








Direct neural perturbations reveal a dynamical mechanism for robust computation

 Daniel J. O'Shea,  Lea Duncker, Werapong Goo,  Xulu Sun,  Saurabh Vyas,  Eric M. Trautmann,  Ilka Diester,  Charu Ramakrishnan,  Karl Deisseroth,  Maneesh Sahani,  Krishna V. Shenoy

doi: <https://doi.org/10.1101/2022.12.16.520768>

 3  0  4  0

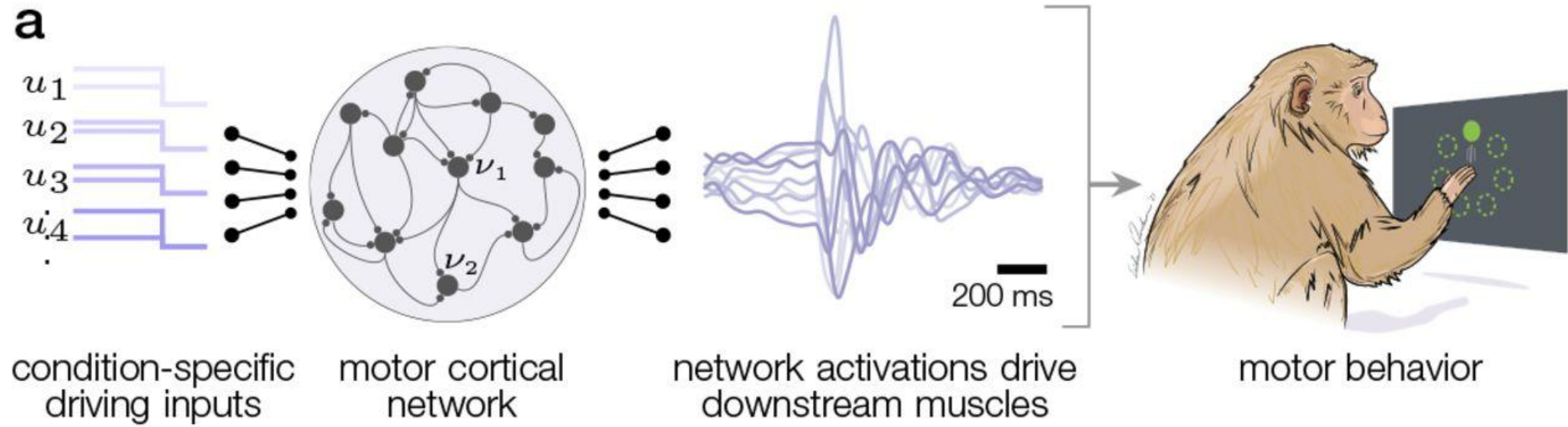
This article is a preprint and has not been certified by peer review [what does this mean?].

Abuzar
Mahmood
Katz Lab

CSN IC 2/15/24

- Intro + Background
 - Hypotheses (total of 3)
 - Experimental Setup
- Results (knock out 1 hypothesis with each section):
 - Optogenetics
 - Intra-Cranial Microstimulation

Introduction

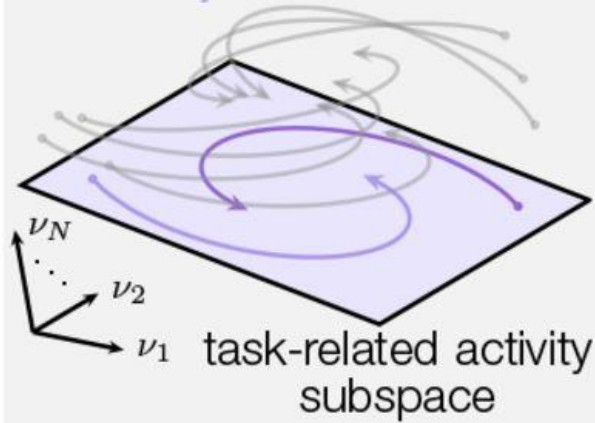


Introduction :

