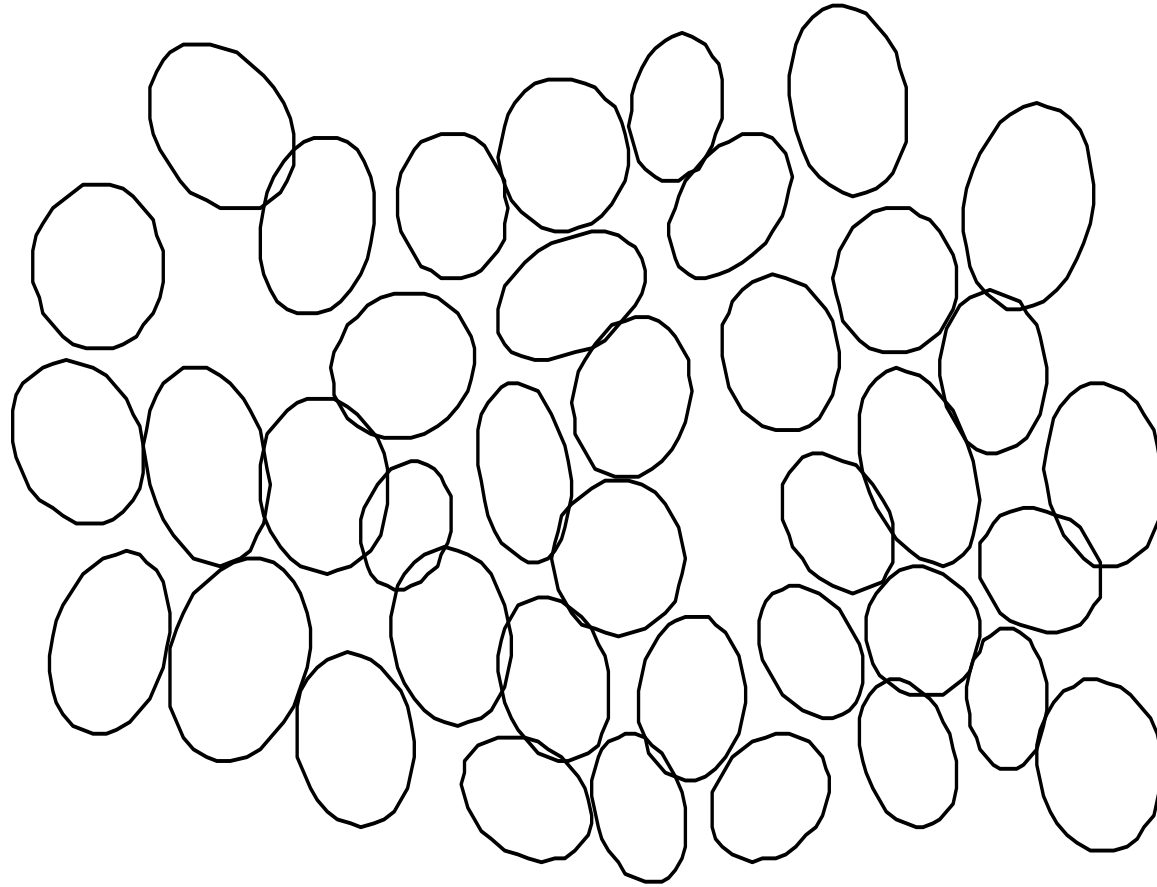


# Decoding



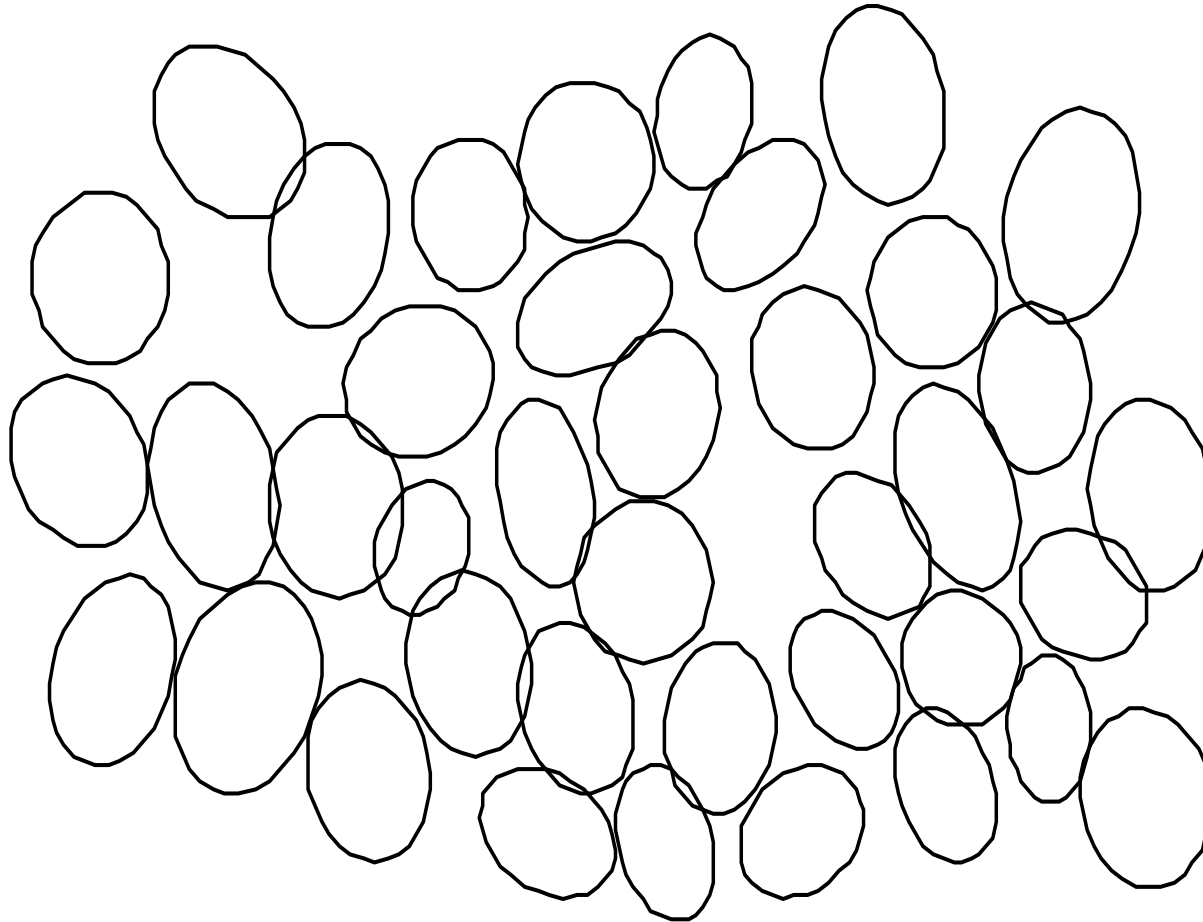
- estimate stimuli from the observed spike times
- tool for comparing different encoding models

# Decode: response 1



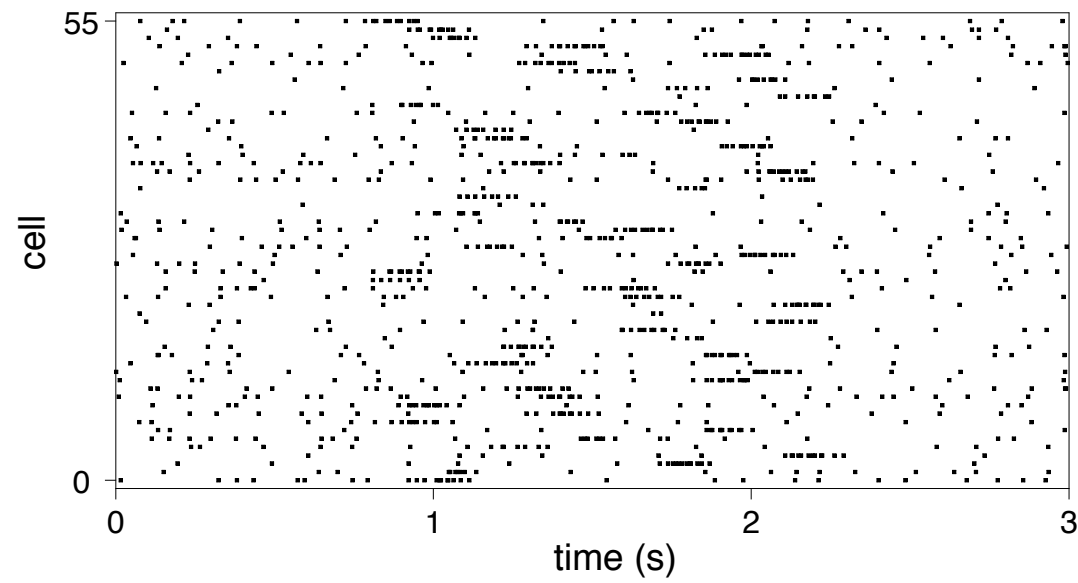
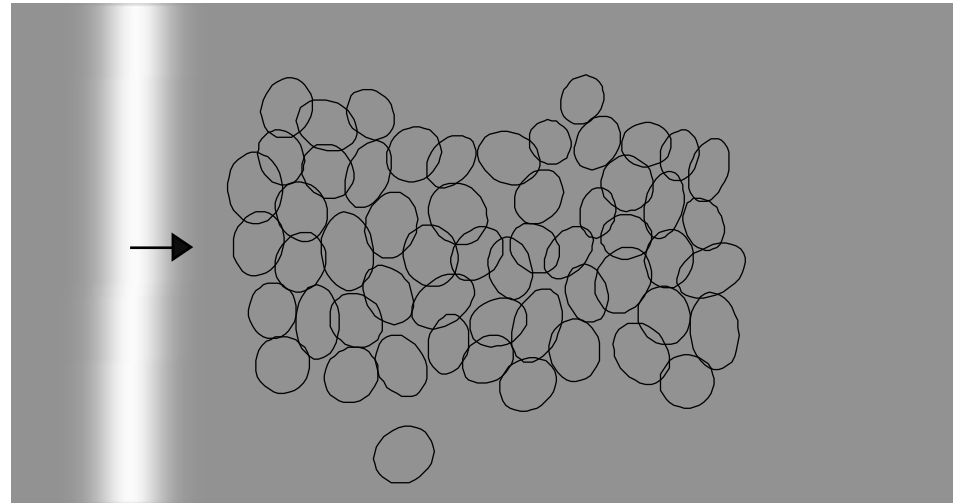
Q: what was the stimulus?

# Decode: response 2

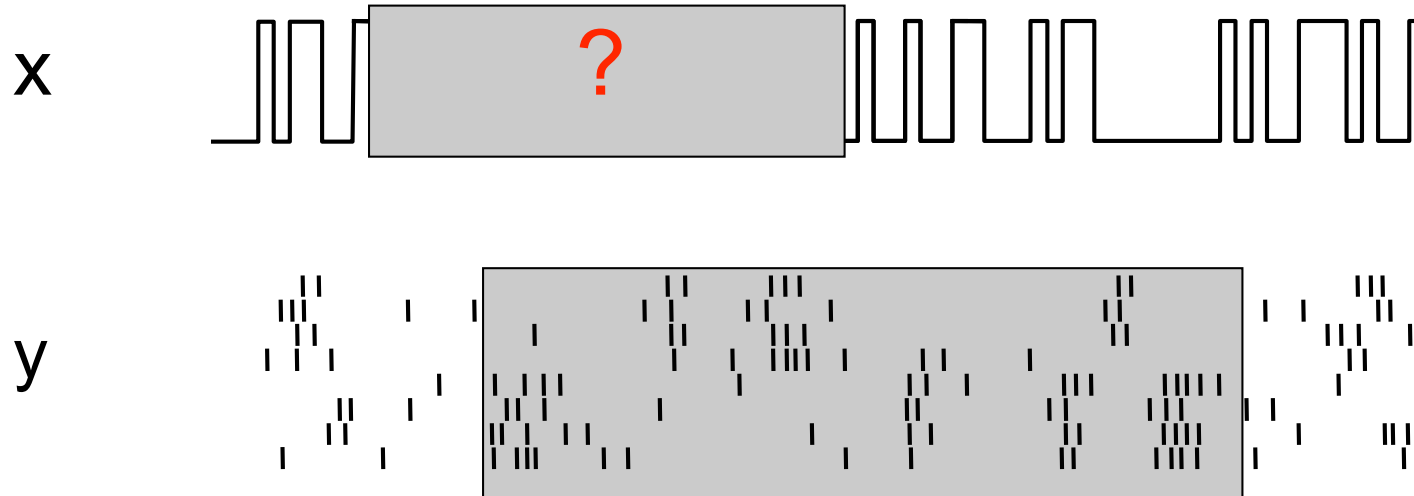


Q: what was the stimulus?

