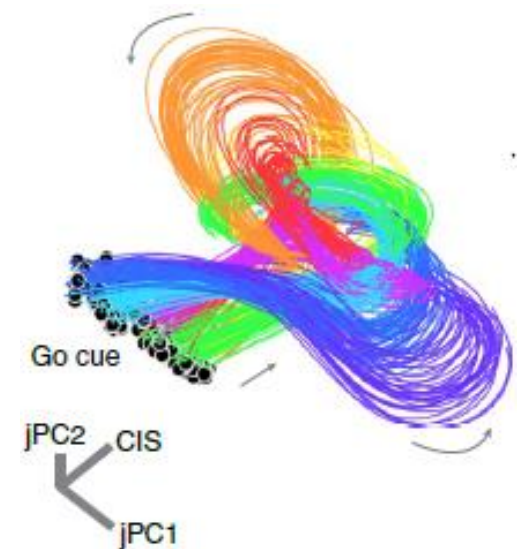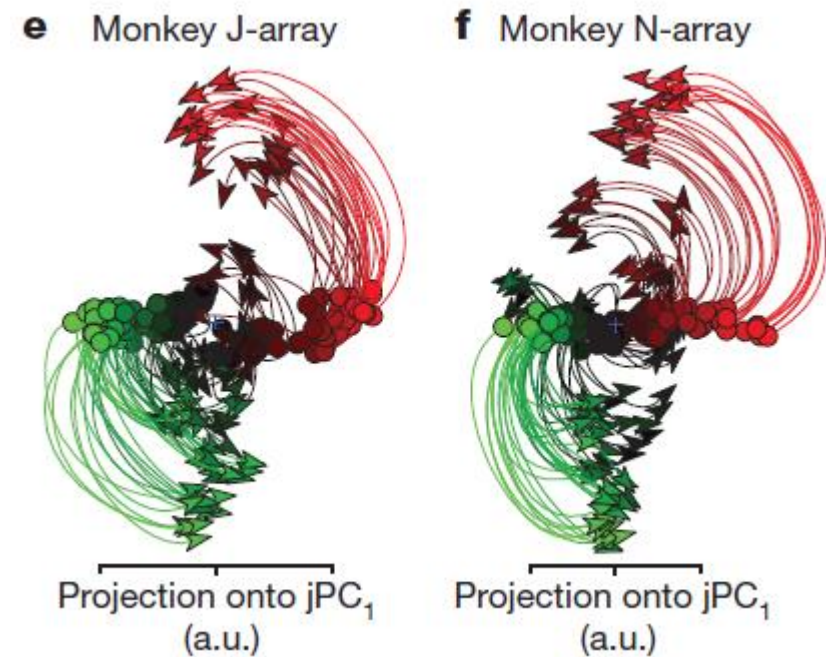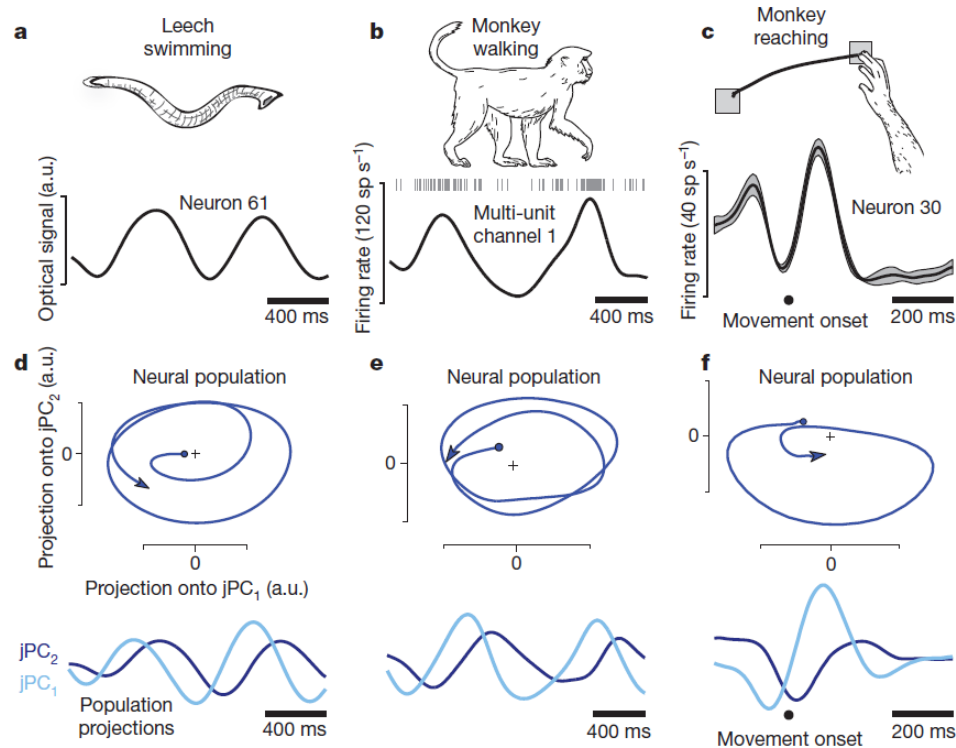# Inferring single-trial neural population dynamics using sequential auto-encoders

Pandarinath et al., 2018, *Nature methods*

*Nhat Le, NeuroComp meeting, Dec 11, 2018*

# Motivation

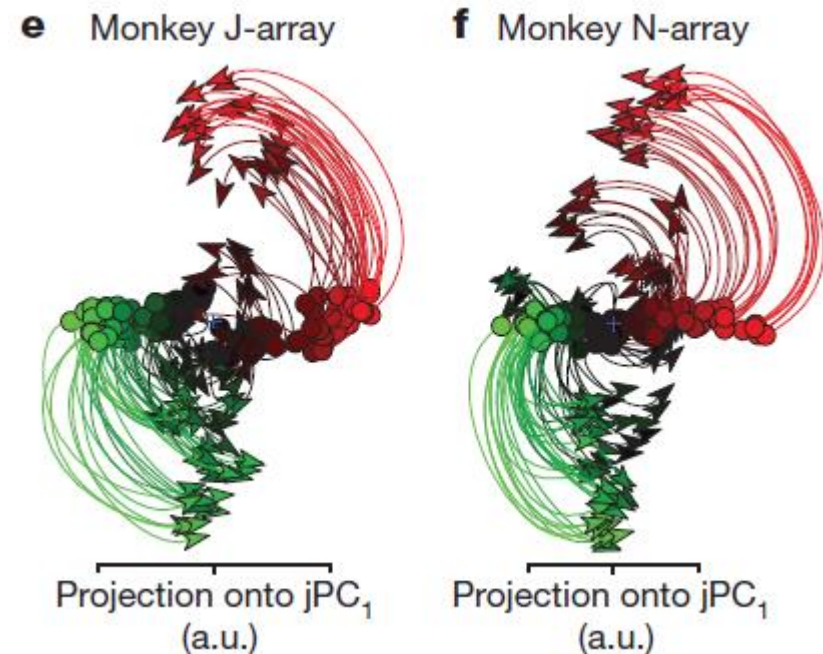- Dynamical systems perspective: dynamical systems underlie the pattern of neural populations

Churchland et al., 2012, *Nature*

# Motivation

- Dynamical systems perspective: dynamical systems underlie the pattern of neural populations

**Problem**: trajectories often computed based on trial averages

→ Can we uncover underlying dynamics of *single trials*?



e  Monkey J-array   f  Monkey N-array

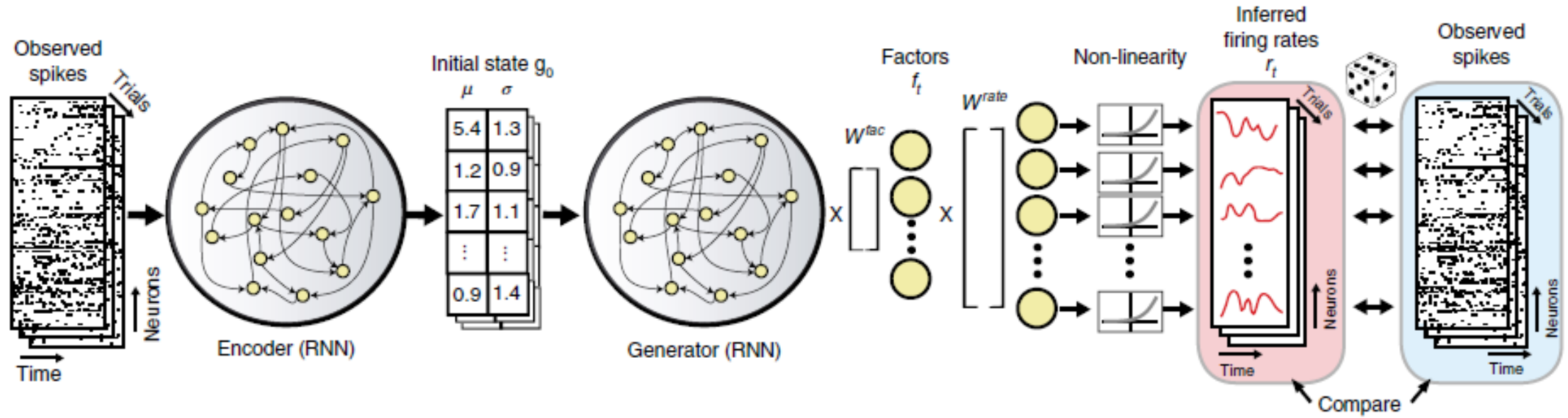Projection onto jPC$_1$ (a.u.)   Projection onto jPC$_1$ (a.u.)