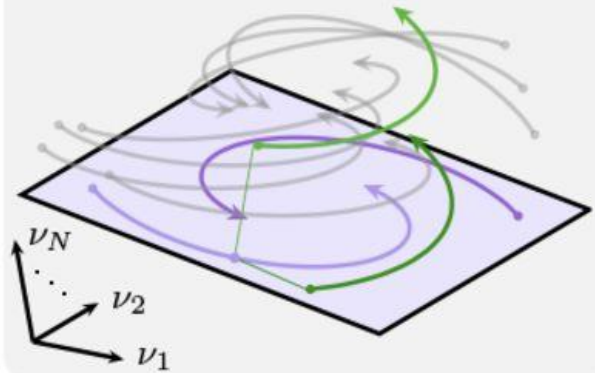
Hypotheses

H1: high-d reservoir dynamics

b rich set of trajectories; subset selected by initial condition

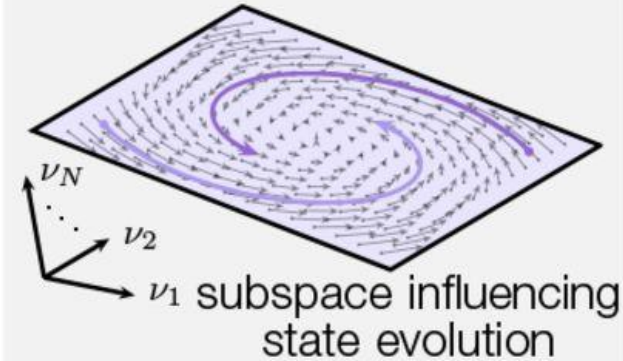


e perturbed trajectory

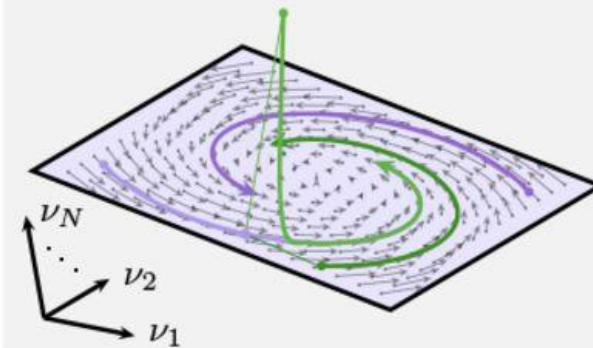


H2: low-d subspace structured dynamics

c dynamics confined to low-d subspace; trajectory selected by initial condition

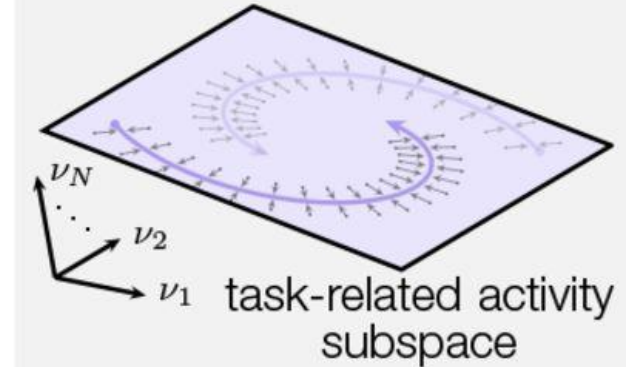


f perturbed trajectory

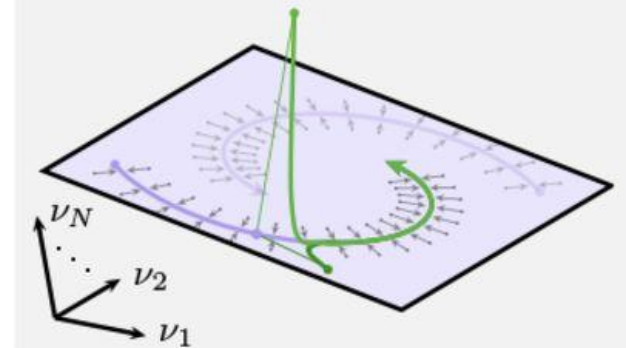


H3: path-following dynamics

d state trajectory tracks externally configured path

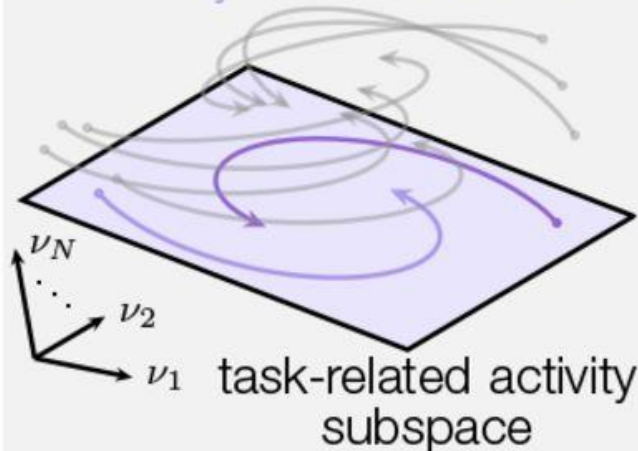