## Responses to Moving Bar



# Bayesian Decoding

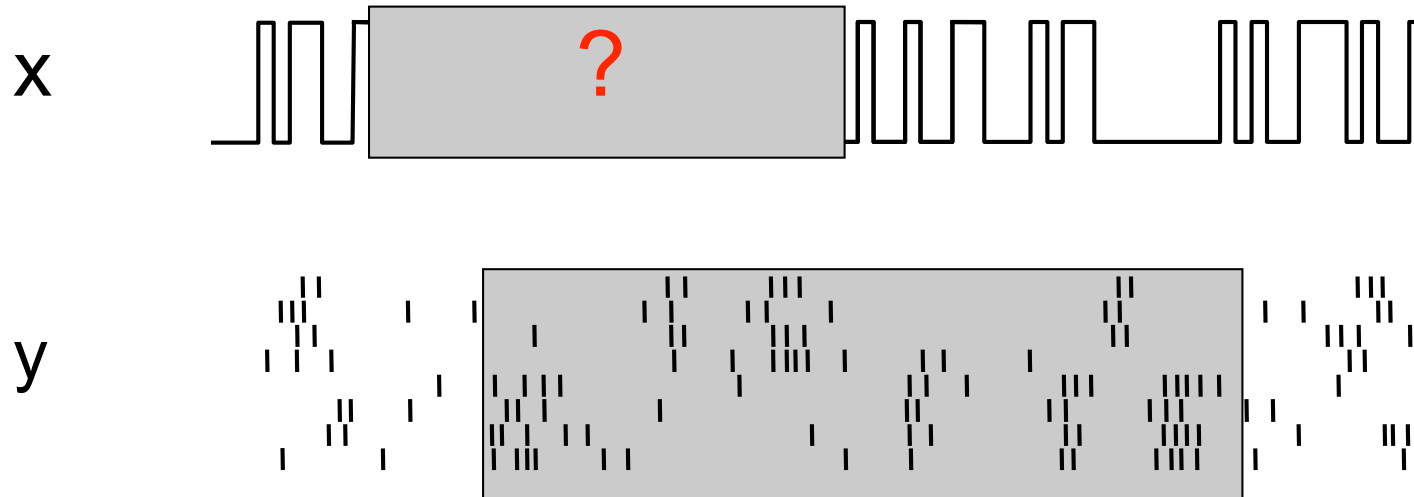


Bayes' rule:  $P(x|y) \propto P(y|x)P(x)$

posterior      likelihood      prior

Arrows point from the labels to the corresponding terms in the equation: posterior to  $P(x|y)$ , likelihood to  $P(y|x)$ , and prior to  $P(x)$ .

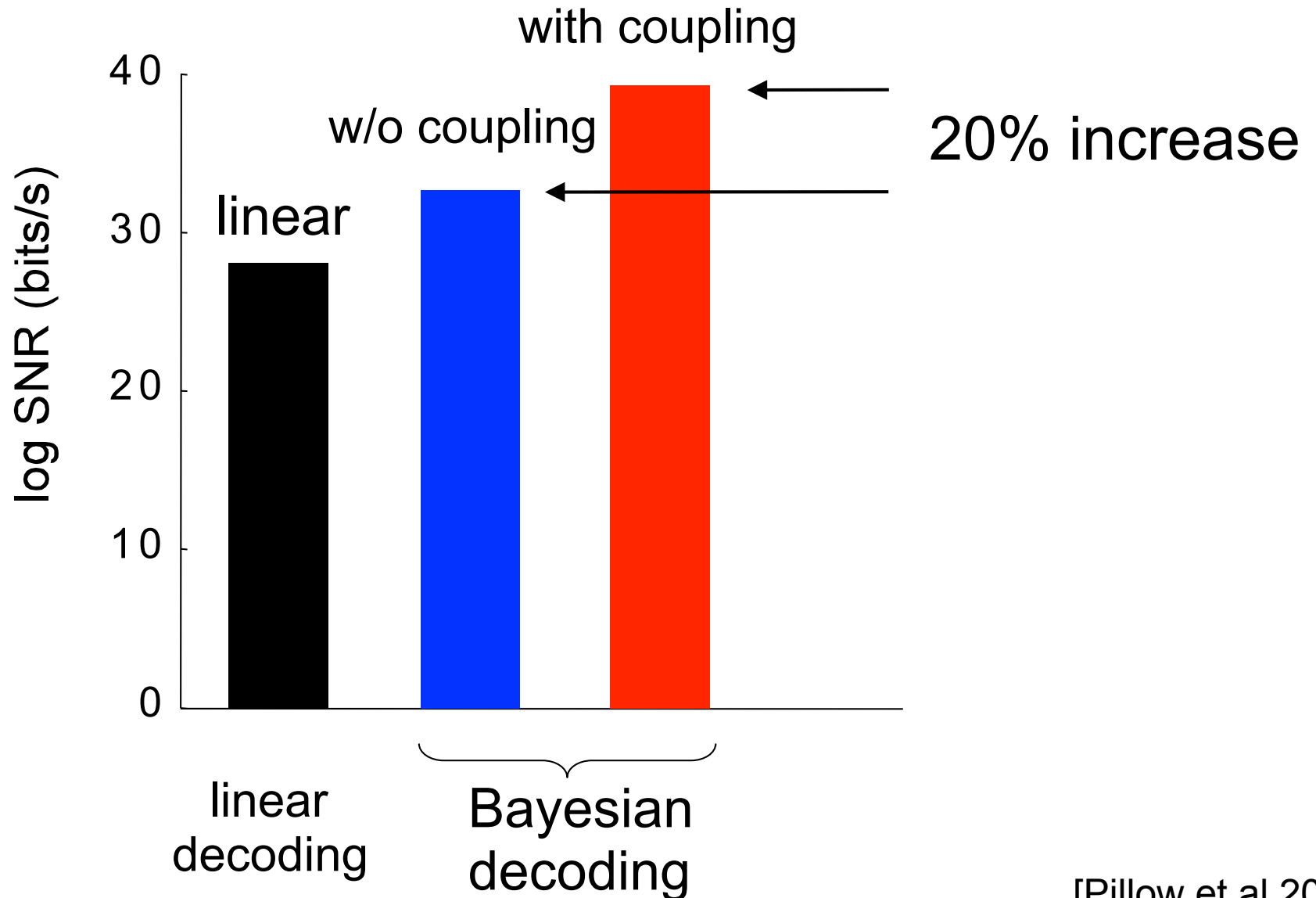
# Bayesian Decoding



Bayes' rule:  $P(x|y) \propto \underbrace{P(y|x)}_{\text{likelihood}} P(x)_{\text{prior}}$

$P(y_1|x) \cdots P(y_n|x)$  “independent” (uncoupled GLM)
 vs.
 $P(y_1, y_2, \dots, y_n|x)$  “joint encoding” (coupled GLM)

# Decoding Comparison





# Regularization

# Modern statistics

- more dimensions than samples  $D \geq N$

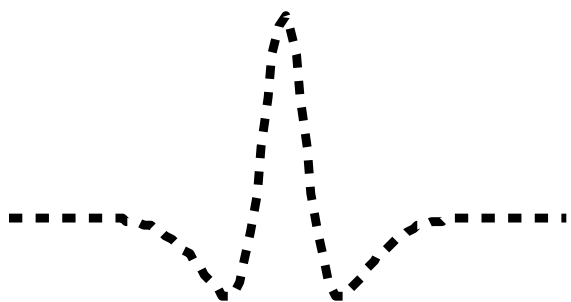
$$\begin{array}{c} \text{N} \\ \text{observations} \end{array} \left\{ \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right. = \begin{array}{c} \text{D regressors} \\ \left[ \begin{array}{ccc} \text{---} & \vec{x}_1 & \text{---} \\ & & \\ & & \\ \text{---} & \vec{x}_N & \text{---} \end{array} \right] \end{array} \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix} + \text{noise}$$

- fewer equations than unknowns!
- no unique solution

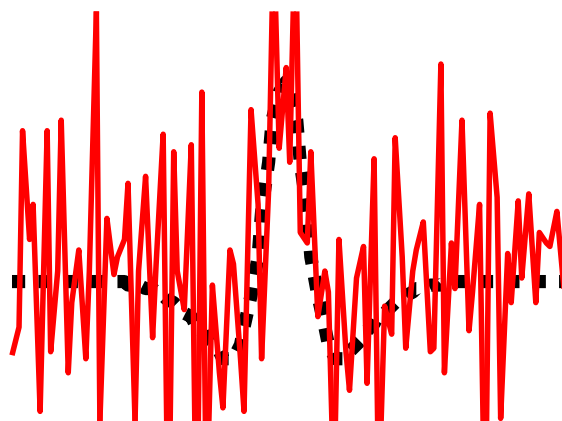
# Simulated Example

- 100-element filter ( $D=100$ )
- 100 noisy samples ( $N=100$ )

true  $\mathbf{w}$



maximum likelihood



*maximize*

$$\log p(data|\mathbf{w})$$

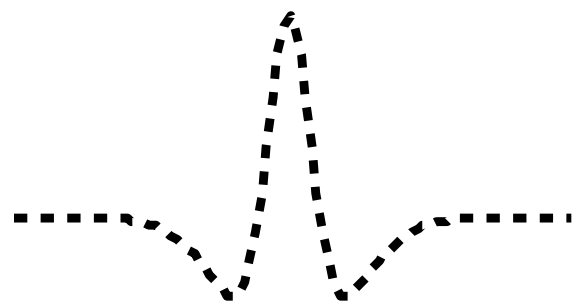
“overfitting” - parameters fit to details in the training data that are not useful for predicting new data

# Simulated Example

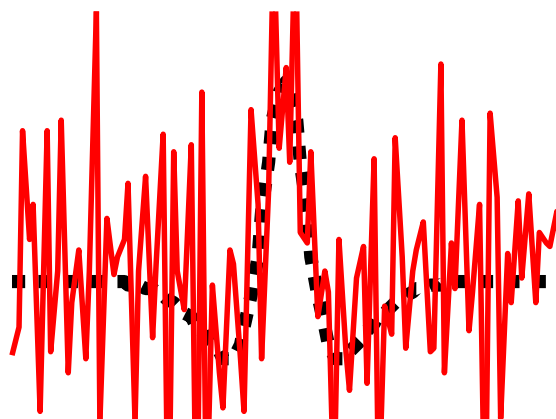
- 100-element filter ( $D=100$ )
- 100 noisy samples ( $N=100$ )

$$\hat{\mathbf{w}} = (X^\top X + \lambda I)^{-1} X^\top Y$$

true  $\mathbf{w}$



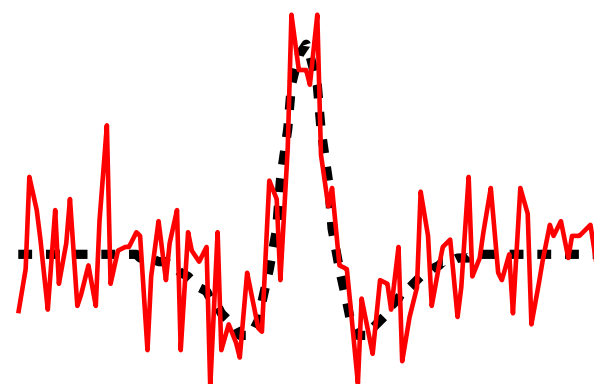
maximum likelihood



*maximize*

$$\log p(\text{data}|\mathbf{w})$$

“ridge regression”



*maximize*

$$\log p(\text{data}|\mathbf{w}) - \underbrace{\lambda \sum w_i^2}_{\text{penalty on big weights}}$$

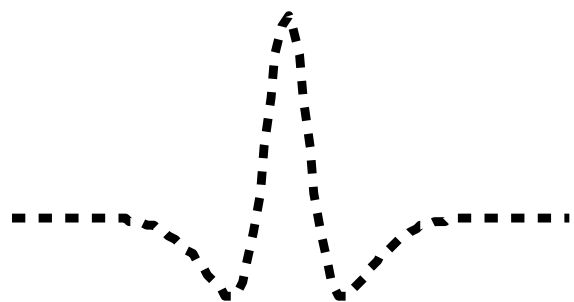
penalty on  
big weights

- biased, but gives improved performance for appropriate choice of  $\lambda$  (James & Stein 1960)

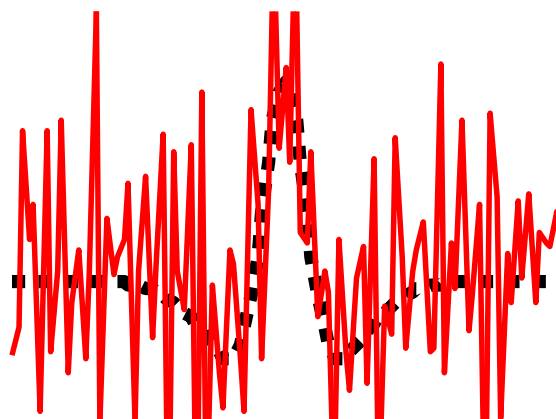
# Simulated Example

- 100-element filter ( $D=100$ )
- 100 noisy samples ( $N=100$ )

true  $\mathbf{w}$

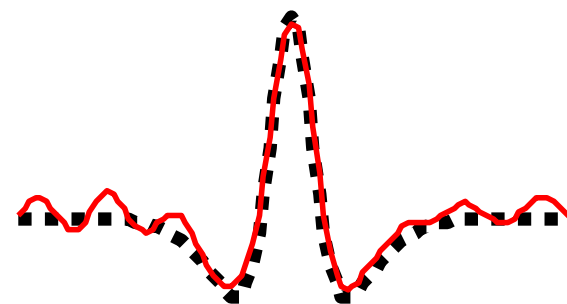


maximum likelihood



*maximize*  
 $\log p(\text{data}|\mathbf{w})$

“smoothed”



*maximize*  
 $\log p(\text{data}|\mathbf{w})$   
 $-\lambda \underbrace{\sum (w_i - w_{i-1})^2}_{\text{smoothness penalty}}$

**Q:** how to set the regularization strength  $\lambda$ ?

**Simplest answer:** use cross-validation!

# GLM tutorial (matlab):

code: <https://github.com/pillowlab/GLMspiketraintutorial>

data: available on request from [pillow@princeton.edu](mailto:pillow@princeton.edu)

- **tutorial1\_PoissonGLM.m** - fitting of a linear-Gaussian GLM and Poisson GLM (aka LNP model) to RGC neurons stimulated with temporal white noise stimulus.
- **tutorial2\_spikehistcoupledGLM.m** - fitting of a Poisson GLM with spike-history and coupling between neurons.
- **tutorial3\_regularization\_linGauss.m** - regularizing linear-Gaussian model parameters using maximum a posteriori (MAP) estimation under two kinds of priors:
  - (1) ridge regression (aka "L2 penalty");
  - (2) L2 smoothing prior (aka "graph Laplacian").
- **tutorial4\_regularization\_PoissonGLM.m** - MAP estimation of Poisson-GLM parameters using same two priors as in tutorial3.