# Part I: The LFADS basic architecture

1. LFADS assumes an underlying ***dynamical system*** (***recurrent neural network***) that generates spike trains

2. LFADS as a form of ***factor analysis***

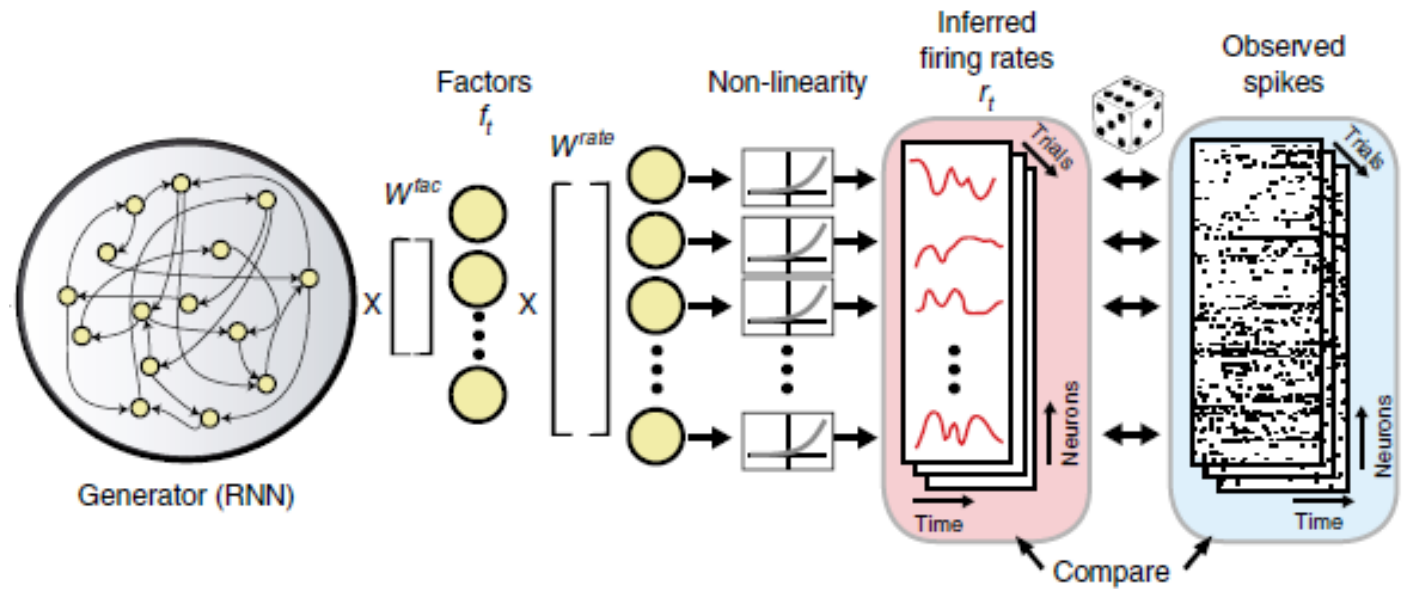3. LFADS as a ***variational autoencoder***

# Part I: The LFADS basic architecture

1. **LFADS assumes an underlying *dynamical system* (*recurrent neural network*) that generates spike trains**

2. LFADS as a form of *factor analysis*

3. LFADS as a *variational autoencoder*

# LFADS basic architecture

# LFADS basic architecture



Underlying
**dynamical system**
$g_{t'}$
which evolves
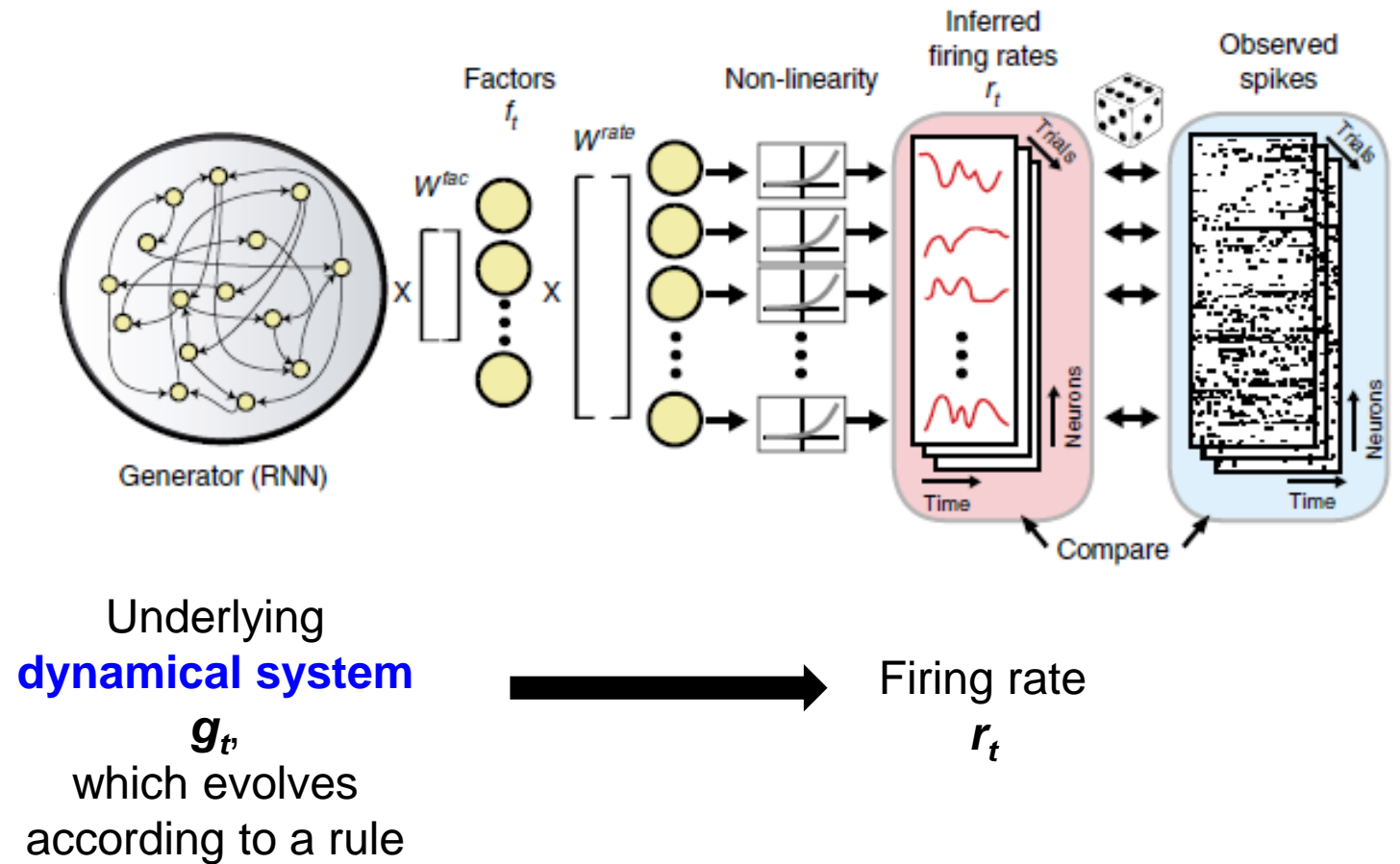according to a rule

Firing rate
$r_t$

$$\dot{g}(t) = F(g(t), u(t))$$

# LFADS basic architecture

Observed spikes depend on
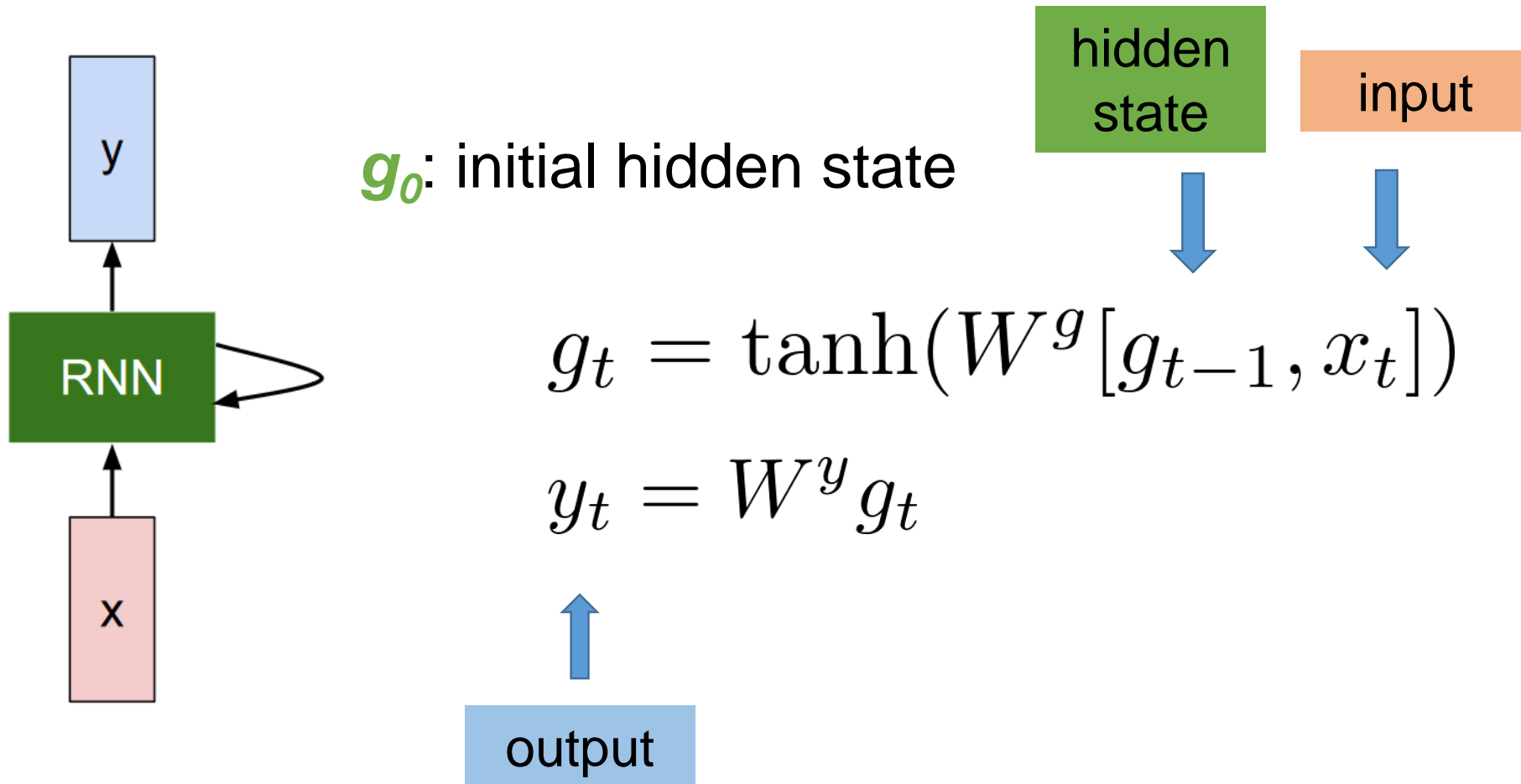
(1) Underlying dynamics ($F$)

(2) Initial conditions ($g_0$)

(3) Inputs from other brain areas ($u$)