g perturbed trajectory

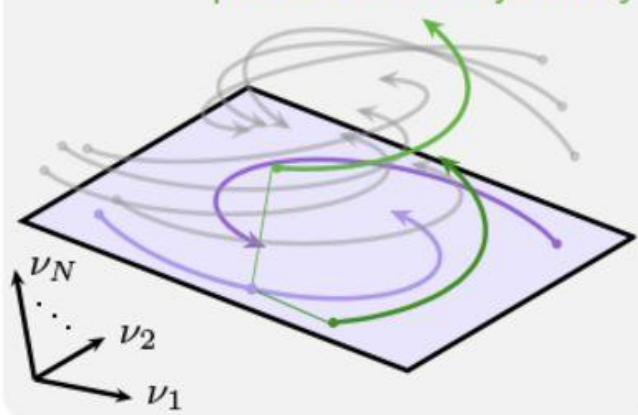


H1: high-d reservoir dynamics

b rich set of trajectories;
subset selected
by initial condition



e perturbed trajectory



Hypothesis

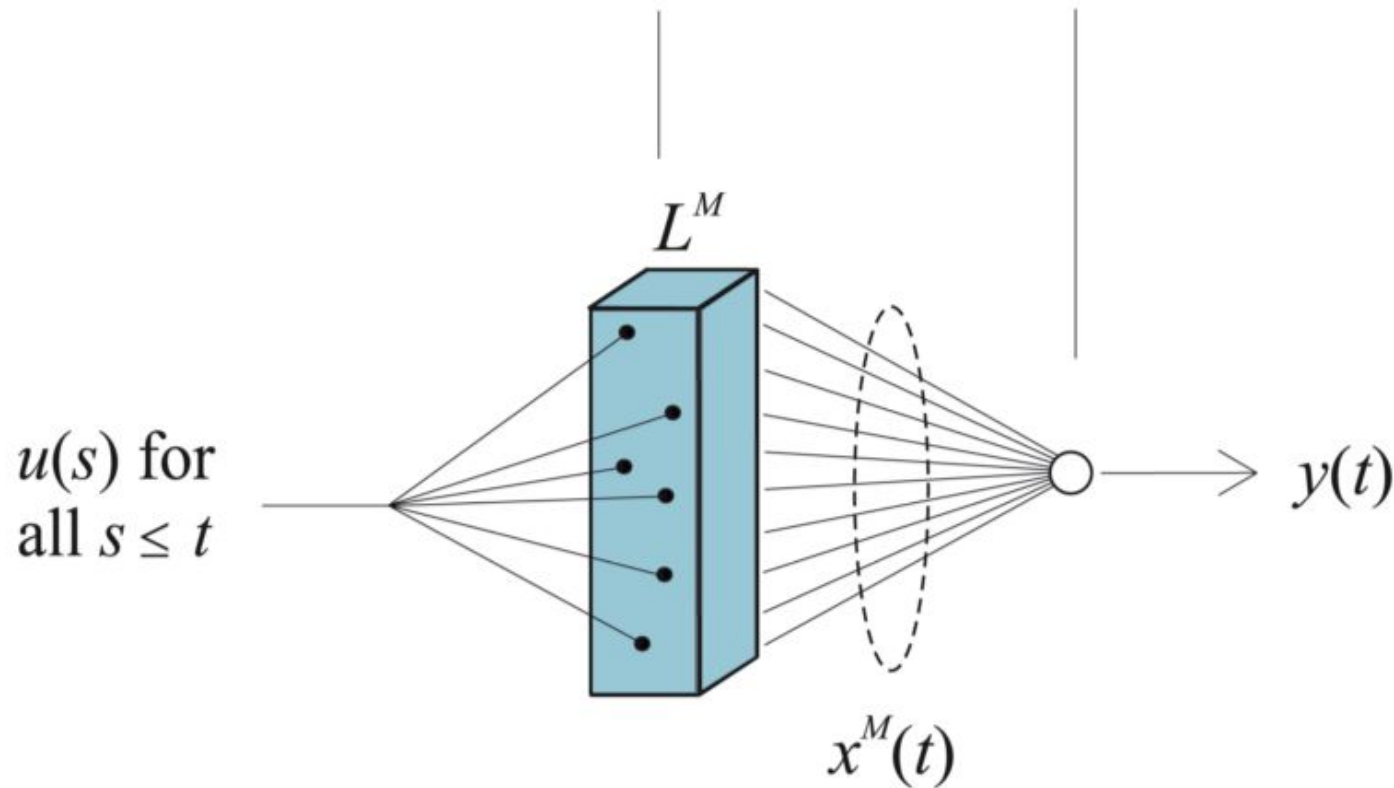
1

- In reservoir networks, which retain the concurrent potential to express basis patterns in many different dimensions, perturbations of neural state should engage new dynamical modes and thus evoke complex, long-lasting transients (**Fig. 1e**)

Background : Reservoir Computing

a bank of basis filters
(or some more general dynamical system)