# GLM summary

- linear (“dim reduction”) + nonlinear + noise
- incorporate spike-history via “spike history” filter
- rich dynamical properties: refractoriness, bursting, adaptation
- incorporate correlations between neurons via “coupling” filters
- flexible tool for encoding & decoding analyses
- regularize to reduce overfitting (essential w/ correlated stimuli)

# Beyond GLM



# polynomial models

Lee & Schetzen 1965  
Marmarelis & Naka 1972  
Korenberg & Hunter 1986

## Volterra / Wiener Kernels

Taylor series expansion of a function  $f(\mathbf{x})$  in  $n$  dimensions

$$y = k_0 + \vec{k}_1 \cdot \vec{x} + \vec{x}^t K_2 \vec{x} + K_3 \cdot \vec{x}^3 + \dots$$

↓  
const



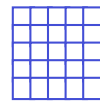
# parameters: 1

↓  
vector



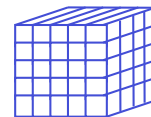
$n$   
(20)

↓  
matrix



$n^2$   
(400)

↓  
3-tensor

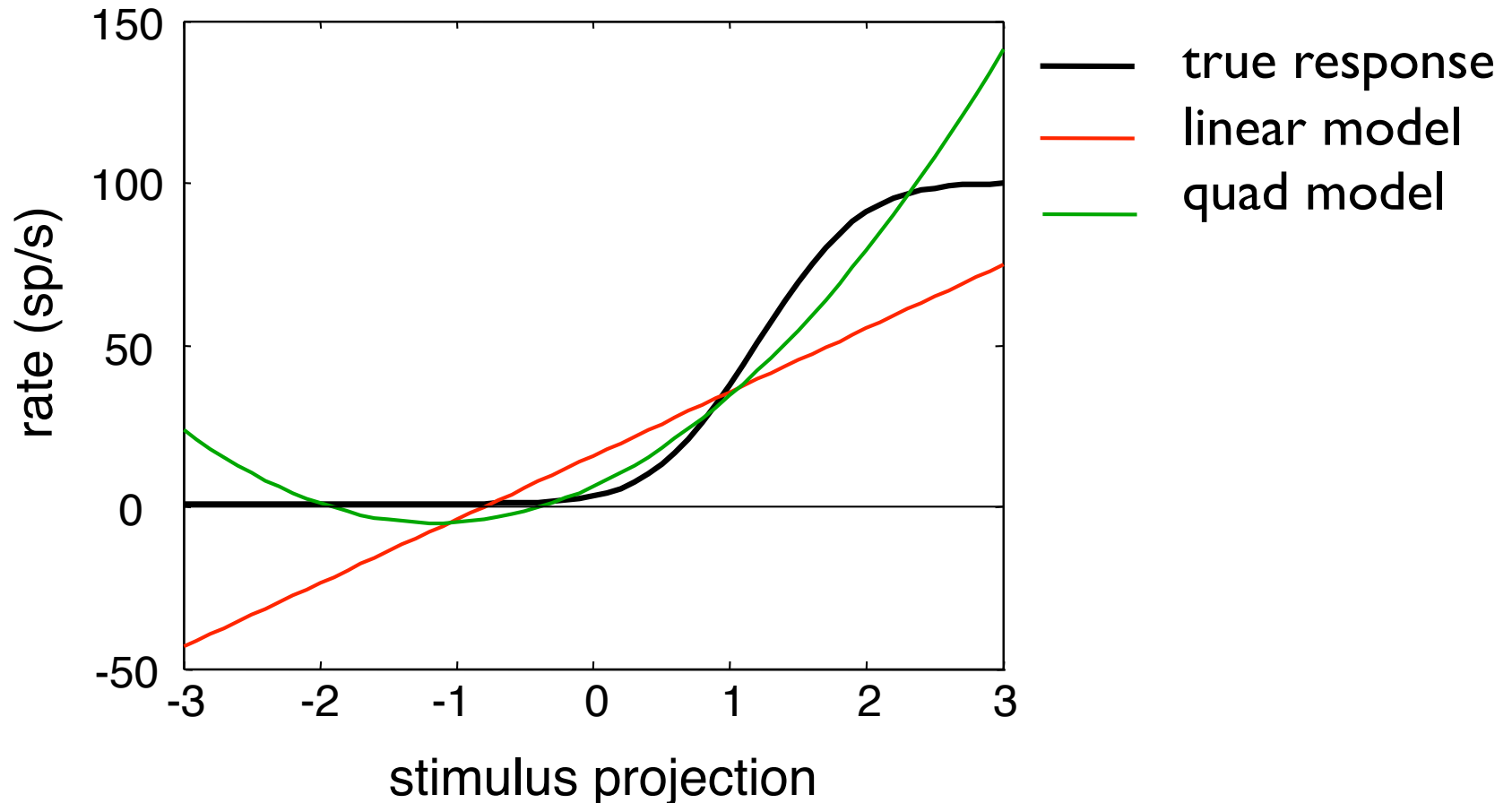


$n^3$   
(8000)

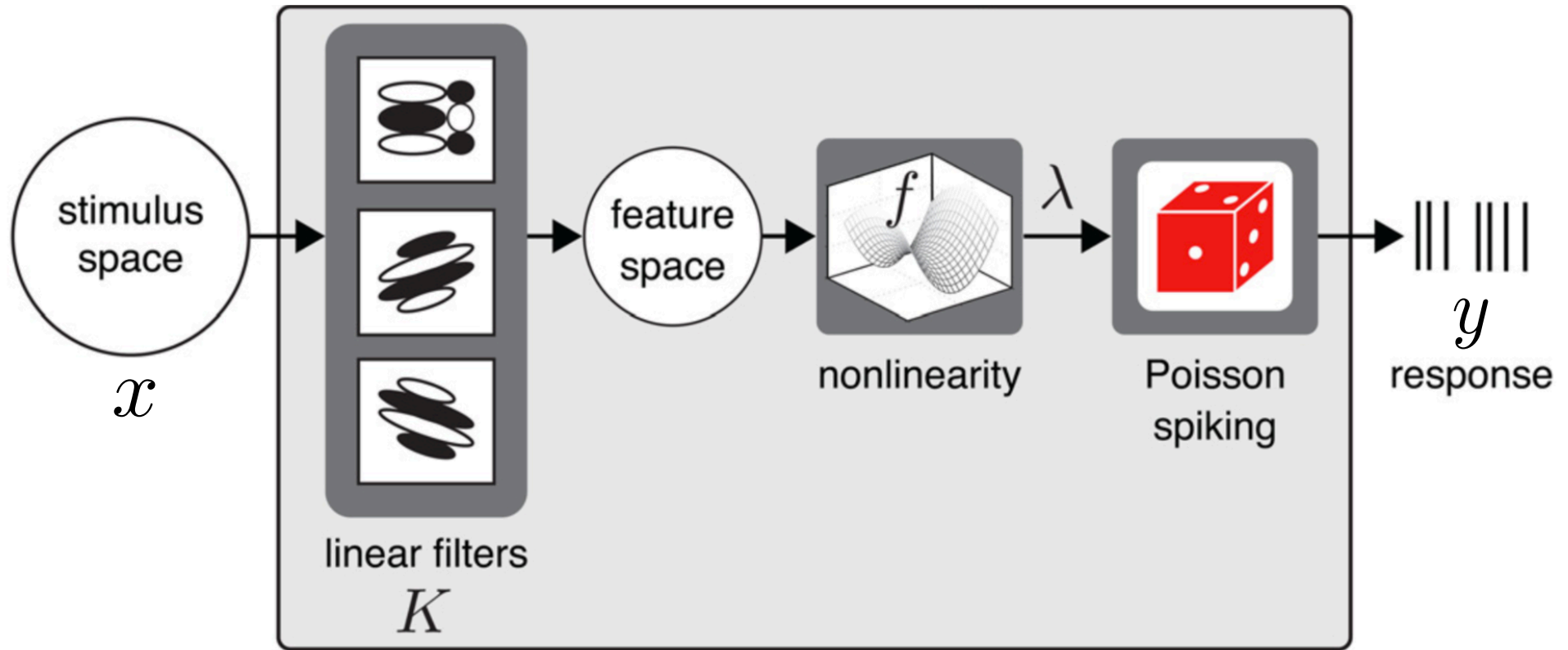
- from “systems identification” literature (1960s-70s)
- white noise stimuli
- estimate kernels using moments of spike-triggered stimuli

# Why are Volterra/Wiener models (generally) bad?

- no output nonlinearity
- polynomials give poor fit to neural nonlinearities (e.g., rectifying, saturating)

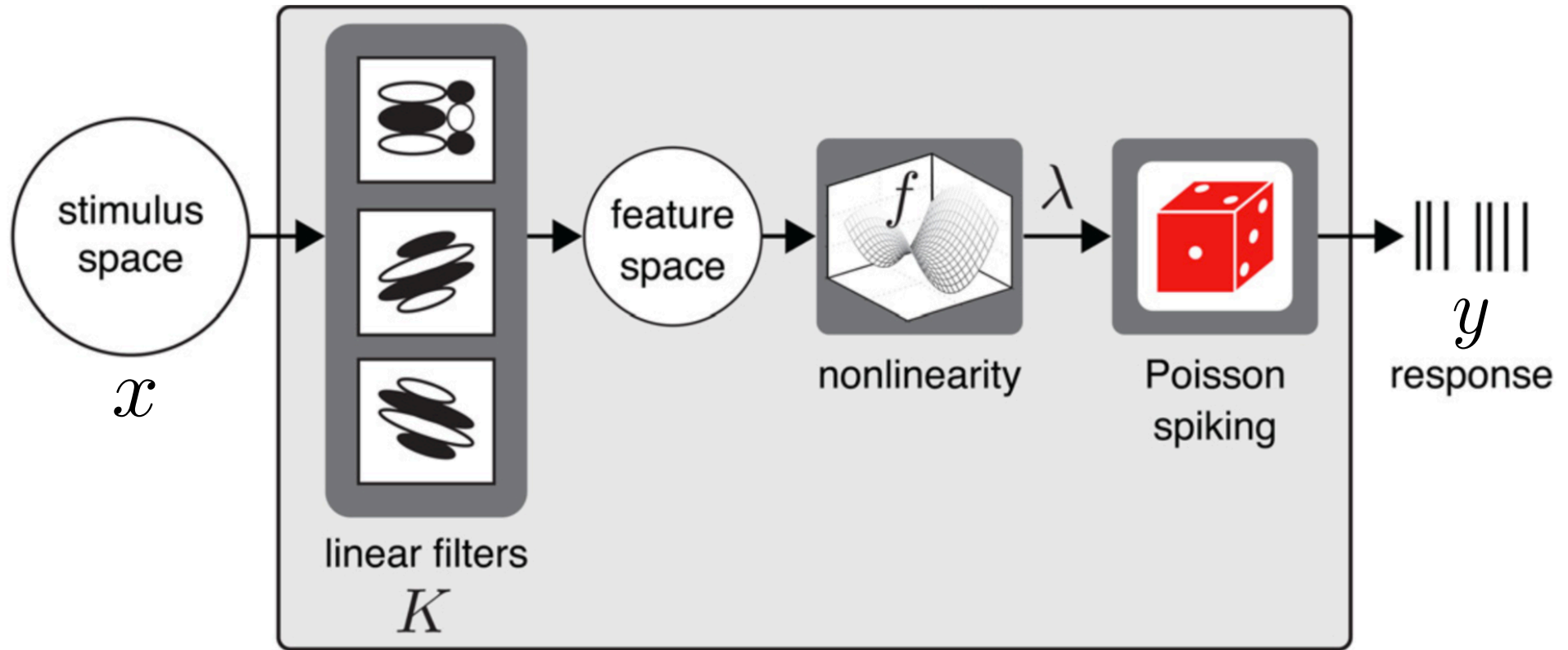


# multi-filter LNP



- responses may depend on more than one projection of stimulus!
- emphasis on *dimensionality reduction*
- no longer technically a GLM if fitting nonlinearity  $f$