(4) Spiking variability



Underlying **dynamical system** $g_t$, which evolves according to a rule

$\longrightarrow$ Firing rate $r_t$

$$\dot{g}(t) = F(g(t), u(t))$$

# To model underlying dynamics:
Recurrent neural networks (simple)



$g_0$: initial hidden state

$$g_t = \tanh(W^g[g_{t-1}, x_t])$$

$$y_t = W^y g_t$$

# RNN variant: Gated Recurrent Unit (GRU)

**Simple**                                    **GRU**
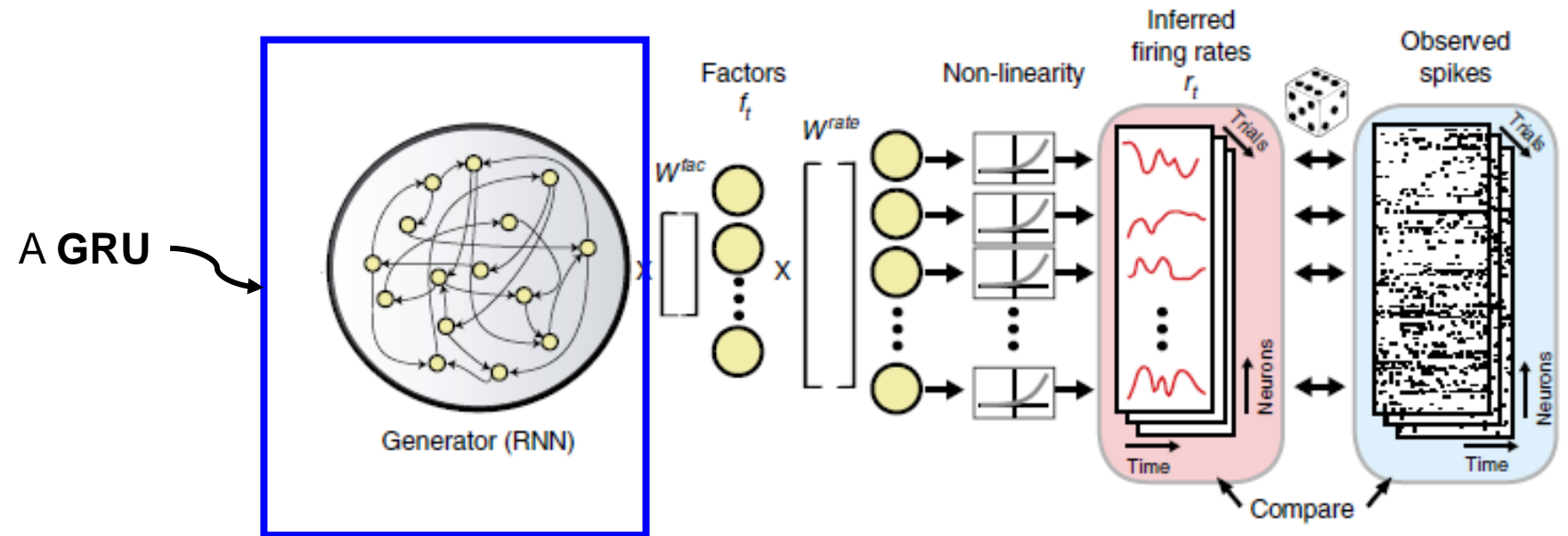
$$g_t = \tanh(W^g[g_{t-1}, x_t]) \quad c_t = \tanh(W^c[x_t, r_t \odot g_{t-1}])$$

$$g_t = u_t \odot \boxed{g_{t-1}} + (1 - u_t) \odot \boxed{c_t}$$

*- A combination of previous state '**memory**' and **updated state***

*- Prevents vanishing gradients*

# LFADS basic architecture

A **GRU**



Underlying
**dynamical system**
$g_{t'}$
which evolves
according to a rule

Firing rate
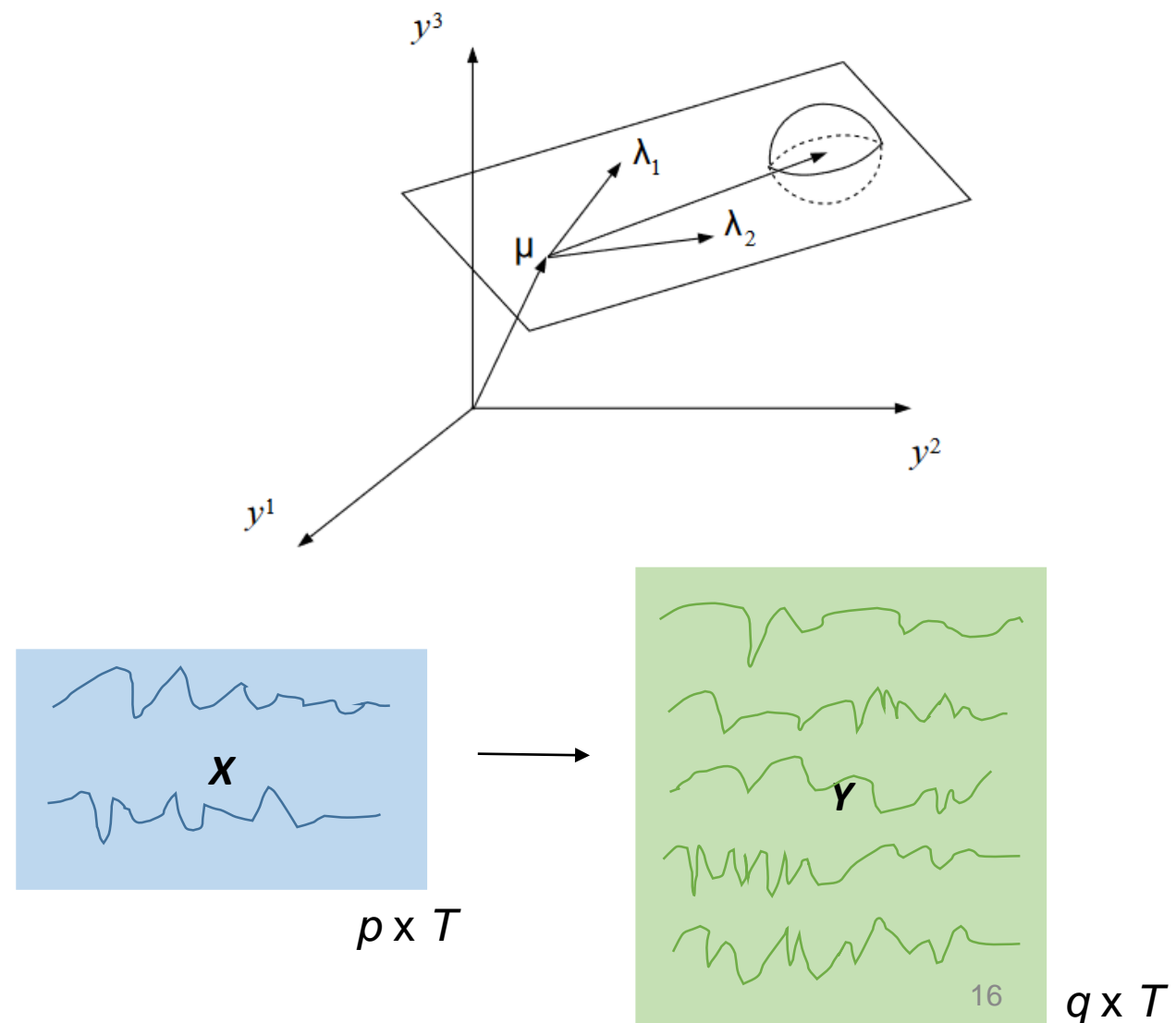$r_t$

$$\dot{g}(t) = F(g(t), u(t))$$

# Part I: The LFADS basic architecture

1. LFADS assumes an underlying *dynamical system* (*recurrent neural network*) that generates spike trains

2. LFADS as a form of *factor analysis*

3. LFADS as a *variational autoencoder*

# Factor Analysis

Assumes data $y$ (dimension $q$) is generated by a latent variable $x$ (dimension $p$) where $p < q$

$$\mathbf{y}_{:,t} \mid \mathbf{x}_{:,t} \sim \mathcal{N}\left(C\mathbf{x}_{:,t} + \mathbf{d},\ R\right),$$



X

p x T

Y

q x T

Roweis, 2004. Factor analysis and PCA

# Previous methods to uncover latent factors
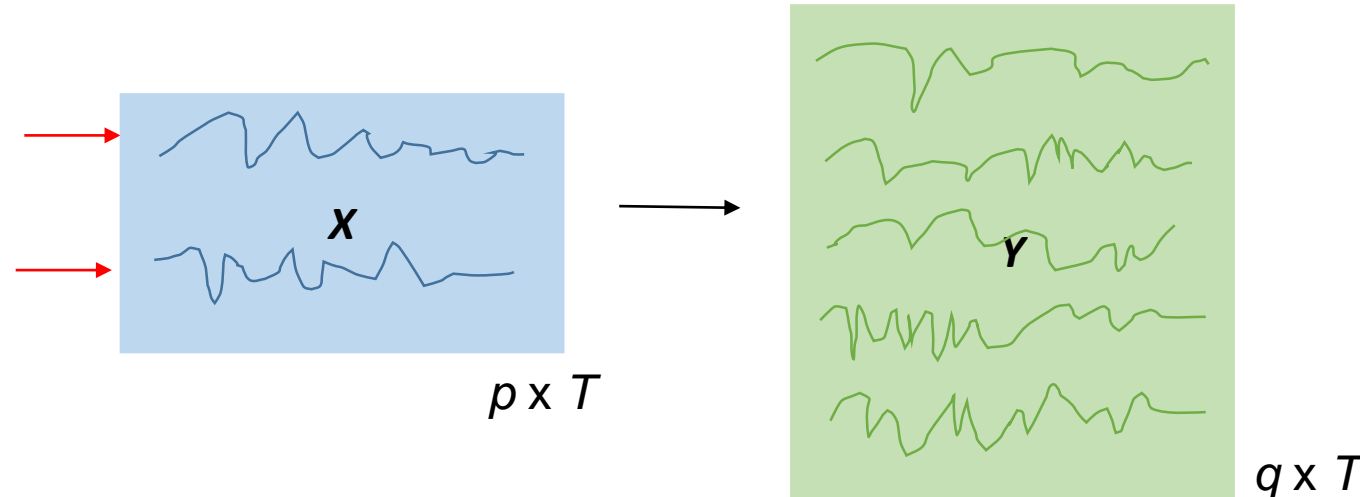
$$\mathbf{y}_{:,t} \mid \mathbf{x}_{:,t} \sim \mathcal{N}\left(C\mathbf{x}_{:,t} + \mathbf{d},\ R\right),$$