memoryless readout,
trained for a specific task

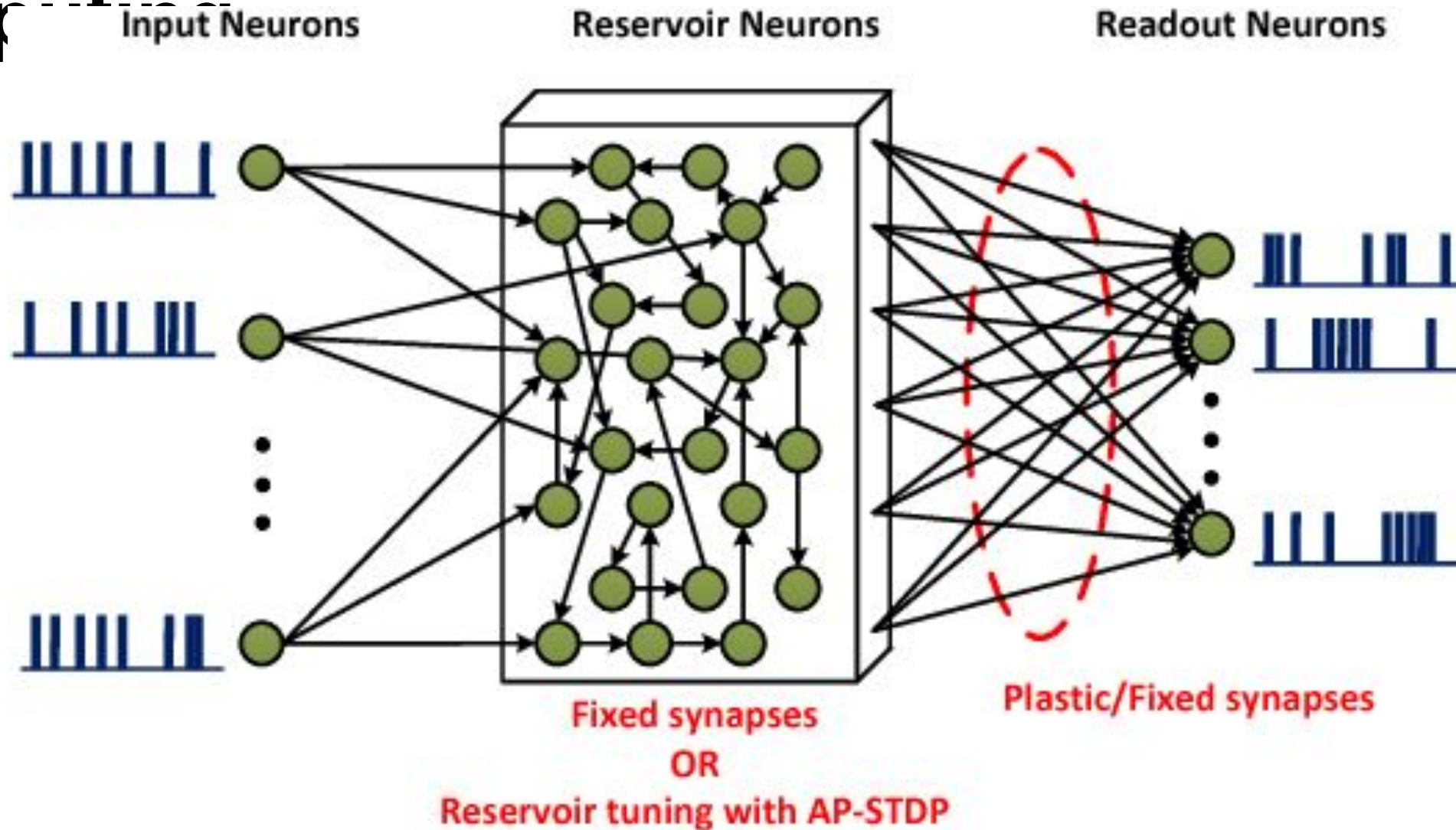


$$x^M(t) = (L^M u)(t)$$
$$y(t) = f^M(x^M(t))$$

= liquid state of the
Liquid State Machine

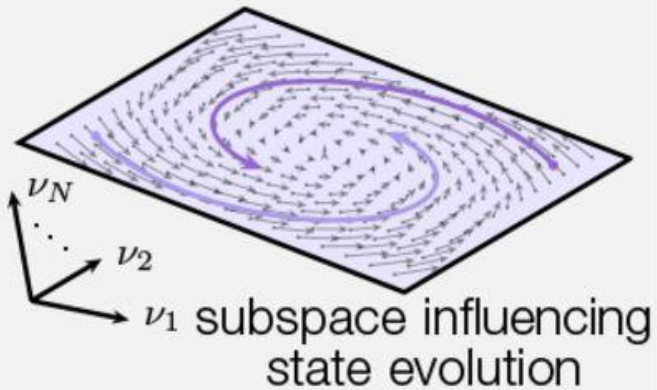
Background : Reservoir

Comp

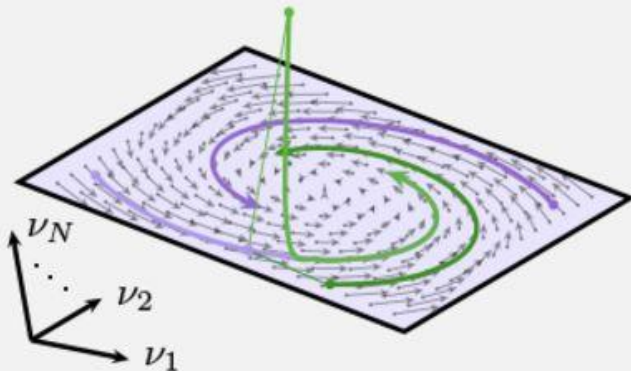


H2: low-d subspace structured dynamics

c dynamics confined to low-d subspace; trajectory selected by initial condition



f perturbed trajectory



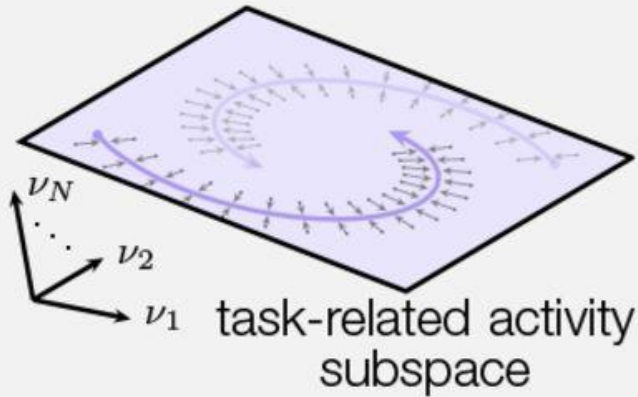
Hypothesis 2

- In subspace structured networks, where activity is driven by the state within a low-dimensional subspace, only perturbations that affect this subspace should elicit complex, long-lasting effects, while perturbations along all other dimensions should fail to engage with the circuit dynamics (**Fig. 1f**).