# multi-filter LNP

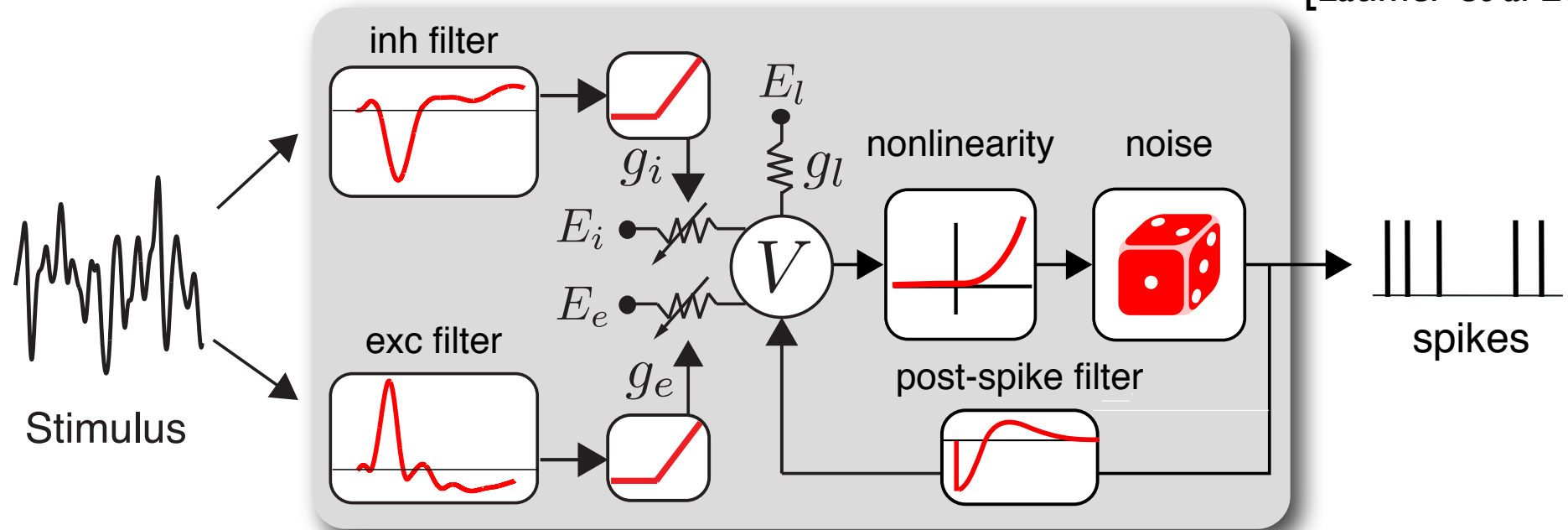


## Estimators:

- Spike-triggered covariance (STC) [de Ruyter & Bialek 1998, Schwartz et al 2006]
- Generalized Quadratic Model (GQM) [Park & Pillow 2011; Park et al 2013; Rajan et al 2013]
- maximally informative dimensions (MID) / maximum likelihood  
[Sharpee et al 2004] [Williamson et al 2015]

# extending GLM to conductance-based model

[Latimer et al 2014]



conductances

$$g_e(t) = f_c(k_e \cdot \mathbf{x}(t))$$

$$g_i(t) = f_c(k_i \cdot \mathbf{x}(t))$$

membrane  
dynamics

$$\frac{dV}{dt} = g_l(E_l - V) + g_e(E_e - V) + g_i(E_i - V)$$

inst. spike rate

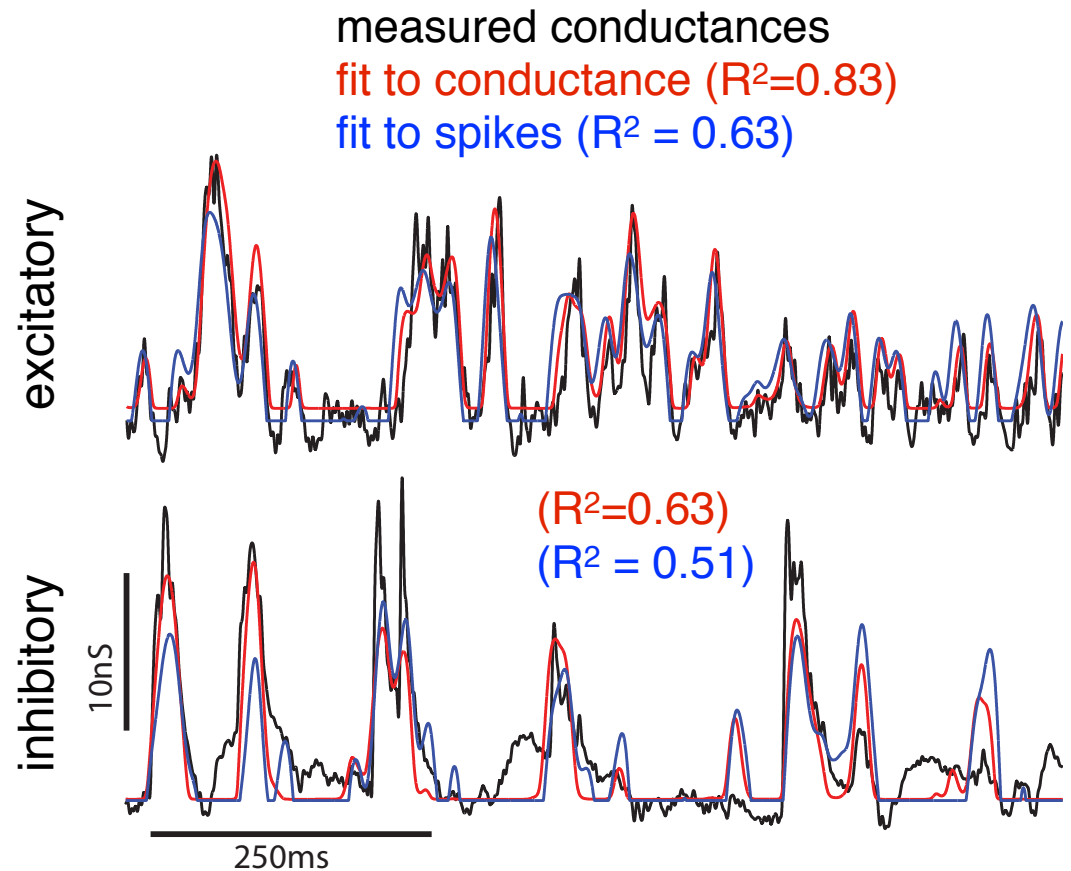
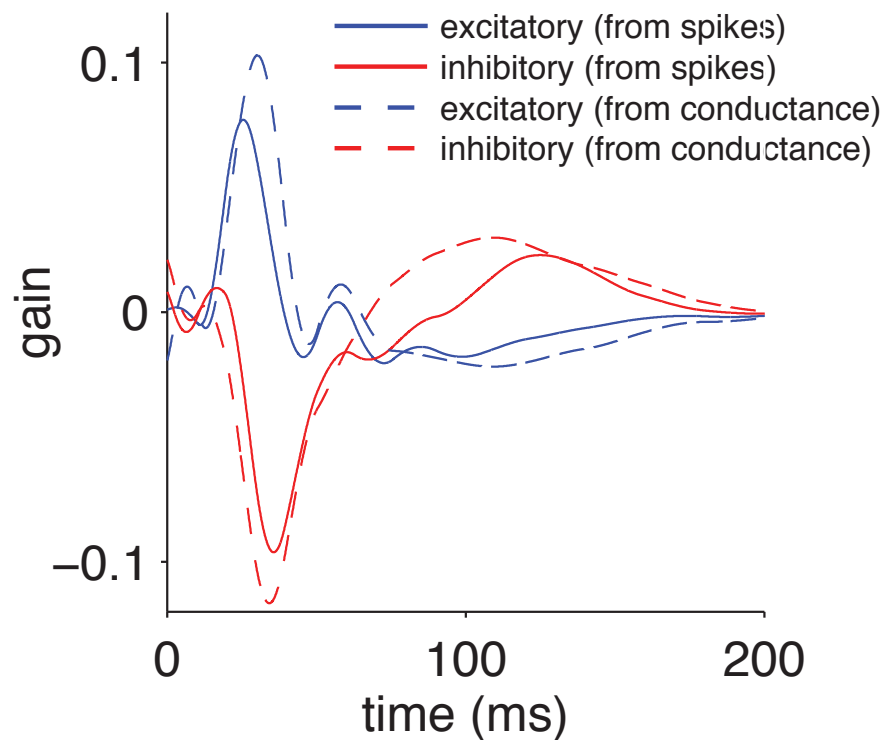
$$\lambda(t) = f(V(t))$$

- shunting inhibition
- adaptive changes in dynamics

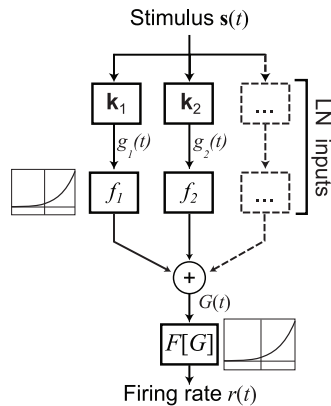
# extending GLM to conductance-based model

- intracellular recordings in macaque parasol RGCs (Fred Rieke)

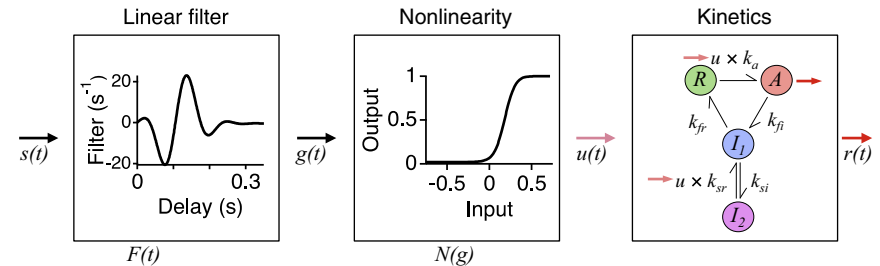
Linear filters



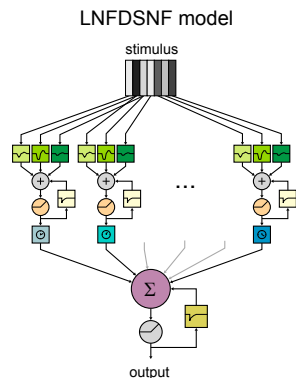
# many other biophysically oriented extensions



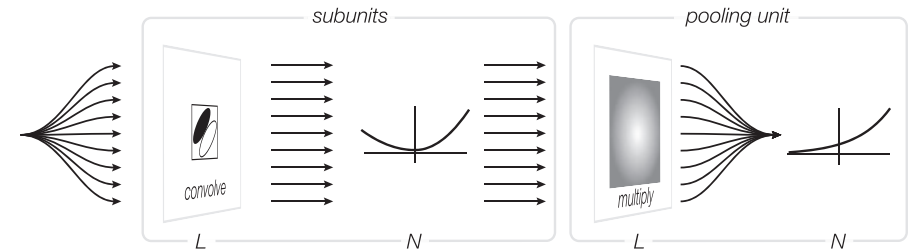
Nonlinear input model (NIM)  
[McFarland, Cui, & Butts 2013]



Linear-Nonlinear-Kinetics (LNK)  
[Ozuysal & Baccus 2014]