- *Gaussian process factor analysis*: each dimension of x is a Gaussian process
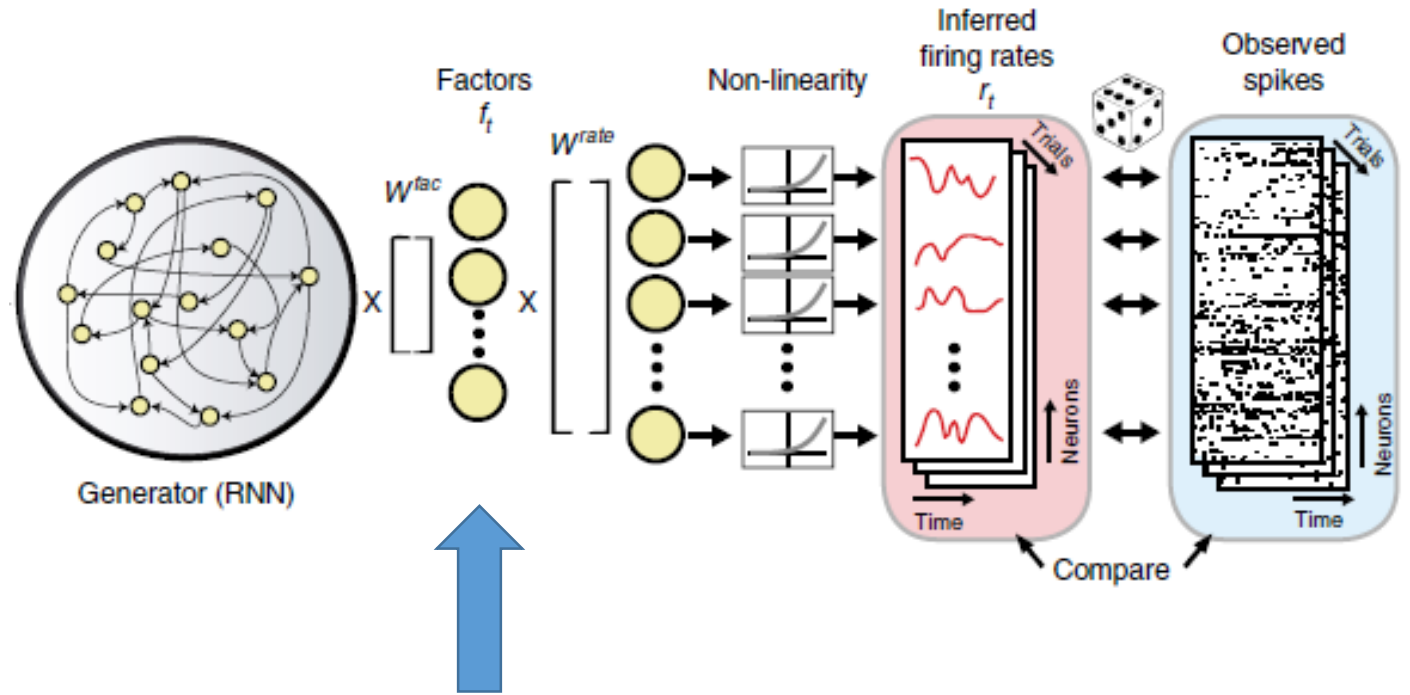
$$\mathbf{x}_{i,:} \sim \mathcal{N}\left(\mathbf{0},\ K_i\right)$$

- *Poisson feed-forward linear dynamical system (PfLDS):* y is generated from x through a rate $\lambda$ and a noise model P

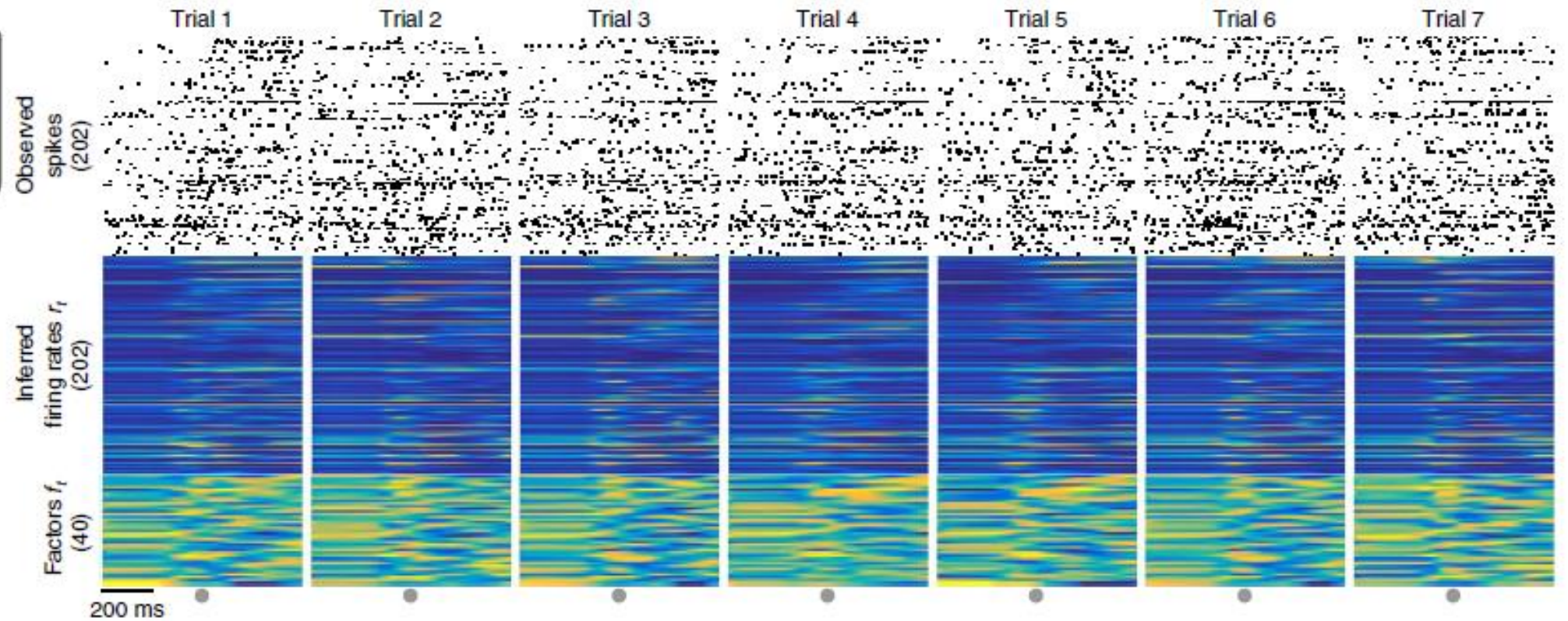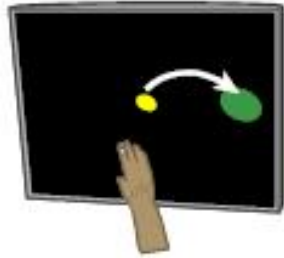$$x_{rti} \mid \mathbf{z}_{rt} \sim \mathcal{P}_\lambda\left(\lambda_{rti} = [f_\psi(\mathbf{z}_{rt})]_i\right)$$



*X*

*p x T*

*Y*

*q x T*
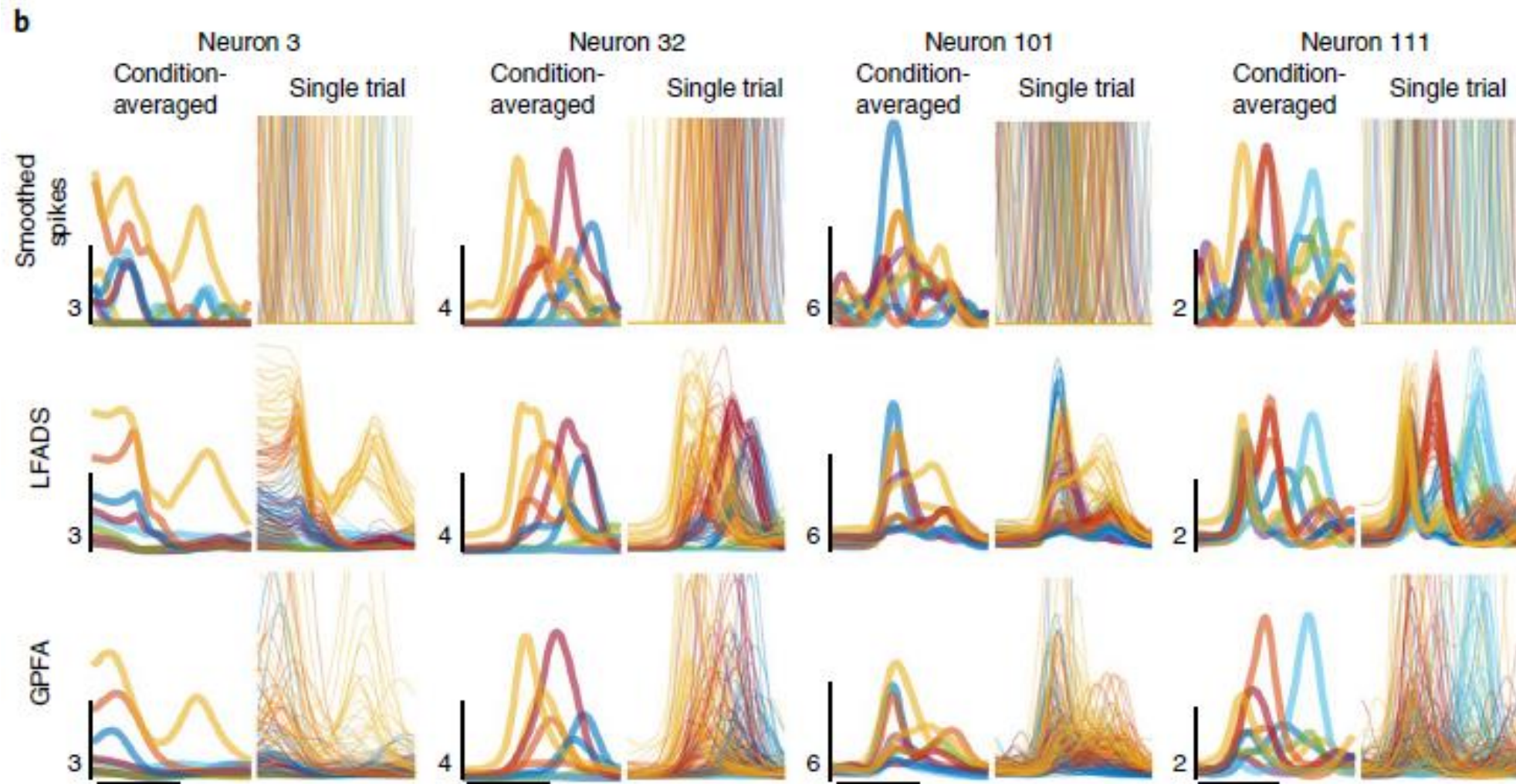
Roweis, 2004. Factor analysis and PCA

17

# LFADS basic architecture



Low-dimensional factors which generate firing rates

# LFADS can recover low-dimensional factors

**b**

| Neuron 3 | Neuron 32 | Neuron 101 | Neuron 111 |
|---|---|---|---|
| Condition-averaged / Single trial | Condition-averaged / Single trial | Condition-averaged / Single trial | Condition-averaged / Single trial |