H3: path-following dynamics

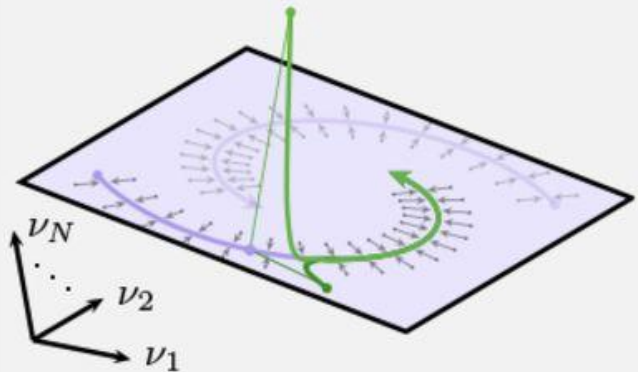
d

state trajectory tracks
externally configured path



g

perturbed trajectory



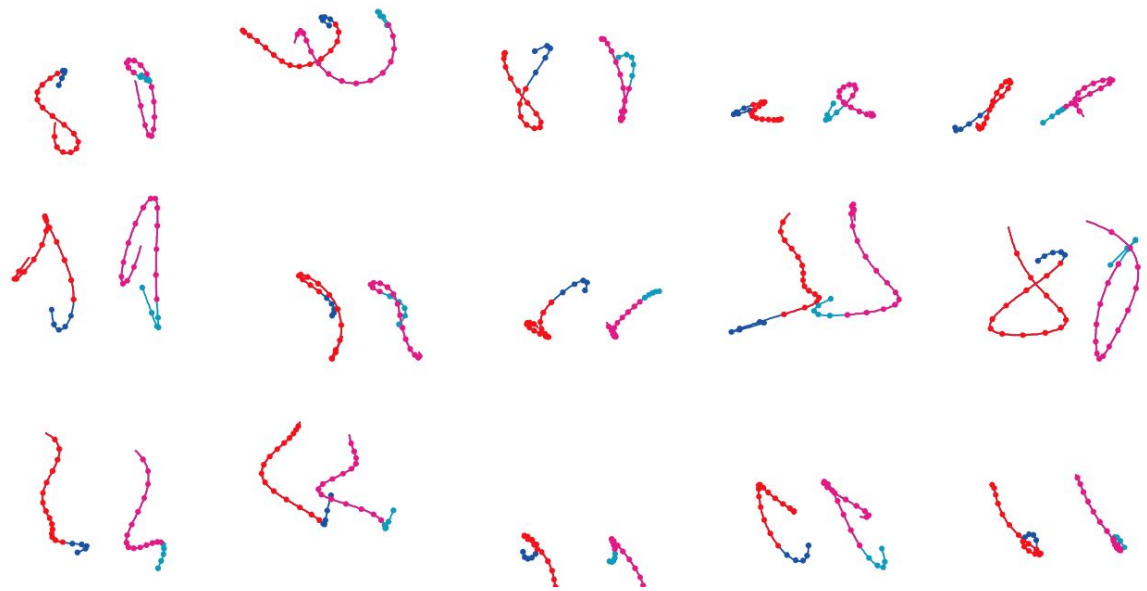
Hypothesis 3

- Lastly, in a path-following network, all perturbations away from the externally-configured trajectory will decay back rapidly (Fig. 1g).