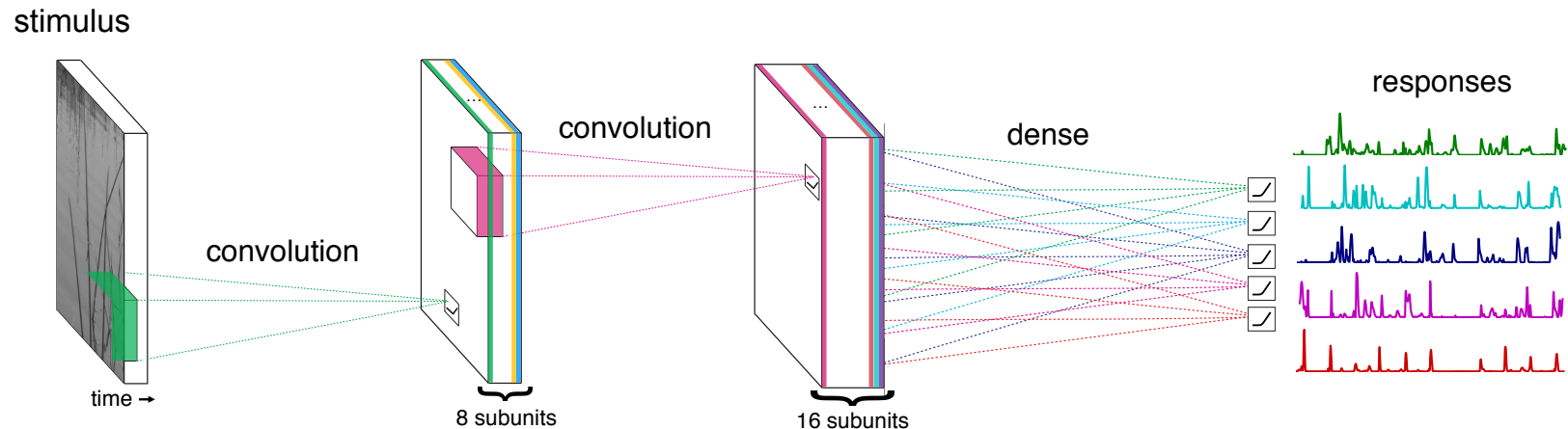


LNFDSNF model  
linear-nonlinear-feedback-delayed-sum-nonlinear-feedback  
[Real, Asari, Gollisch & Meister 2017]



convolutional subunit model  
[Vintch, Movshon & Simoncelli 2015,  
Wu, Park & Pillow 2014]

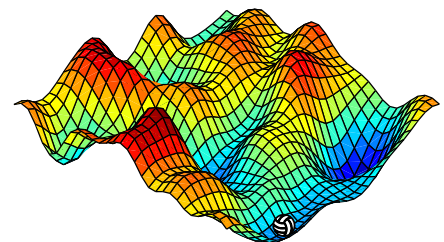
# deep learning / deep neural networks (DNNs)



[Yamins et al 2014, McIntosh et al 2016, Maheswaranathan et al 2017, Benjamin et al 2017, ...]

If you understand GLMs... you understand DNNs!

- stack many LNs on top of each other: LN LN LN LN P
- use gradient ascent to maximize likelihood
- use software (tensorflow, theano) to compute gradients (no more computing gradients by hand!)
- use a bunch of tricks (batches, noise, SGD, dropout, ....)
- do NOT worry about local maxima!



[credit: Jakob Macke]

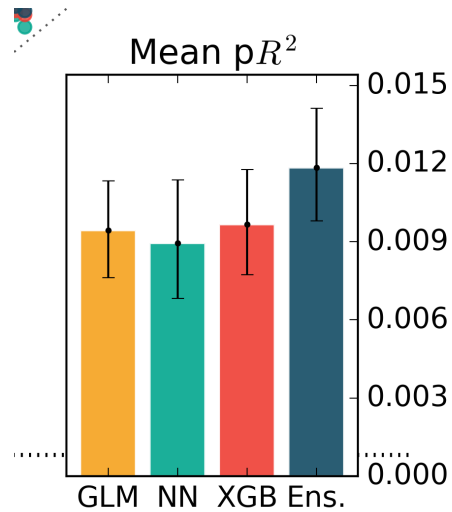


# Modern machine learning far outperforms GLMs at predicting spikes

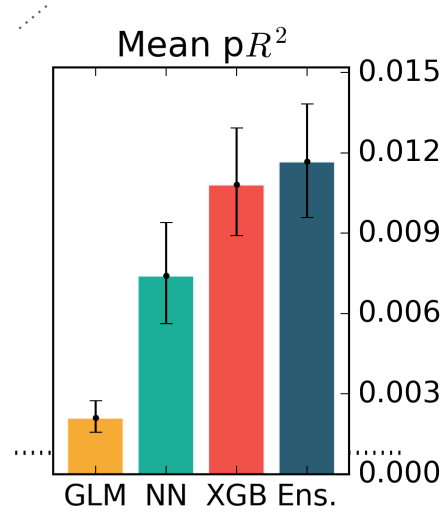
Ari S. Benjamin<sup>1</sup>, Hugo L. Fernandes<sup>2</sup>, Tucker Tomlinson<sup>3</sup>, Pavan Ramkumar<sup>2,4</sup>, Chris VerSteeg<sup>1</sup>, Lee Miller<sup>1,2,3</sup>, Konrad Paul Kording<sup>1,2,3</sup>

macaque M1

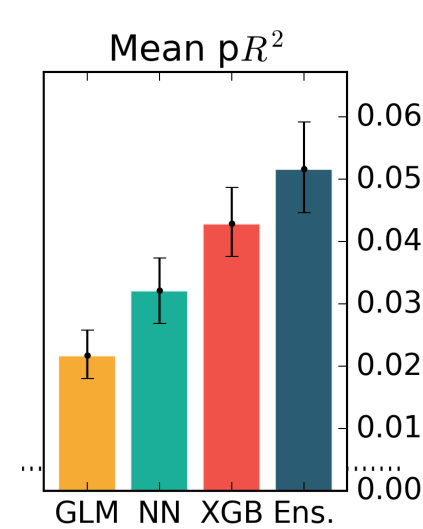
**Fig 2**



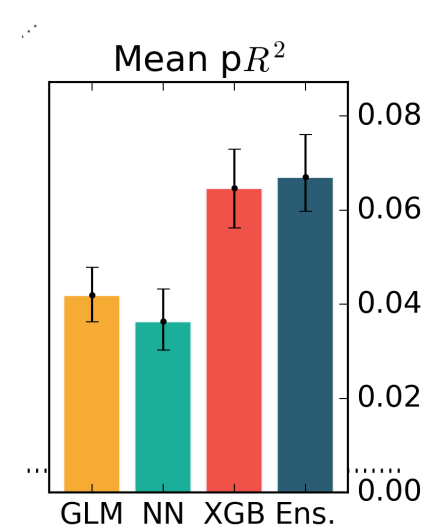
**Fig 3**



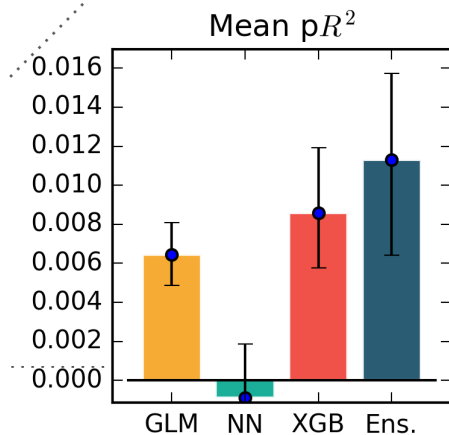
**Fig 4**



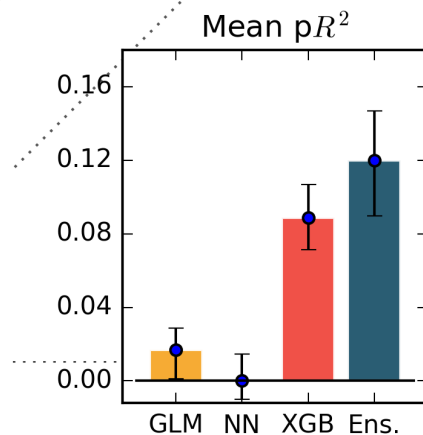
**Fig 5**



**Fig 6**



macaque S1



hippocampus

GLMs

NNs

~~Modern machine learning far~~ outperforms ~~GLMs~~ at predicting spikes ?

Ari S. Benjamin<sup>1</sup>, Hugo L. Fernandes<sup>2</sup>, Tucker Tomlinson<sup>3</sup>, Pavan Ramkumar<sup>2,4</sup>, Chris VerSteeg<sup>1</sup>, Lee Miller<sup>1,2,3</sup>, Konrad Paul Kording<sup>1,2,3</sup>

Fig 2

Mean  $pR^2$

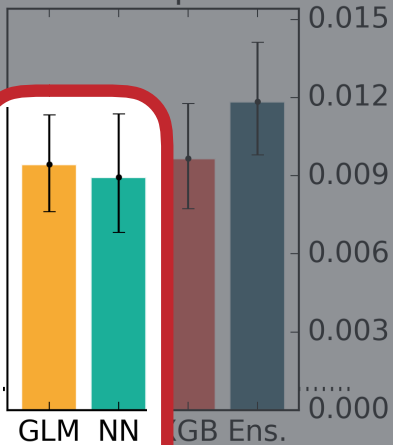


Fig 3

Mean  $pR^2$

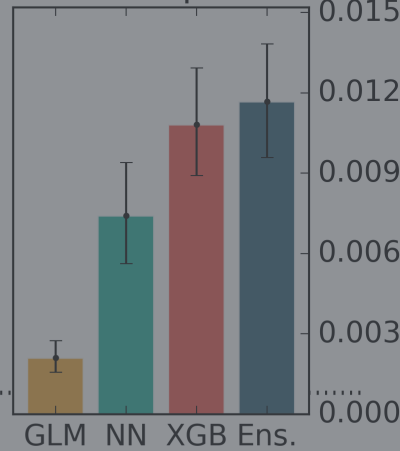


Fig 4

Mean  $pR^2$

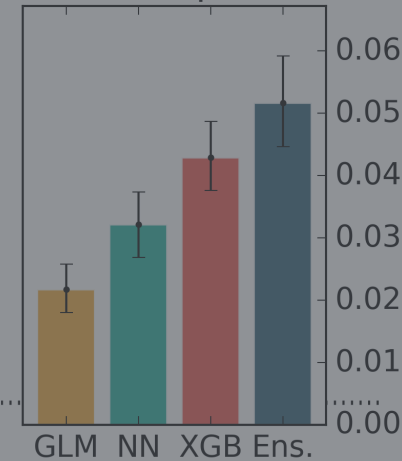


Fig 5

Mean  $pR^2$

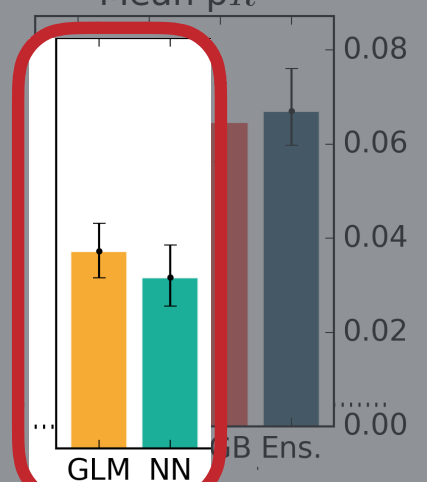
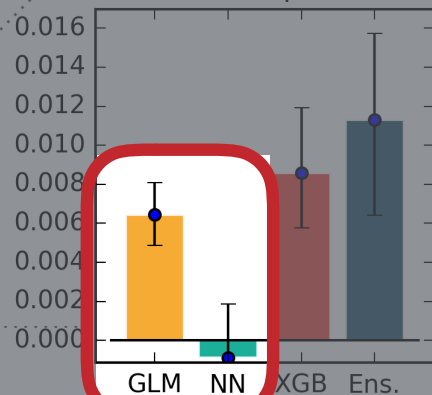