Smoothed spikes

LFADS

GPFA

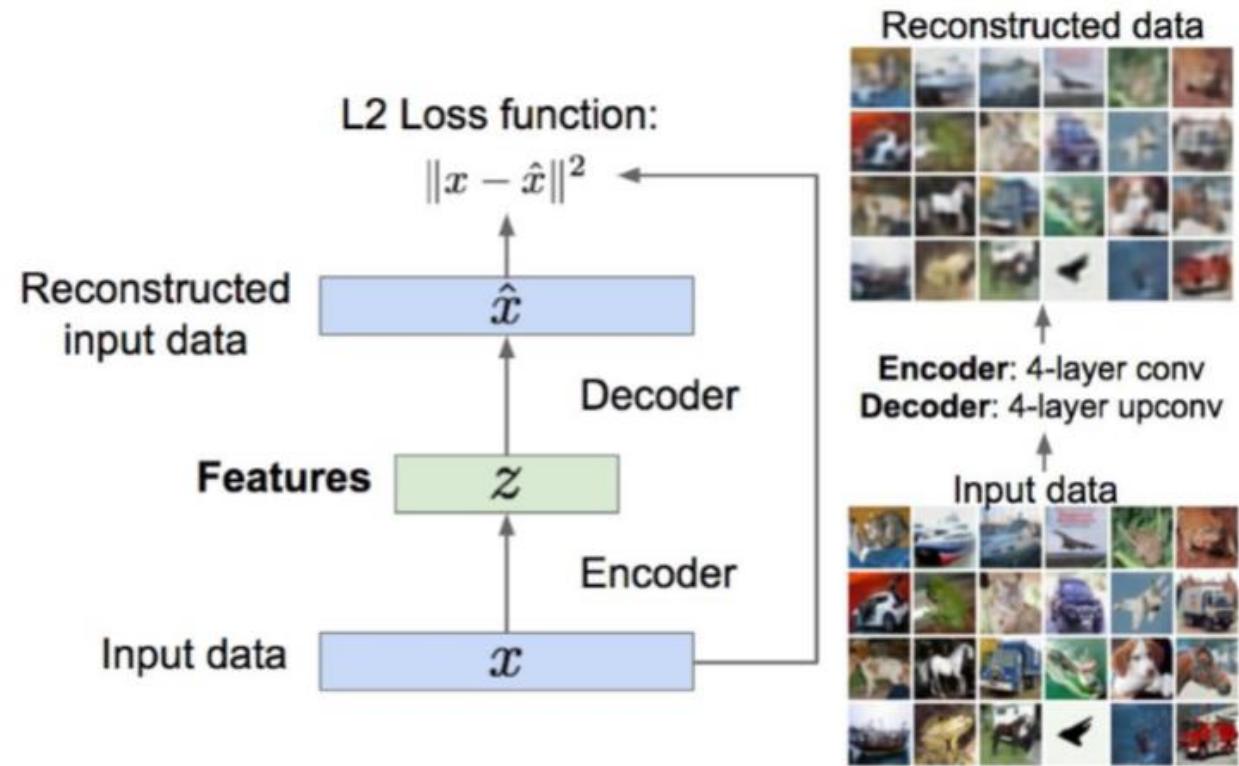# Part I: The LFADS basic architecture

1. LFADS assumes an underyling *dynamical system* (*recurrent neural network*) that generates spike trains

2. LFADS as a form of *factor analysis*

3. LFADS as a *variational autoencoder*

# Autoencoders

A form of **unsupervised learning**

Uncover hidden structure of data (in LFADS: spike trains)

After training, features (factors) will represent a compressed structure of the input data



L2 Loss function:
$$\|x - \hat{x}\|^2$$

Reconstructed input data $\hat{x}$

Decoder

**Features** $z$

Encoder

Input data $x$

Reconstructed data

**Encoder**: 4-layer conv
**Decoder**: 4-layer upconv

Input data

Autoencoders
(Feature learning)

# Variational Autoencoder (VAE)

- Goal: try to learn the probability distribution $p(x)$ that generates training data $x$

VAEs define intractable density function with latent **z**:

$$p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$$

Input data

$x$

Latent variable

$z$

Reconstruction

$\hat{x}$

**z** follows some
parameterized distribution

# Variational Autoencoder (VAE)

- Goal: try to learn the probability distribution $p(x)$ that underlie training data $x$

VAEs define intractable density function with latent $\mathbf{z}$:

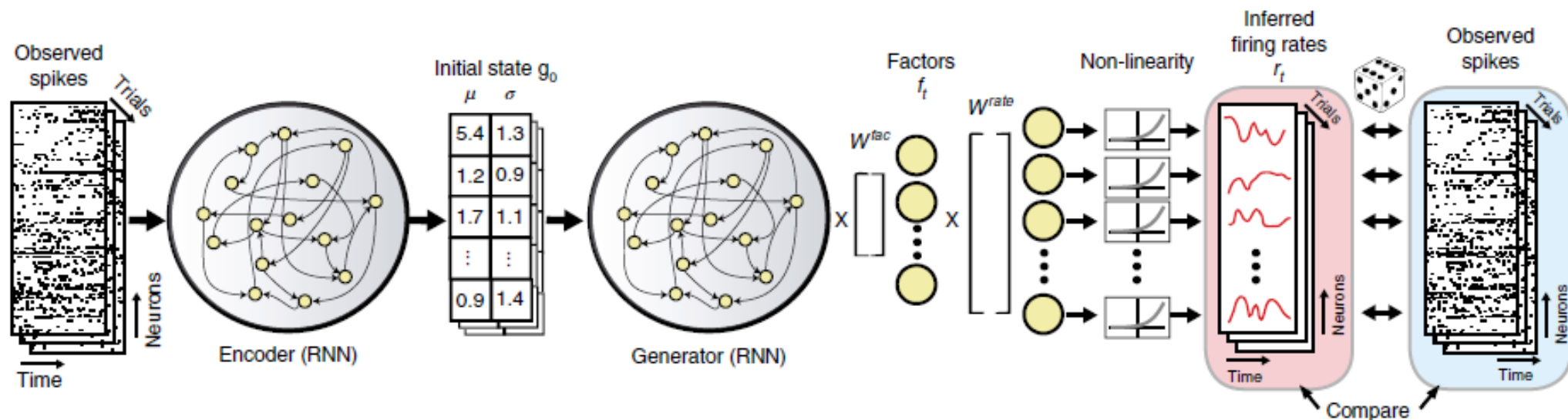$$p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$$

Input data

Latent variable

Reconstruction

$$x$$
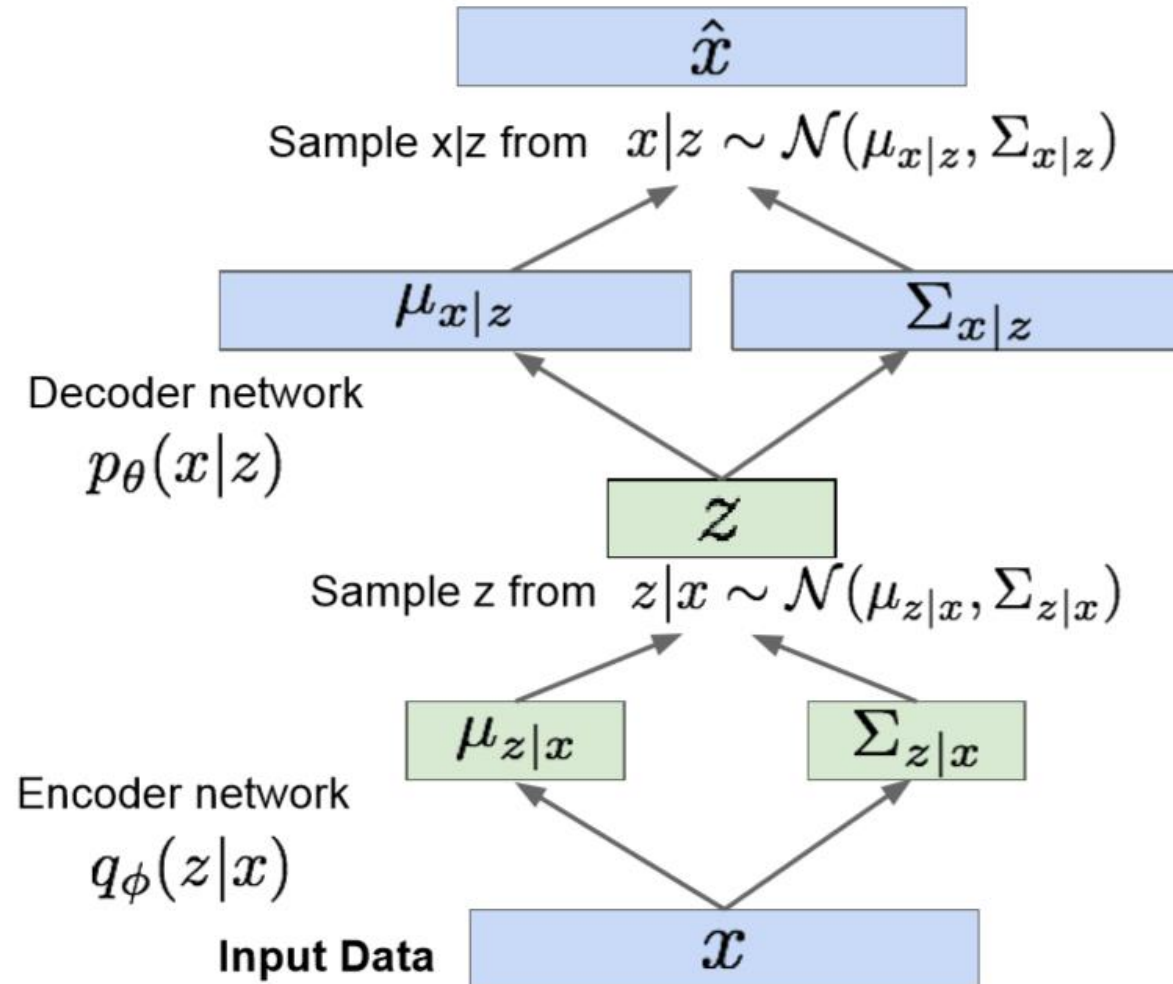
$$z$$

$$\hat{x}$$

# Variational Autoencoder (VAE): Training

Reconstructed spike trains

RNN (decoder)

$g_0$

RNN (encoder)

Spike trains



Sample x|z from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

Decoder network $p_\theta(x|z)$

Sample z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Encoder network $q_\phi(z|x)$