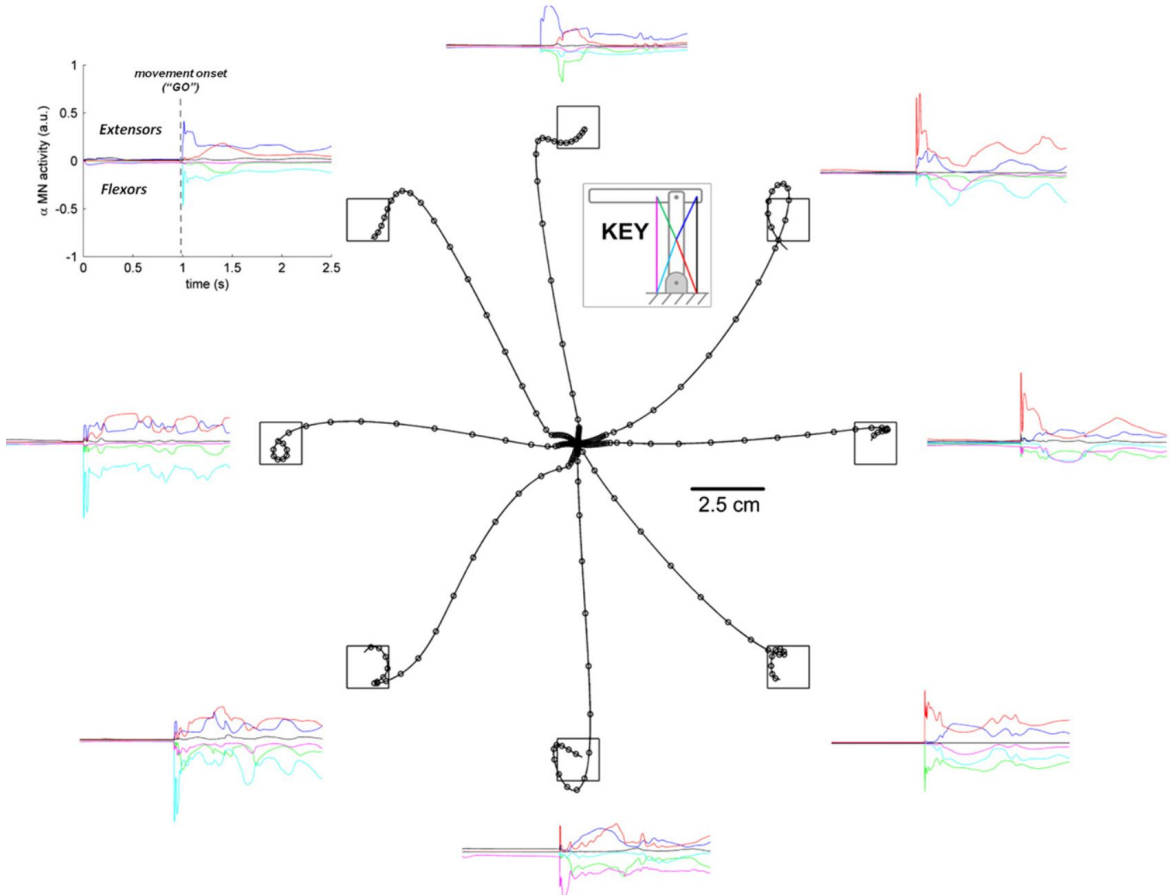
Background : Path-following Dynamics

- *path-following dynamics* (**H3**),
 - in which the neural state is constrained to move along an externally configured path. In the context of motor control, previous work has proposed that the motor cortex might serve to activate specific motor programs implemented by recurrent circuitry in the spinal cord.

A



Hatsopoulos NG, Xu Q, Amit Y. Encoding of Movement Fragments in the Motor Cortex. J Neurosci. 2007 May 9;27(19):5105–14.

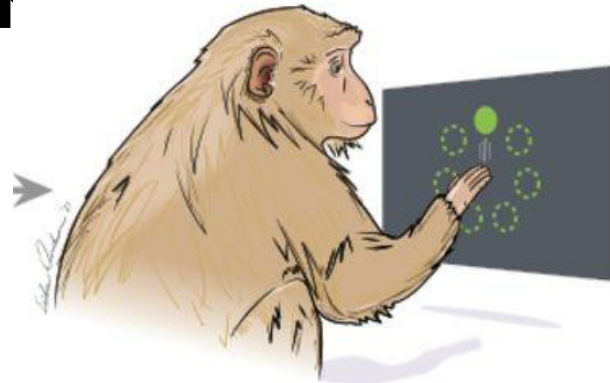


Tsianos GA, Goodner J, Loeb GE. Useful properties of spinal circuits for learning and performing planar reaches. J Neural Eng. 2014 Oct 1;11(5):056006.

- Full extent of hypotheses justification :/
 - No references on expected effects of perturbations
 - But...model for H1 perturbation later on.