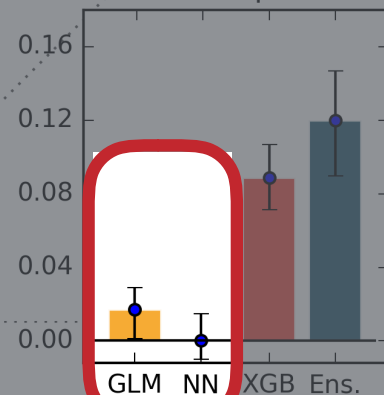
Fig 6

Mean  $pR^2$



macaque S1

Mean  $pR^2$

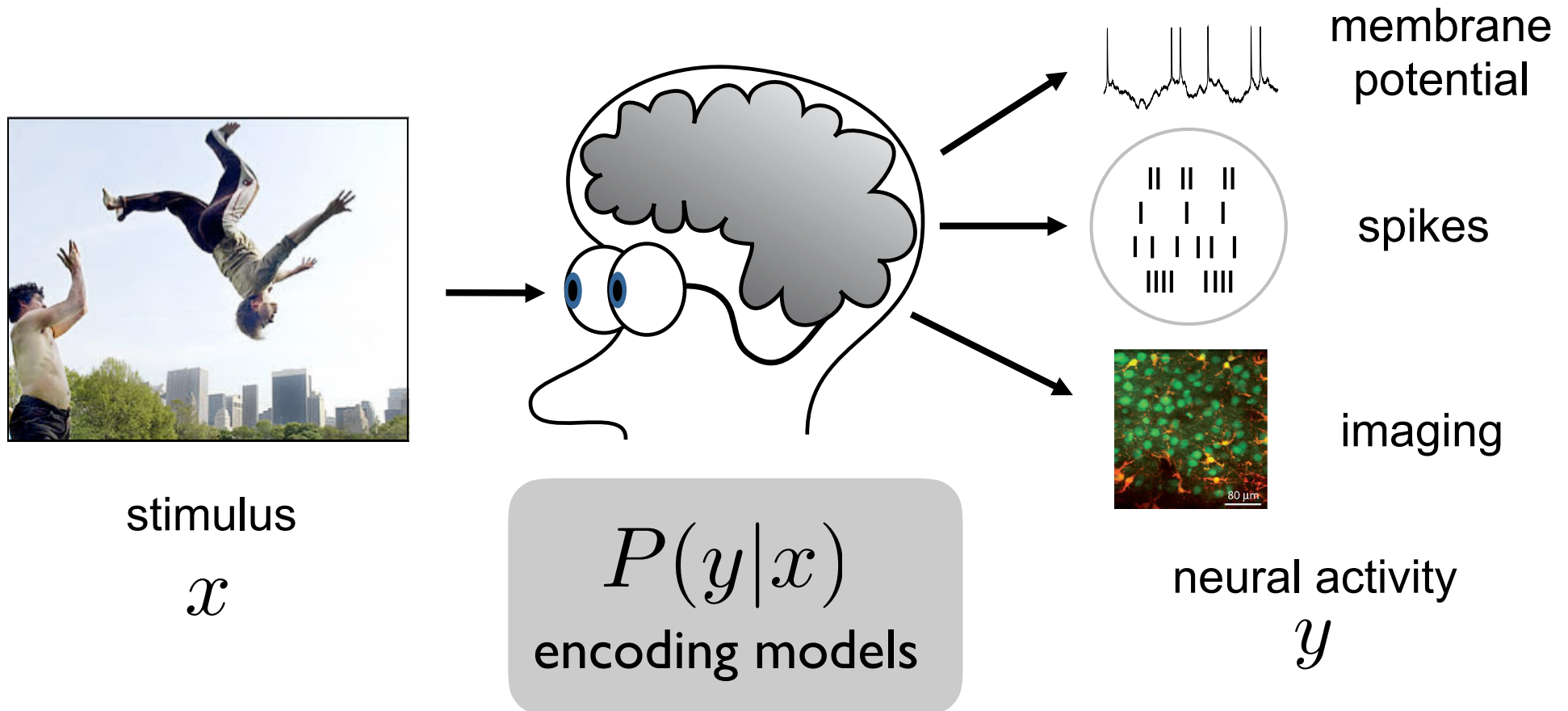


hippocampus

(No of course not!)

GLM is a special case of  
NN!

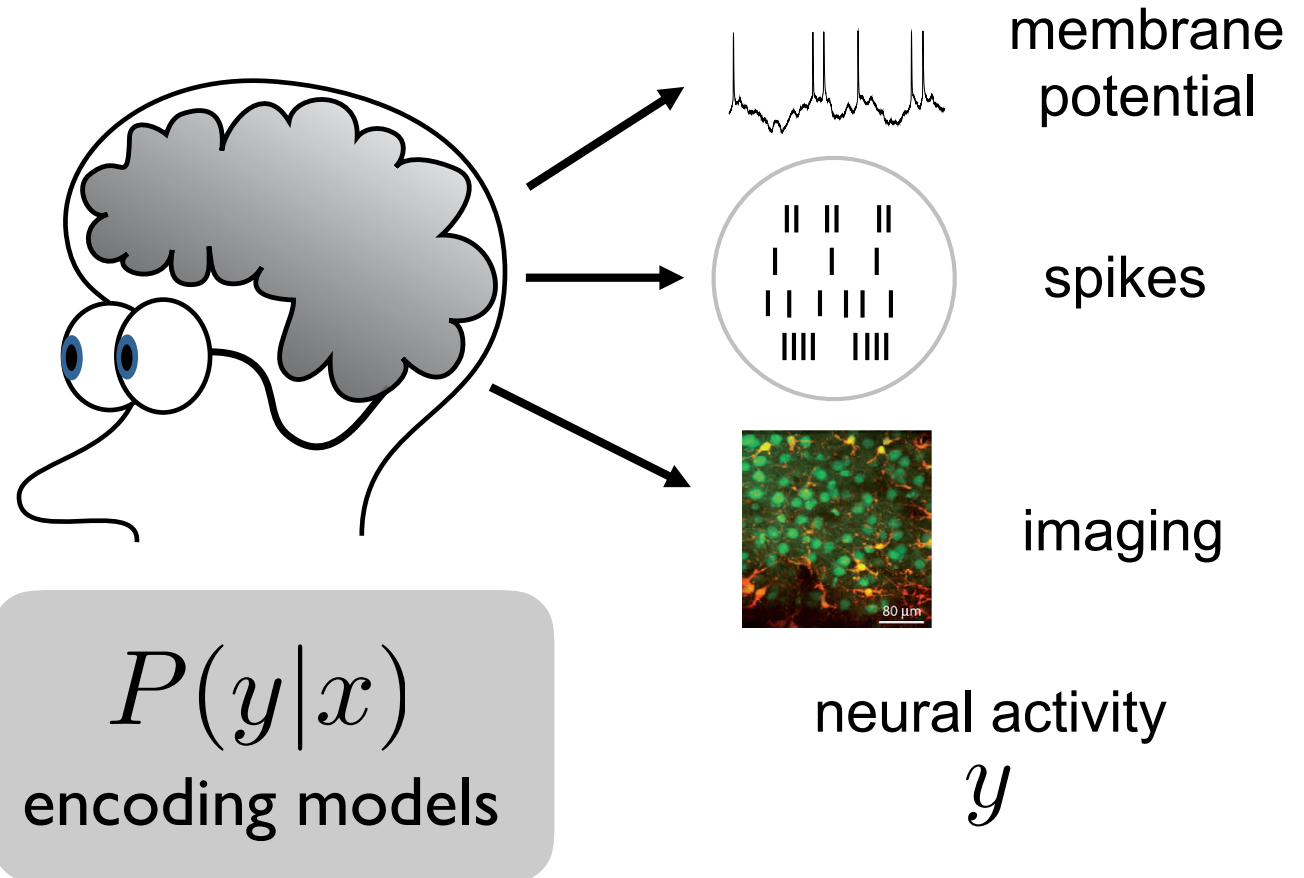
# encoding models



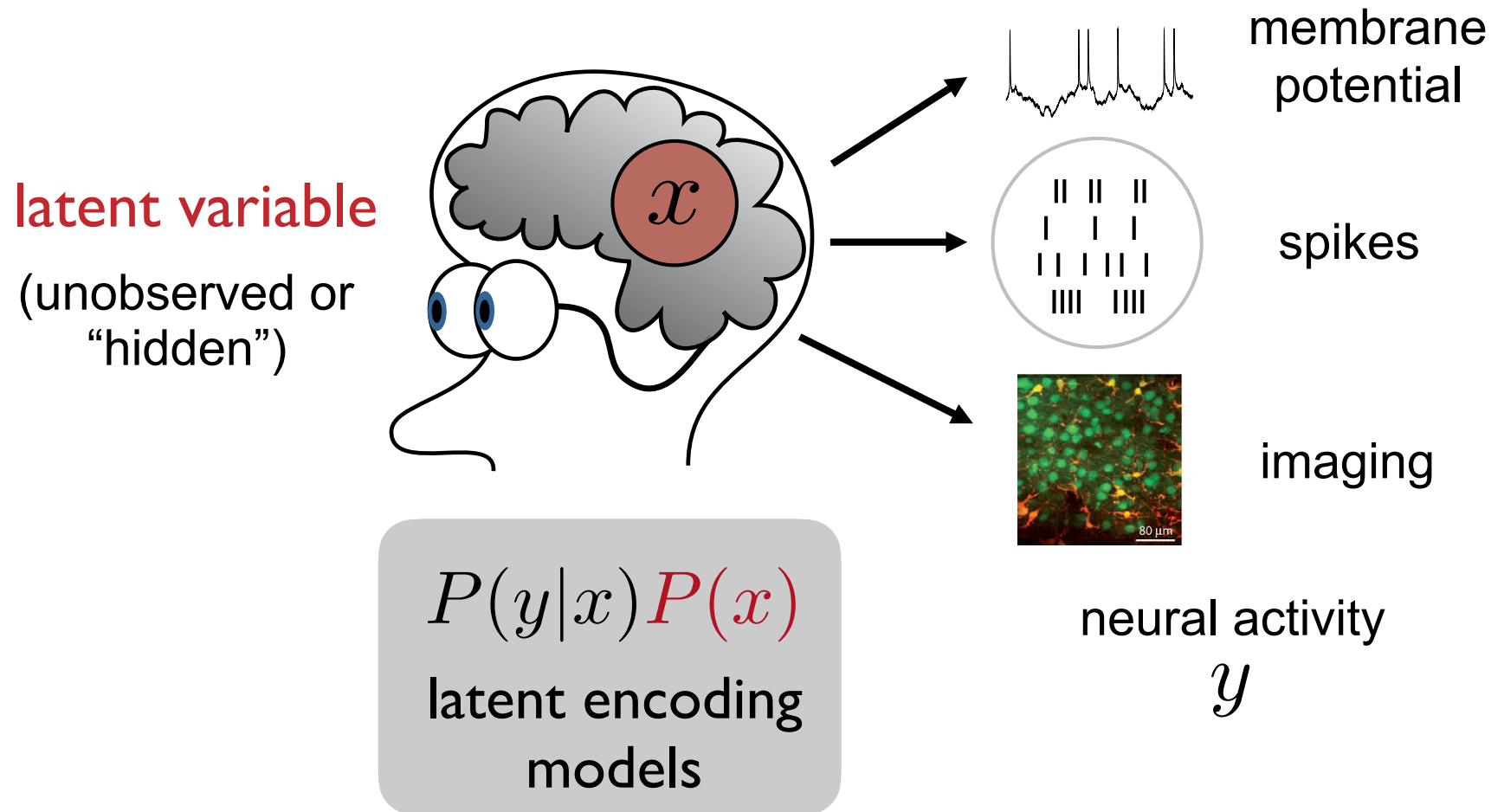
**What if there's  
no stimulus?**

# encoding models

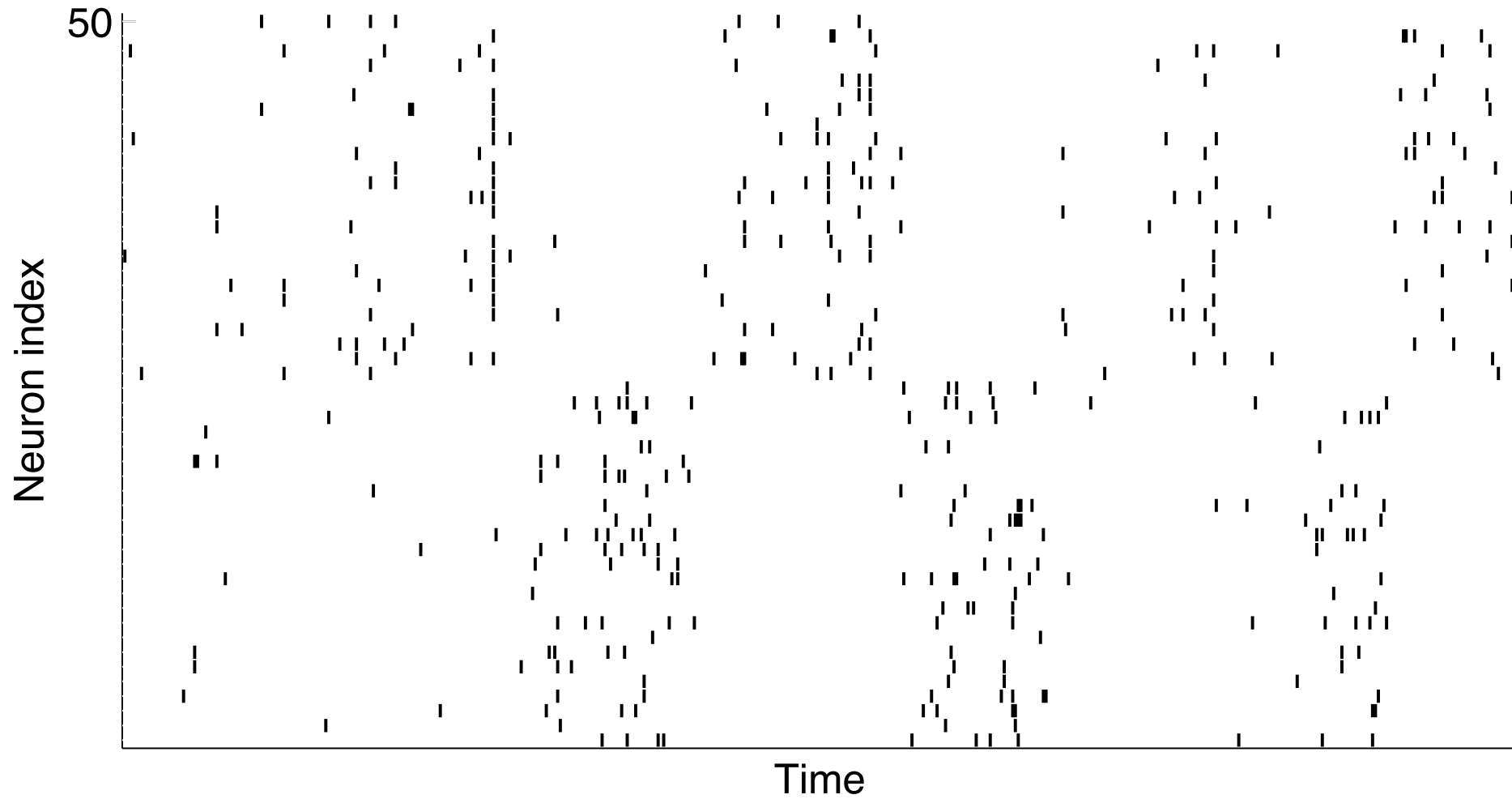
What if there's  
no stimulus?