**Input Data**

*Yeung, CS231n (Stanford) 2017, lecture13*

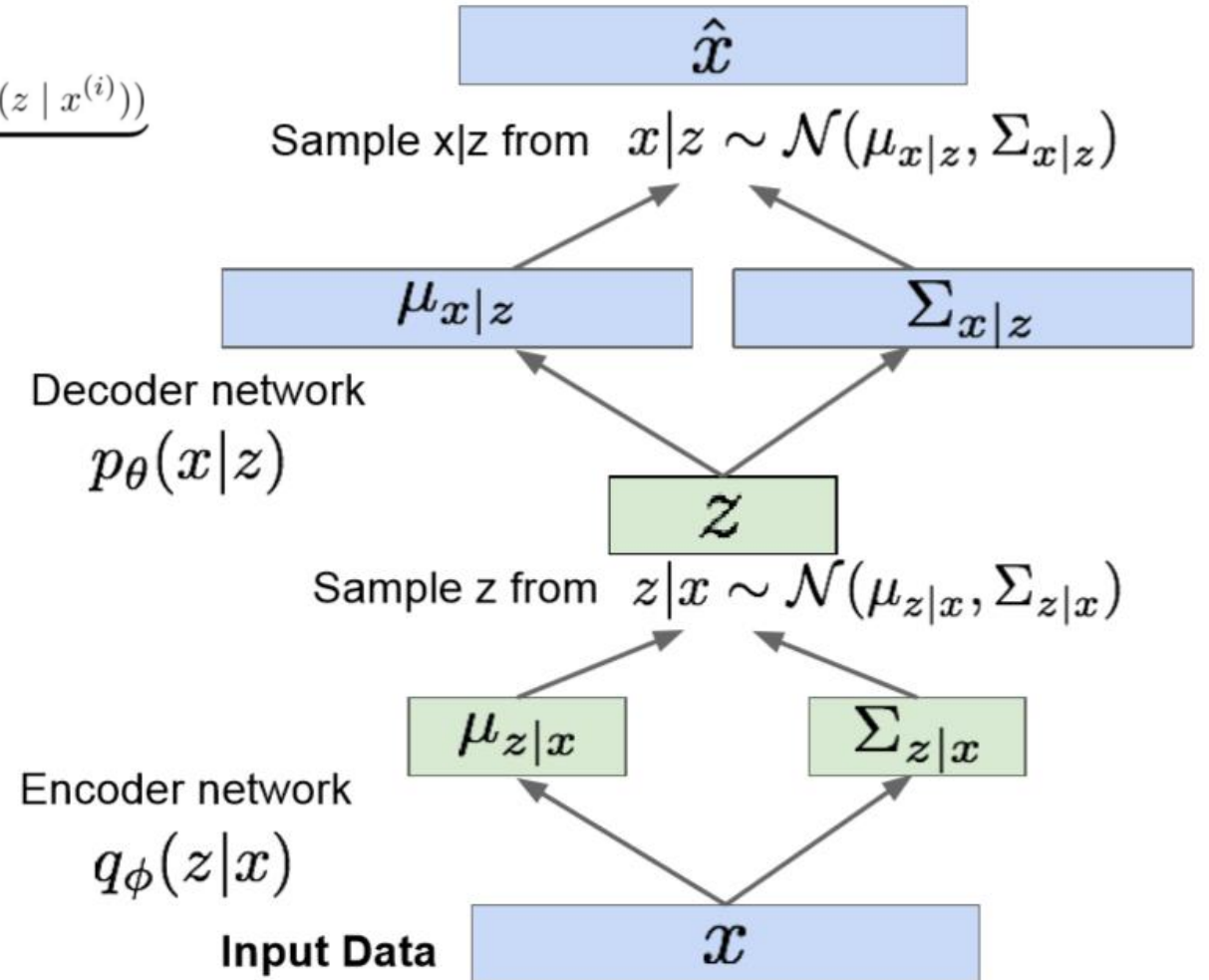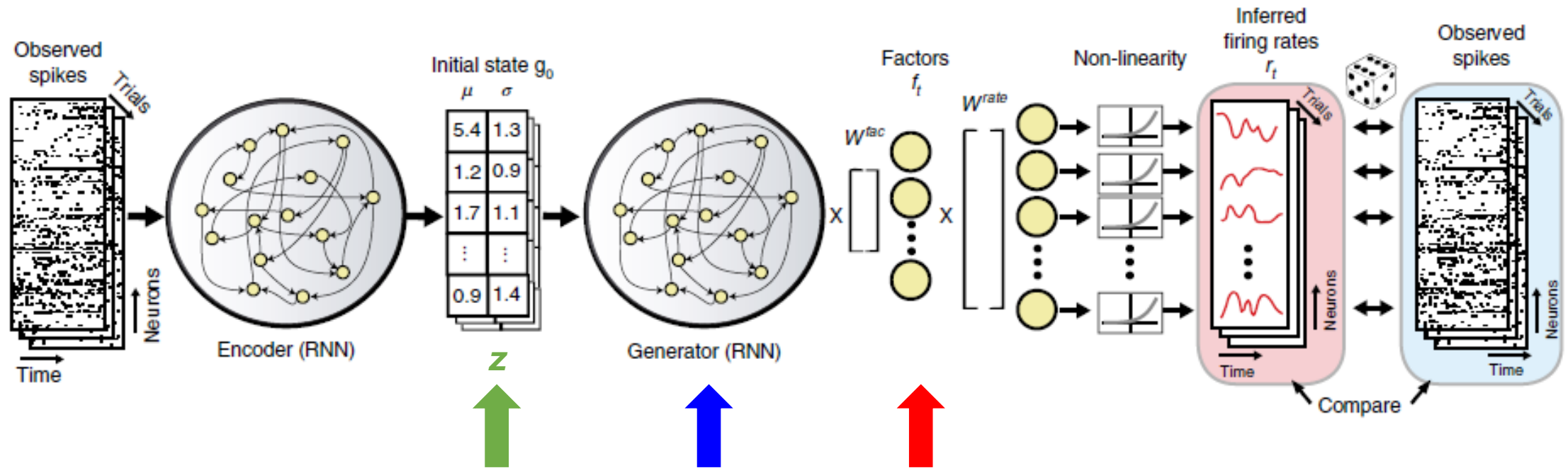# Variational Autoencoder (VAE): Training

$$\log p_\theta(x^{(i)}) =$$

$$= \underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))}_{\geq 0}$$

**Tractable lower bound** which we can take gradient of and optimize! ($p_\theta$(x|z) differentiable, KL term differentiable)



Sample x|z from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

Decoder network $p_\theta(x|z)$

Sample z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Encoder network $q_\phi(z|x)$

**Input Data**
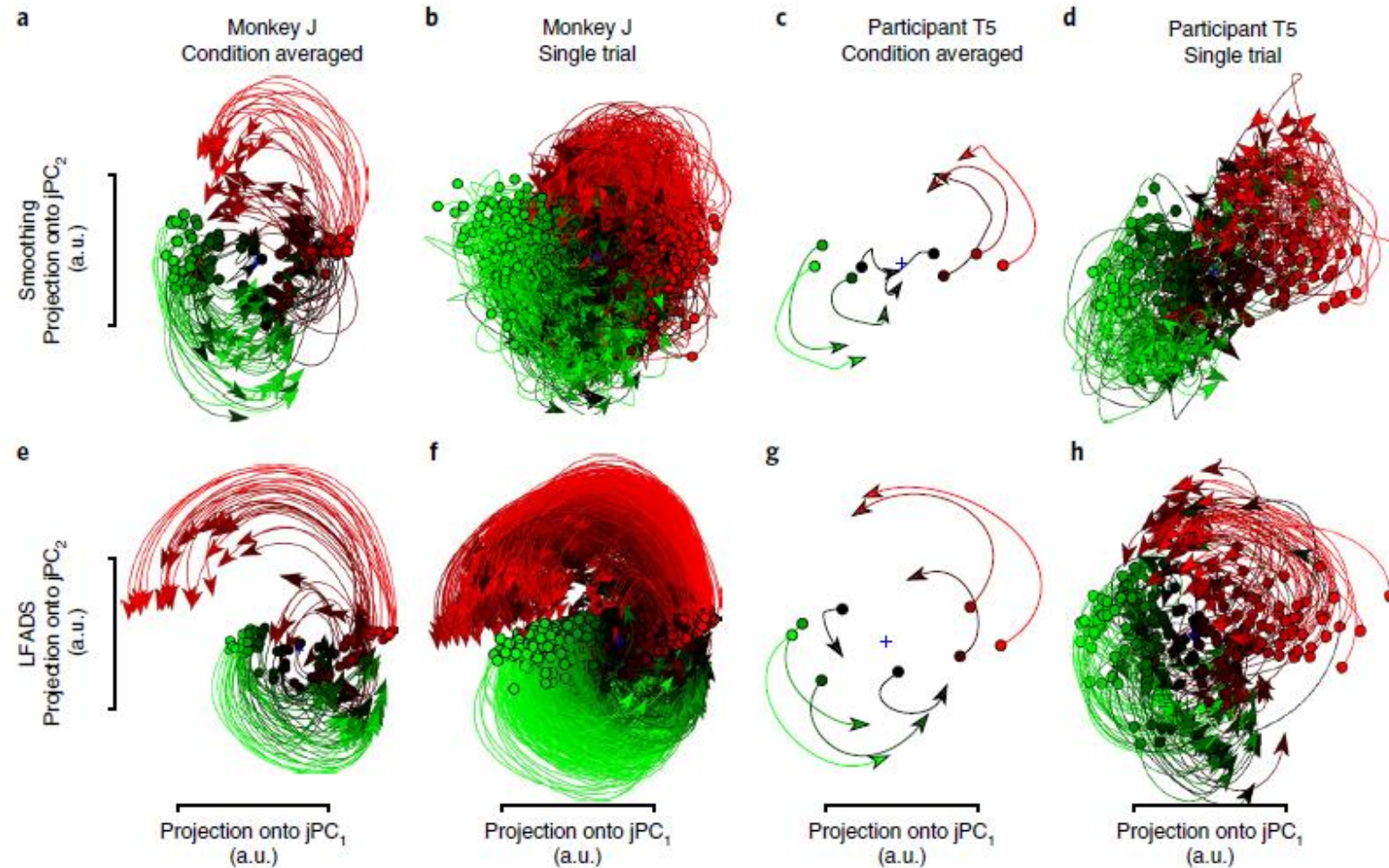
*Yeung, CS231n (Stanford) 2017, lecture13*

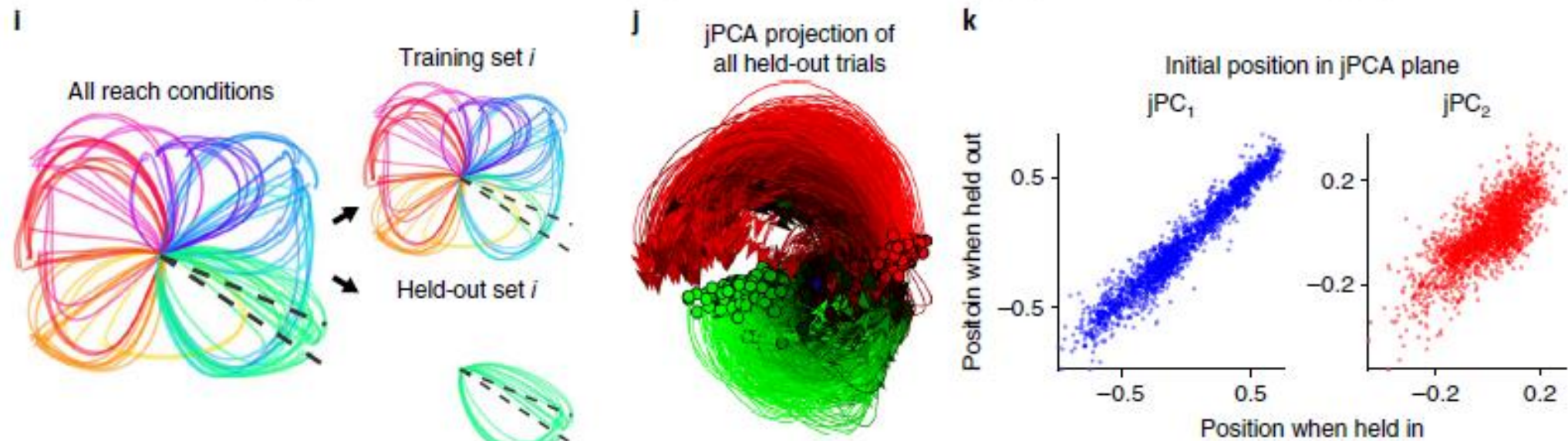# Summary: the LFADS basic architecture



1. LFADS assumes a **dynamical system** underlying spike trains
2. LFADS uncovers **low-dimensional factors** that underlie firing rates
3. LFADS is trained as a variational autoencoder (**VAE**)

28

# LFADS uncovers single-trial rotation dynamics



**a** Monkey J Condition averaged  
**b** Monkey J Single trial  
**c** Participant T5 Condition averaged  
**d** Participant T5 Single trial

**e** Smoothing Projection onto jPC₂ (a.u.) / LFADS Projection onto jPC₂ (a.u.)

Projection onto jPC₁ (a.u.)