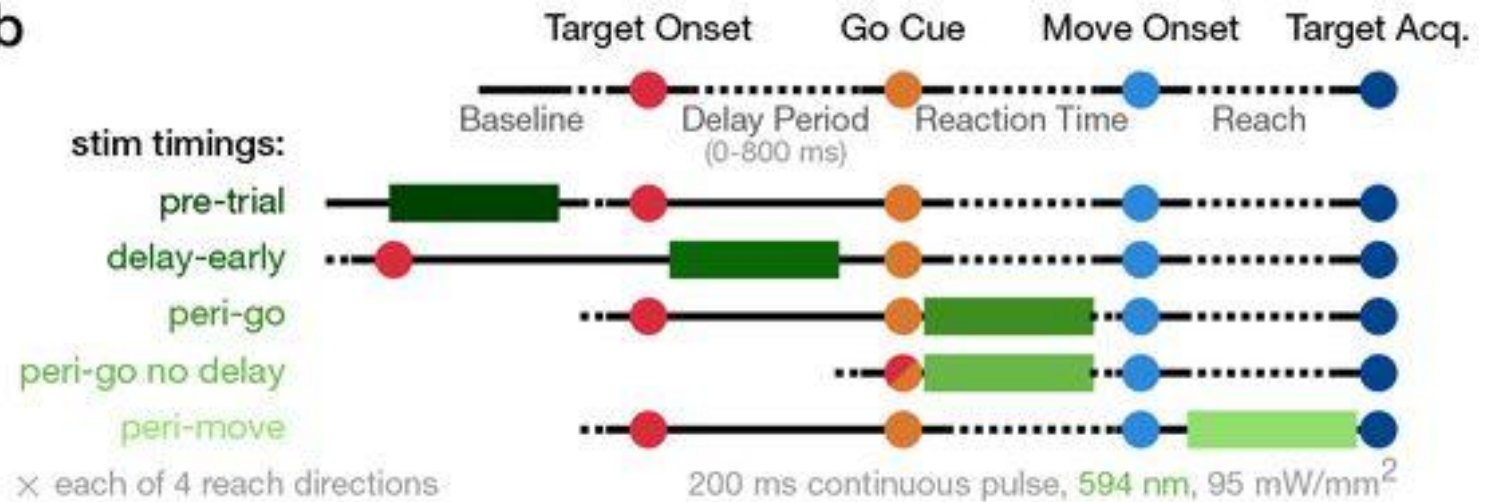
Experimental Setup

(of ...)