# latent variable models

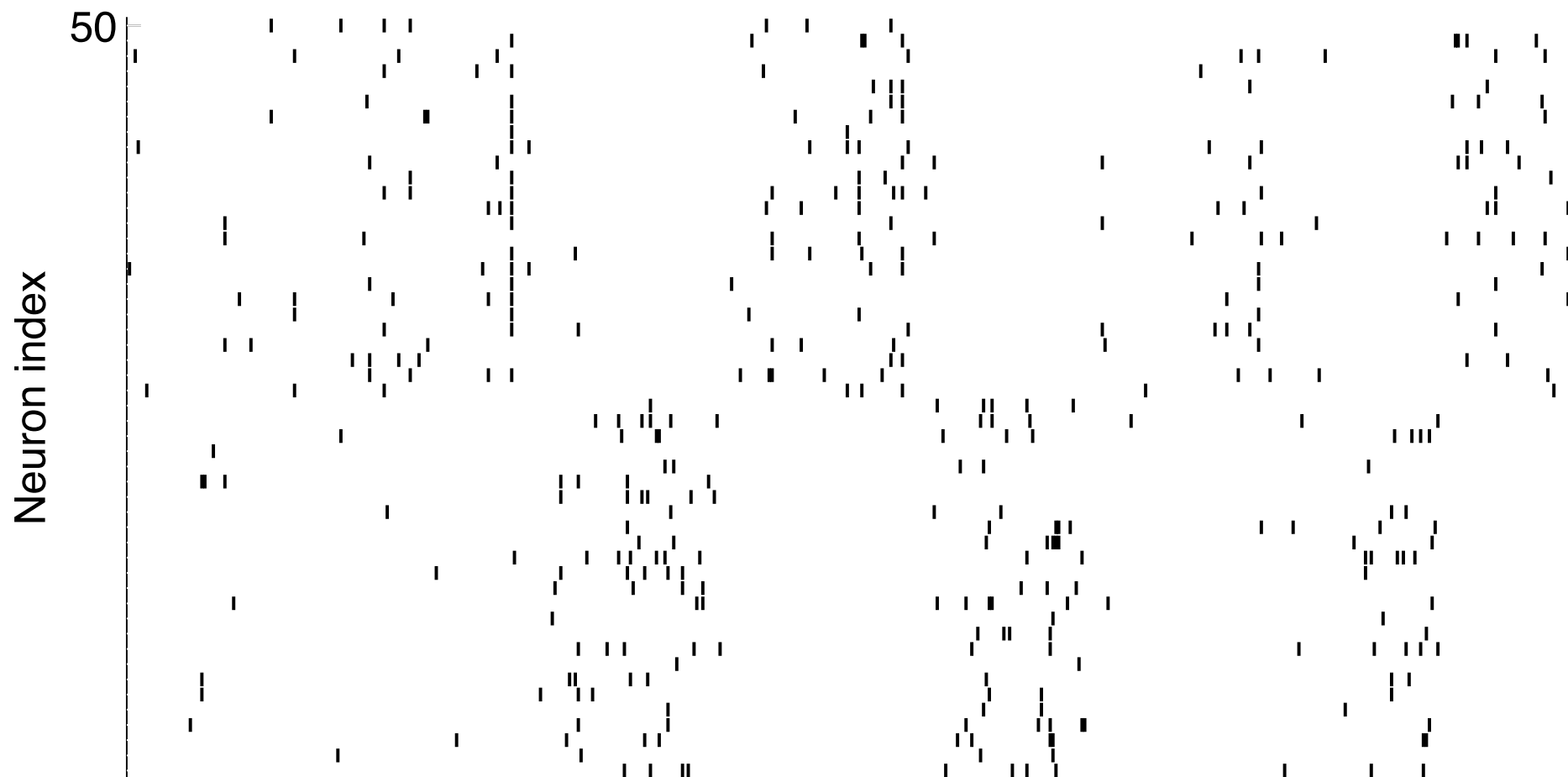


# spike responses

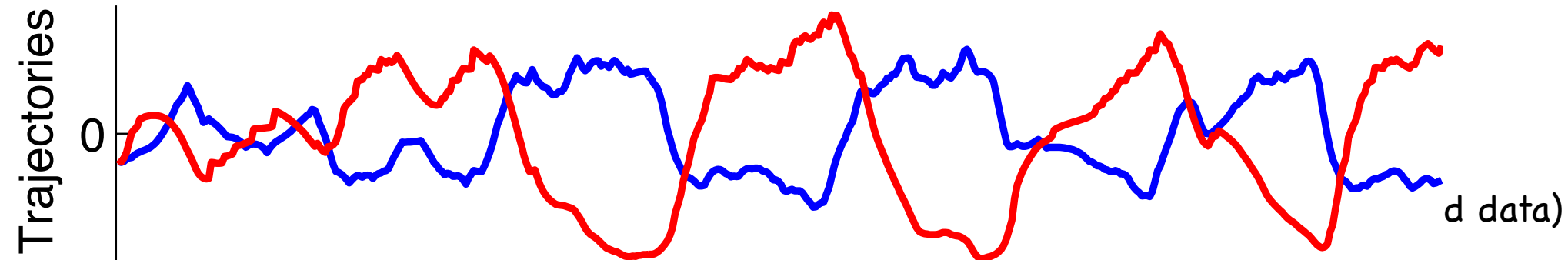


(simulated data)

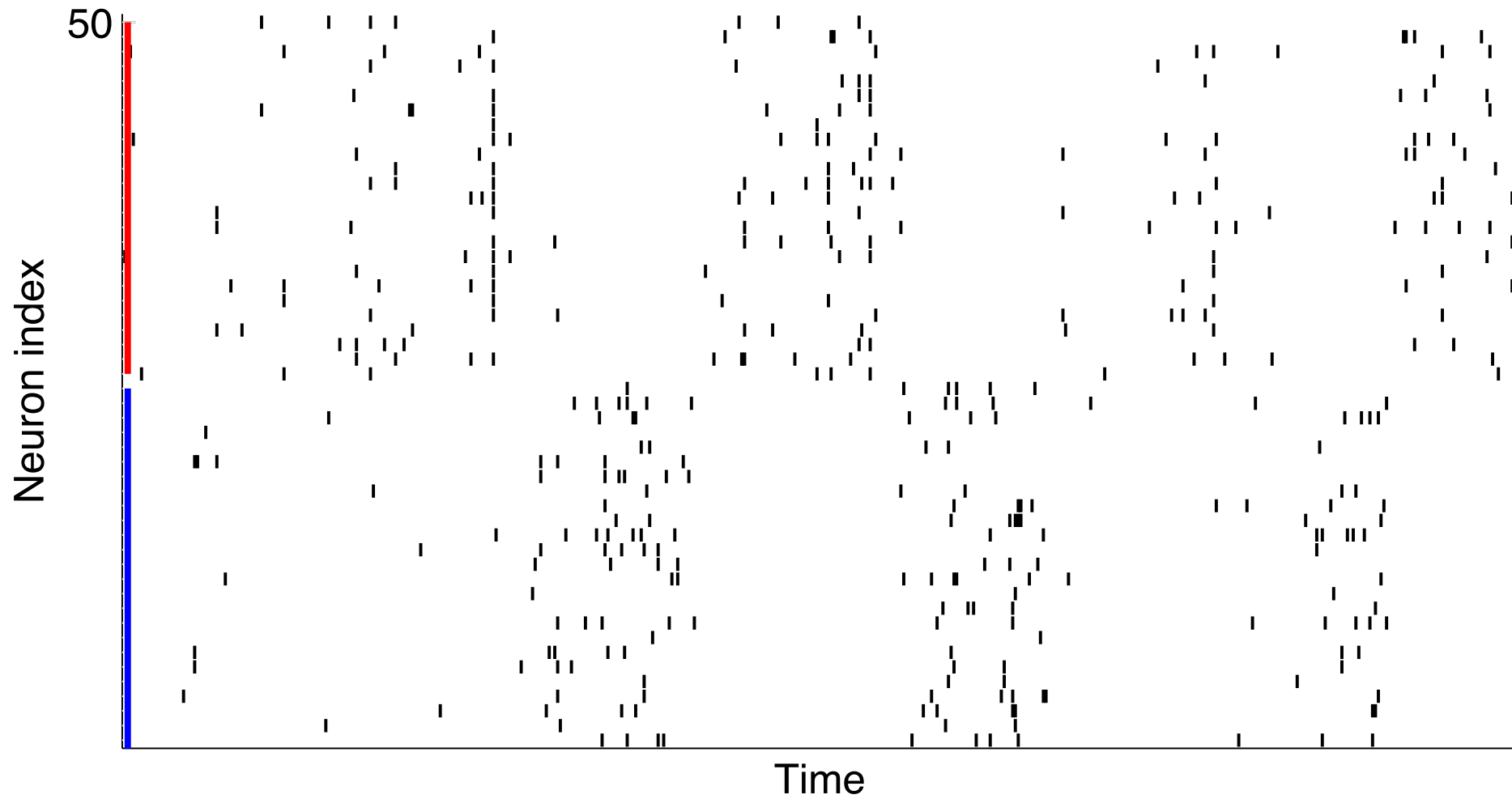
# spike responses



## inferred latent variables

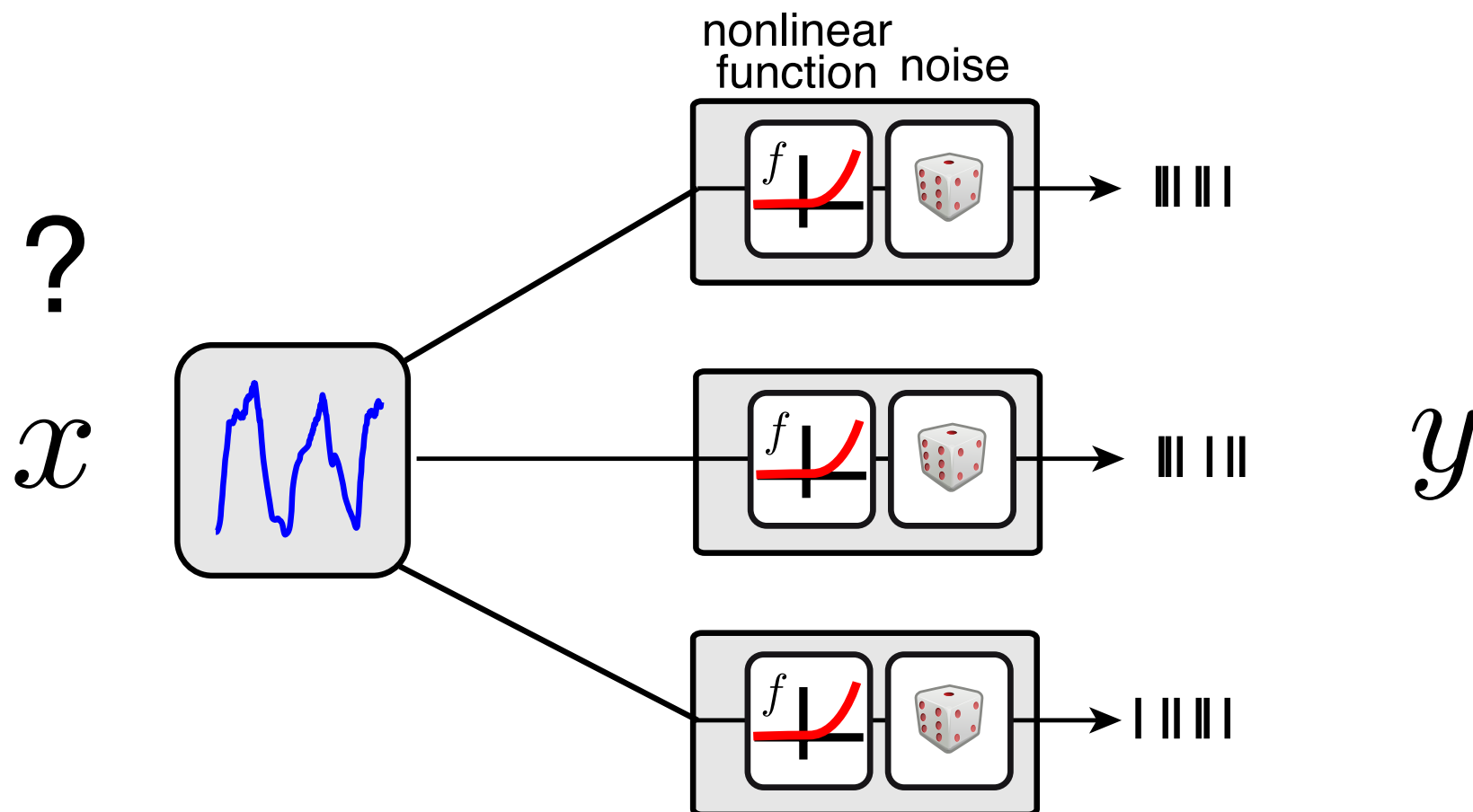


# spike responses





latent variable models = GLMs where we don't know  $x$



**goal:** find shared structure underlying  $y$

**chicken and egg problem**

# Why are latent variable models hard to work with?

- hard to compute likelihood!