# LFADS uncovers single-trial rotation dynamics
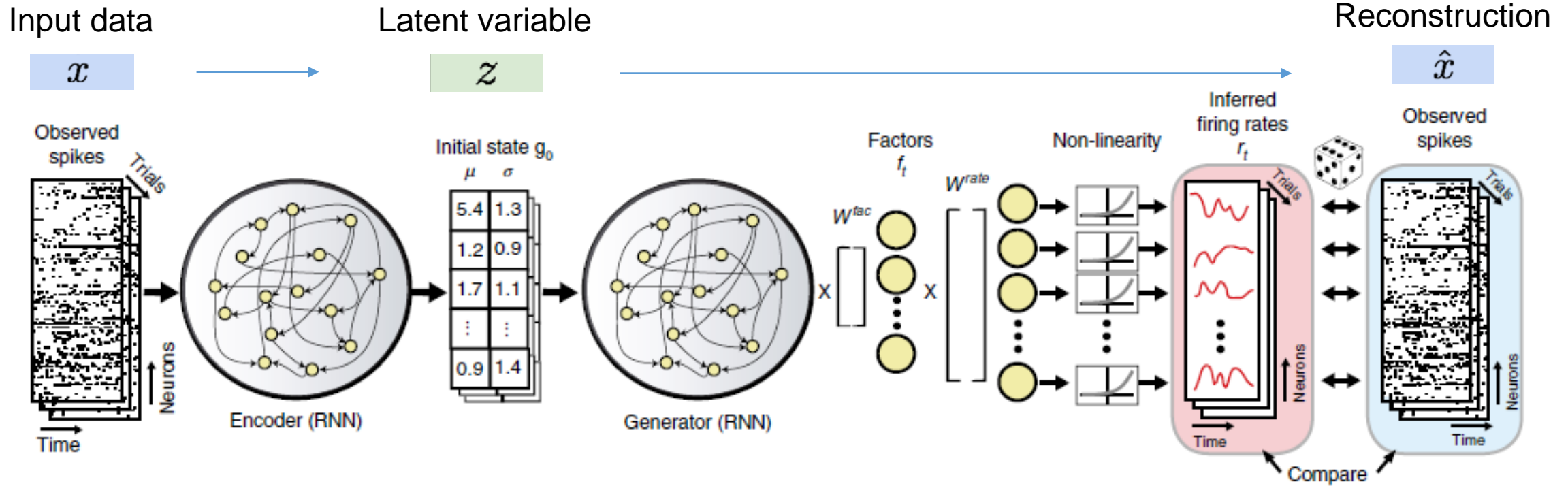


LFADS performs well on held-out trials

# Part II: Variants of LFADS

1. Dynamic neural stitching
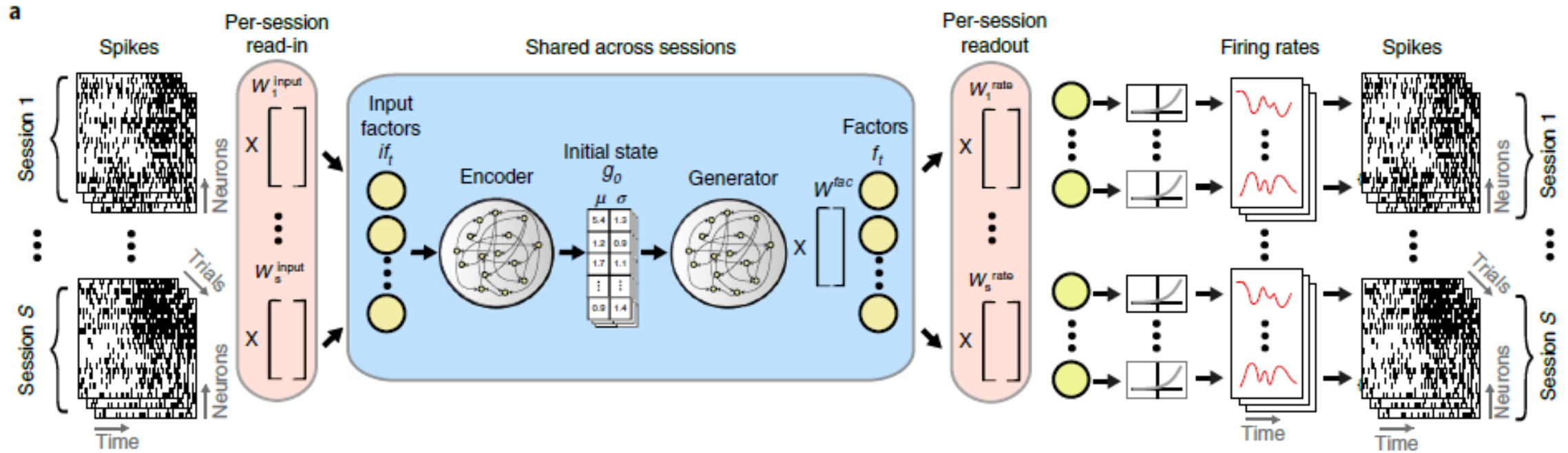
2. LFADS with external perturbations

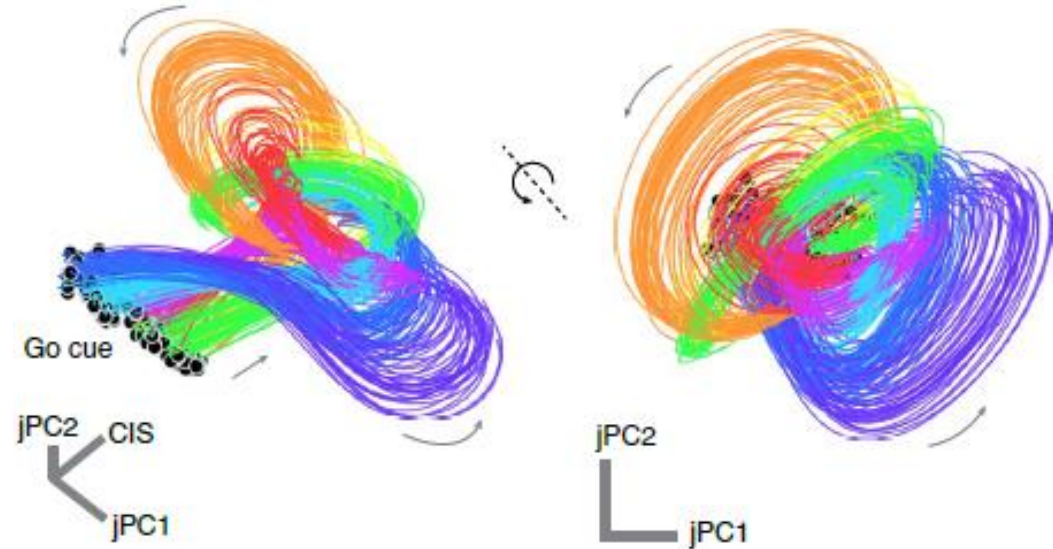# Dynamic neural stitching

# LFADS basic architecture
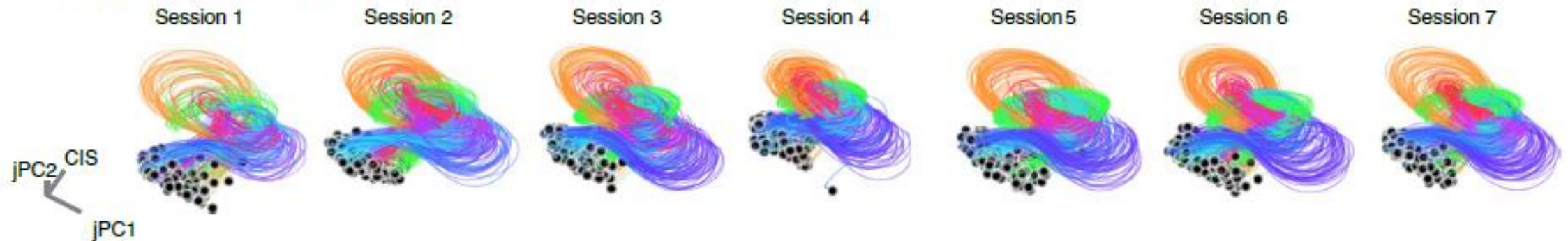
# Dynamic neural stitching

# Dynamic neural stitching

- Consistent trajectories across sessions → consistent with a single underlying dynamical system