# Why are latent variable models hard to work with?

- hard to compute likelihood!

**sensory  
encoding model**

$x$



$y$

Poisson GLM: **fit using:**

$$\log P(Y|X) = \sum_t y_t \log \lambda_t - \lambda_t$$

# Why are latent variable models hard to work with?

- hard to compute likelihood!

**sensory  
encoding model**

$x$



$y$

fit using:

$$\log P(Y|X) = \sum_t y_t \log \lambda_t - \lambda_t$$

**latent  
variable model**

$x$



$y$

fit using:  $\log P(Y)$

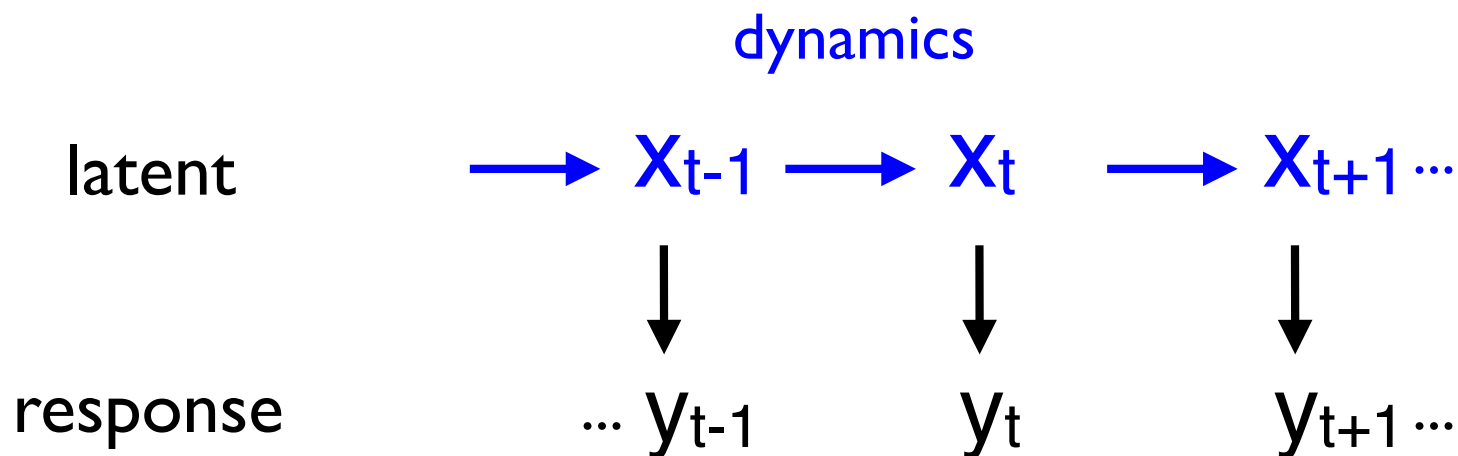
$$= \log \int P(Y|X)P(X)dx$$

requires an integral!

# Why are latent variable models hard to work with?

- hard to compute likelihood!

## latent dynamical model



fit using:  $\log P(Y)$

$$\int \prod_{t=1}^T \left( \overbrace{p(y_t|x_t)}^{\text{observations}} \overbrace{p(x_t|x_{t-1})}^{\text{dynamics}} \right) dx_1 dx_2 \cdots dx_T$$

high-dimensional integral

# Fitting Latent Variable Models

## I. Sampling (“MCMC”) - fully Bayesian inference

- procedure for sampling joint distribution:  $P(\theta, \{\mathbf{x}\} \mid \{\mathbf{r}\}, \{c\})$ 
  - 1) sample latents:  $\mathbf{x}^{(i)} \sim p(\mathbf{x} \mid \mathbf{r}, c, \theta^{(i)})$  conditional over latents
  - 2) sample parameters:  $\theta^{(i+1)} \sim p(\theta \mid \mathbf{r}, c, \mathbf{x}^{(i)})$  conditional over parameters

# Fitting Latent Variable Models

## I. Sampling (“MCMC”) - fully Bayesian inference

- procedure for sampling joint distribution:  $P(\theta, \{\mathbf{x}\} \mid \{\mathbf{r}\}, \{c\})$

1) sample latents:  $\mathbf{x}^{(i)} \sim p(\mathbf{x} \mid \mathbf{r}, c, \theta^{(i)})$  conditional over latents

2) sample parameters:  $\theta^{(i+1)} \sim p(\theta \mid \mathbf{r}, c, \mathbf{x}^{(i)})$  conditional over parameters

## 2. Expectation maximization (EM)

Alternate updating parameters and posterior over latents.

# Fitting Latent Variable Models

## I. Sampling (“MCMC”) - fully Bayesian inference

- procedure for sampling joint distribution:  $P(\theta, \{\mathbf{x}\} \mid \{\mathbf{r}\}, \{c\})$ 
  - 1) sample latents:  $\mathbf{x}^{(i)} \sim p(\mathbf{x} \mid \mathbf{r}, c, \theta^{(i)})$  conditional over latents
  - 2) sample parameters:  $\theta^{(i+1)} \sim p(\theta \mid \mathbf{r}, c, \mathbf{x}^{(i)})$  conditional over parameters

## 2. Expectation maximization (EM)

Alternate updating parameters and posterior over latents.

## 3. Variational inference

Optimize a lower bound on posterior over parameters

Easy with modern probabilistic programming languages

(STAN, Edward)



Latent Variable models are defined by two quantities:

latent $P(x)$	mapping $P(y x)$	→ Model
• discrete	• Gaussian	• Mixture of Gaussians (“clustering”)
• Gaussian	• linear, Gaussian	• Factor analysis (PCA is special case)
• linear Gaussian dynamics	• linear, Gaussian	• Linear Dynamical Systems (LDS) (“Kalman filter”)
• discrete transitions	• any	• Hidden Markov Model (“HMM”)

# variational latent Gaussian process (vLGP)

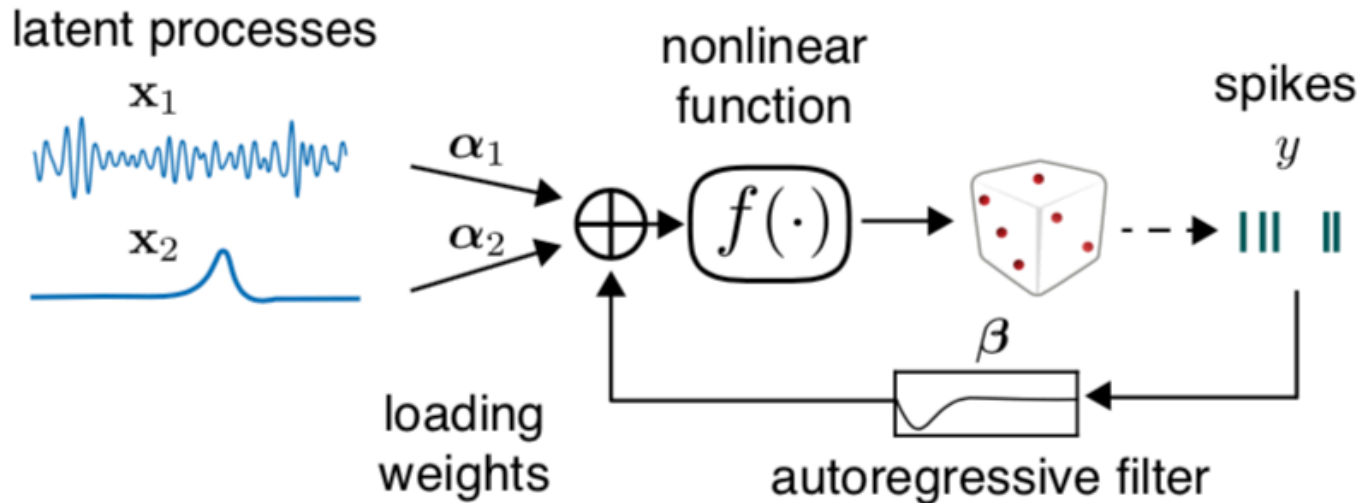
[Zhao & Park 2016]

latent  
 $P(x)$

mapping  
 $P(y|x)$

Gaussian Process

Poisson GLM



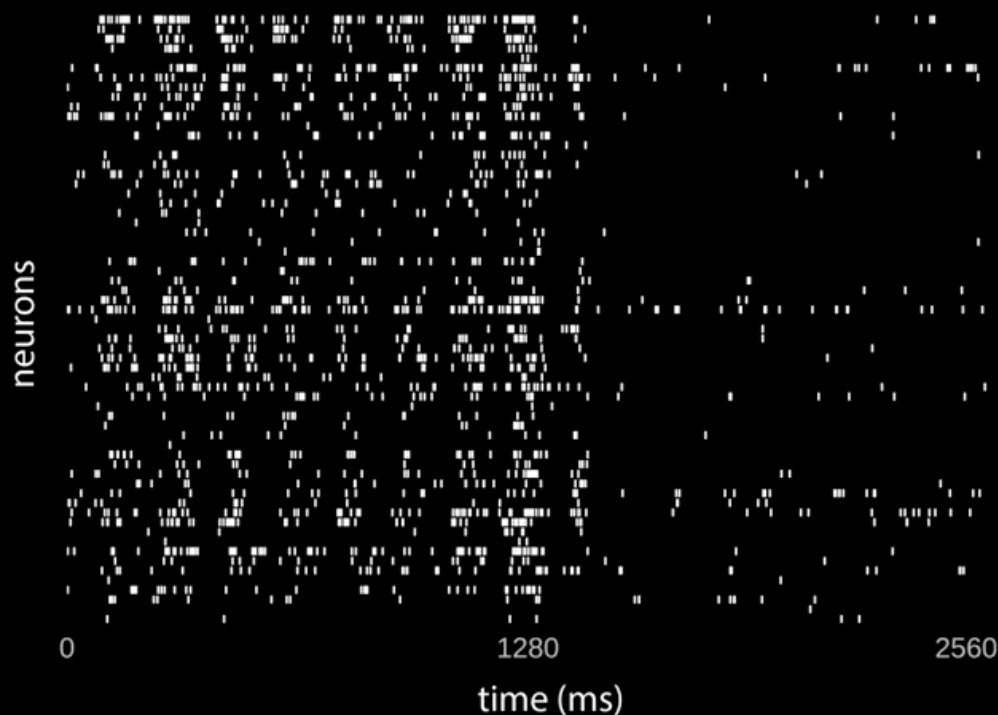
## variational latent Gaussian process (vLGP)

- 63 simultaneously-recorded V1 neurons [Graf et al 2011]
- stimuli: drifting sinusoidal gratings



# Latent dynamics in V1 driven by drifting grating

**arXiv: 1604.03053**



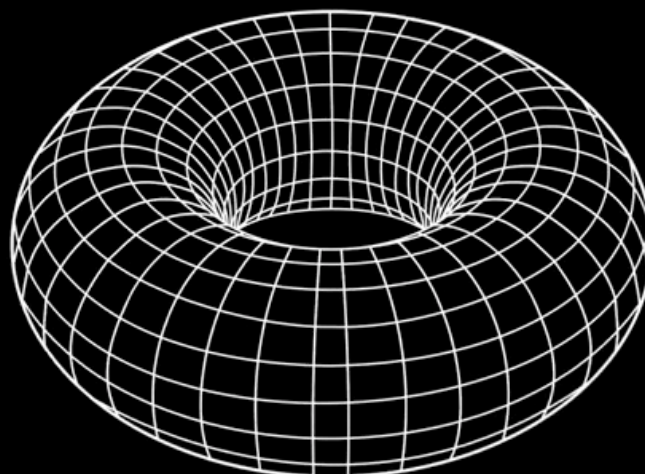
**Yuan Zhao & I. Memming Park**  
**Stony Brook University**

## variational latent Gaussian process (vLGP)

- 63 simultaneously-recorded V1 neurons [Graf et al 2011]
- stimuli: drifting sinusoidal gratings



**Latent trajectories  
for 36 directions  
form a torus**



# Summary

- descriptive statistical “encoding” models
- seek to capture structure in data
- formal tools for comparing models
- encoding and decoding analyses via Bayes rule
- models are modular, easy to build / extend / generalize

# Big Picture

- large-scale recording technology advancing rapidly
- lots of interesting structure in high-D neural data
- big opportunities in computational / statistical for developing new methods and models to find / exploit this structure!