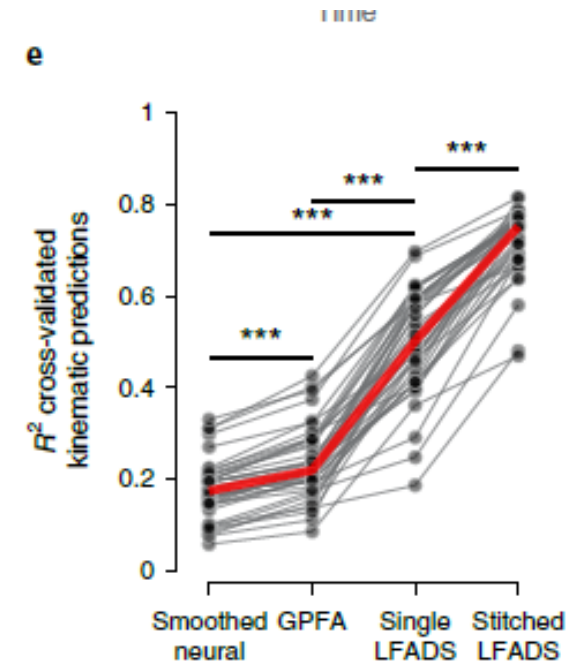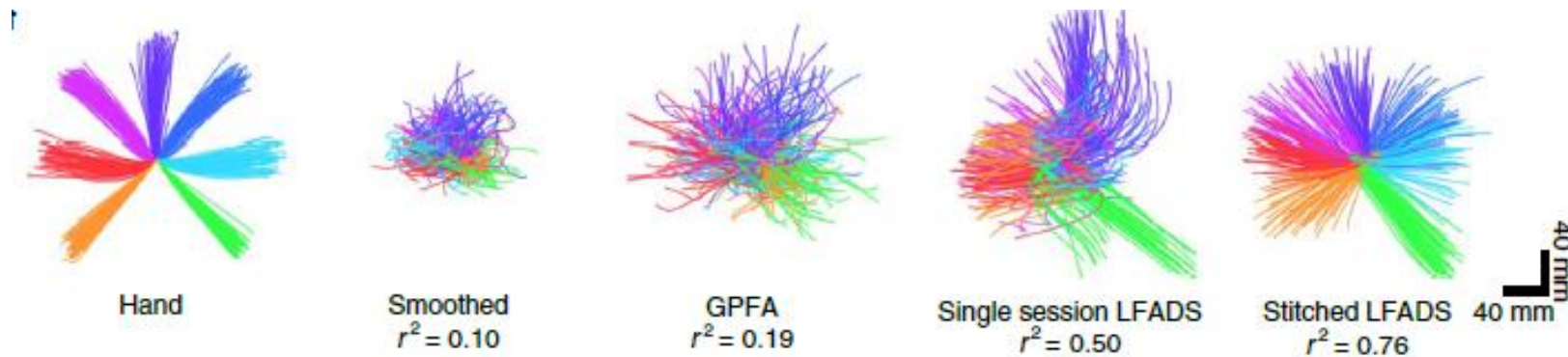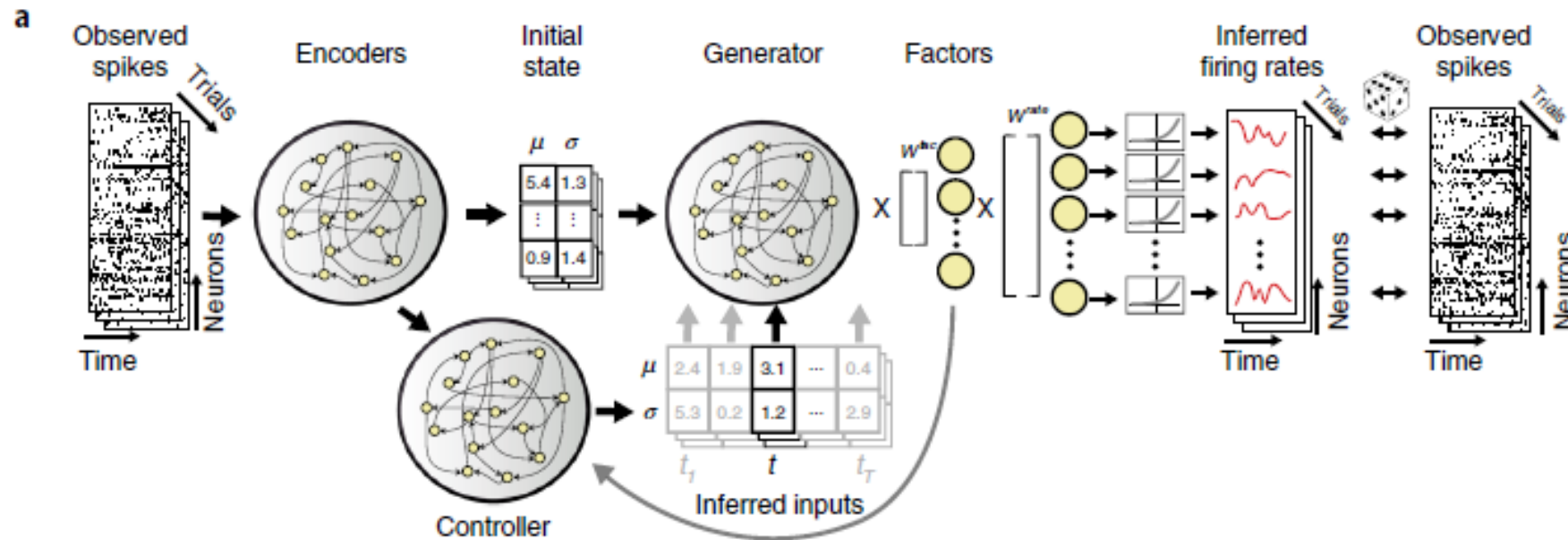**d** Condition-averaged LFADS factor trajectories across sessions

Go cue

jPC2 CIS
jPC1

jPC2
jPC1

**g** Single-trial LFADS factor trajectories

Session 1    Session 2    Session 3    Session 4    Session 5    Session 6    Session 7

jPC2 CIS
jPC1

# Dynamic neural stitching

- Good decoding of kinematic variables using LFADS factors



Hand

Smoothed
$r^2 = 0.10$

GPFA
$r^2 = 0.19$

Single session LFADS
$r^2 = 0.50$

Stitched LFADS
$r^2 = 0.76$

40 mm



e

$R^2$ cross-validated kinematic predictions

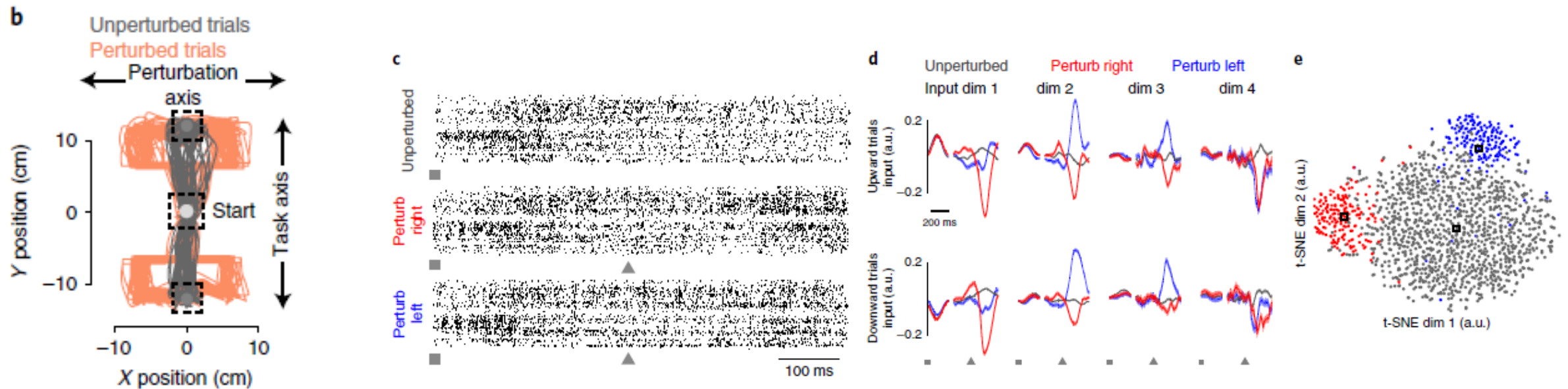Smoothed neural    GPFA    Single LFADS    Stitched LFADS

# LFADS uncover structure of perturbations

- Behavior of dynamical system often disturbed by **external signals** from other brain areas
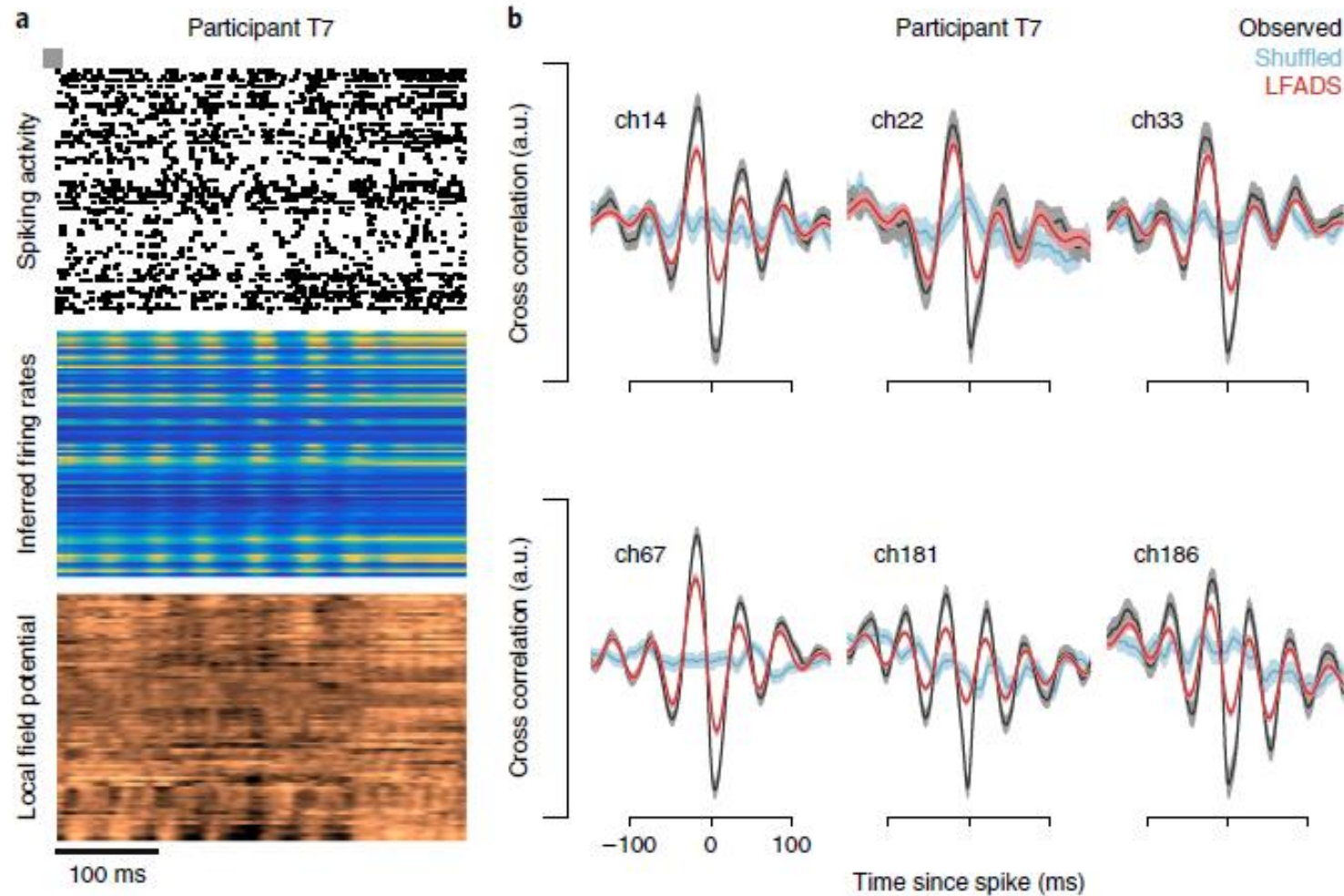→ Can we infer these signals on single-trials?



$g_0$ and $u(t)$ are the latent variables $z$ in the VAE

# LFADS uncover structure of perturbations



External inputs during a cursor manipulation task when a perturbation is given in some of the trials

# LFADS uncover structure of perturbations



Oscillatory inputs before movement initiation, synchronized to LFPs

# Discussion

1. In what context would LFADS be appropriate to model neural data? When will it **not** be appropriate?