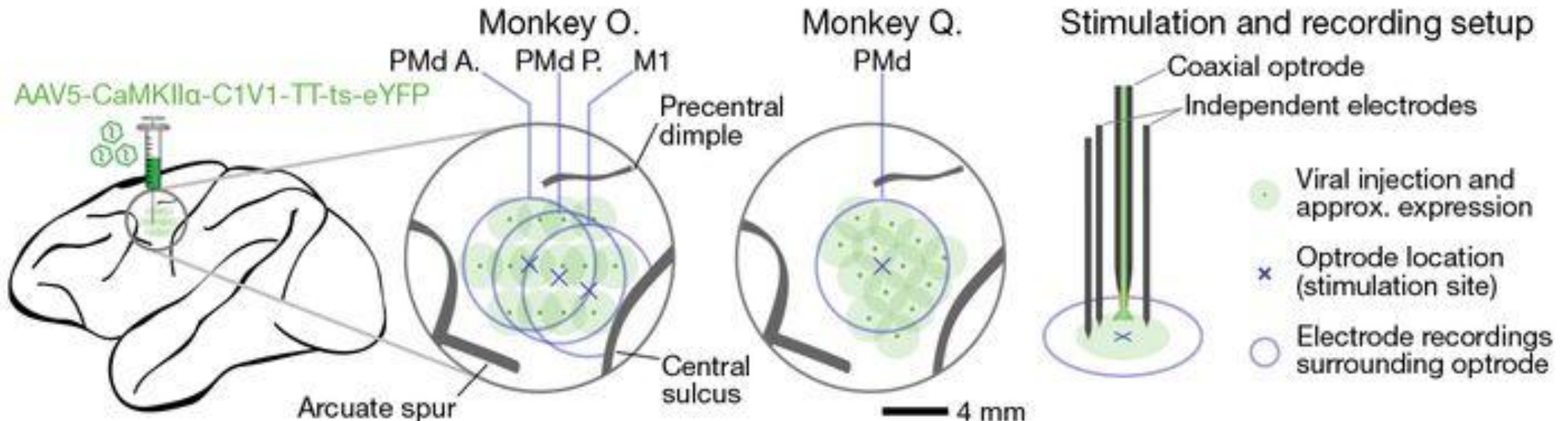


motor behavior

b

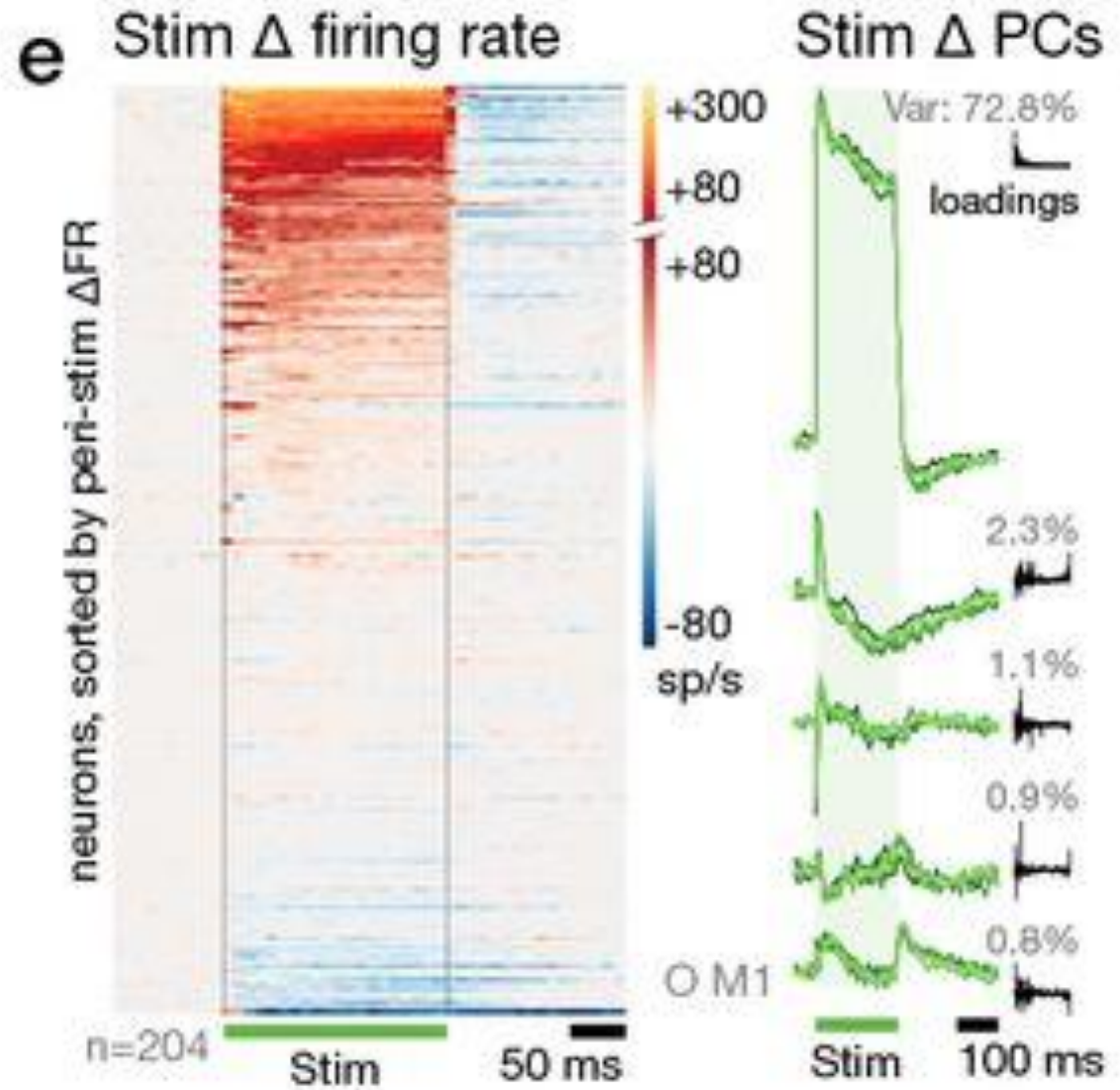
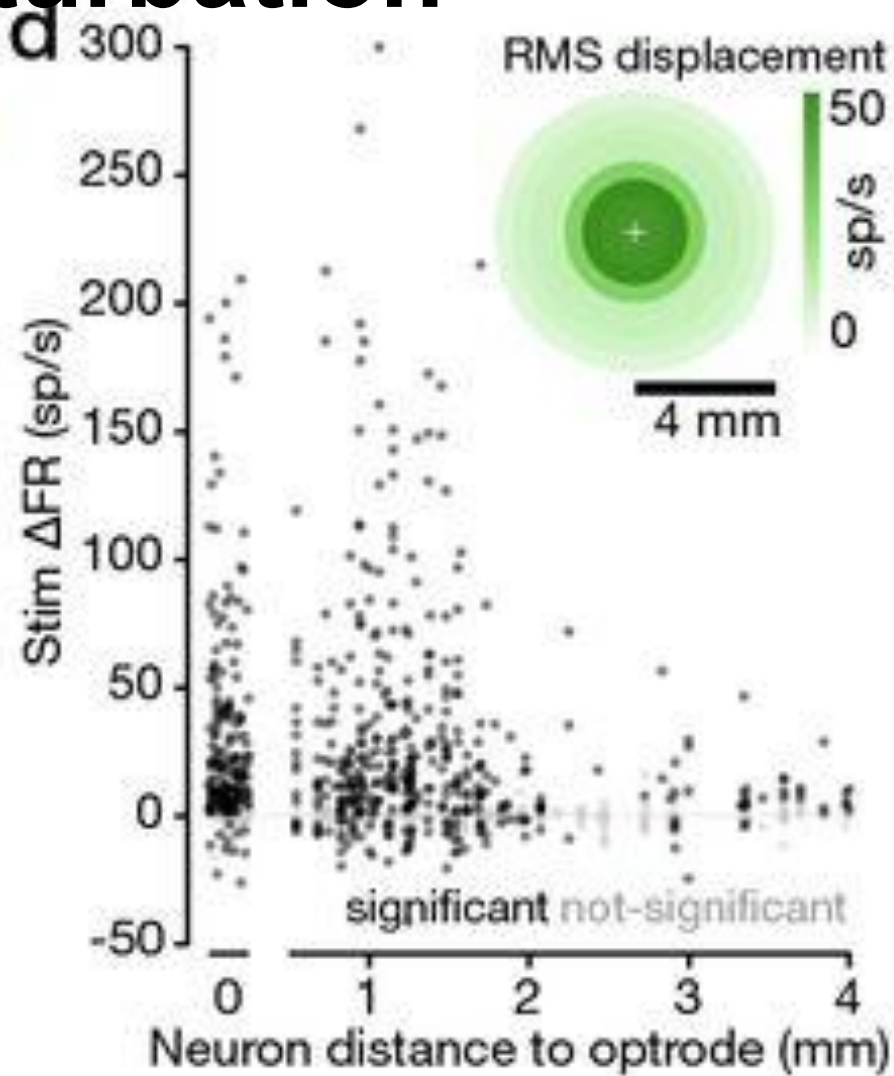


a



- Intro + Background
 - Hypotheses (total of 3)
 - Experimental Setup
- Results (knock out 1 hypothesis with each section):
 - Optogenetics
 - Intra-Cranial Microstimulation

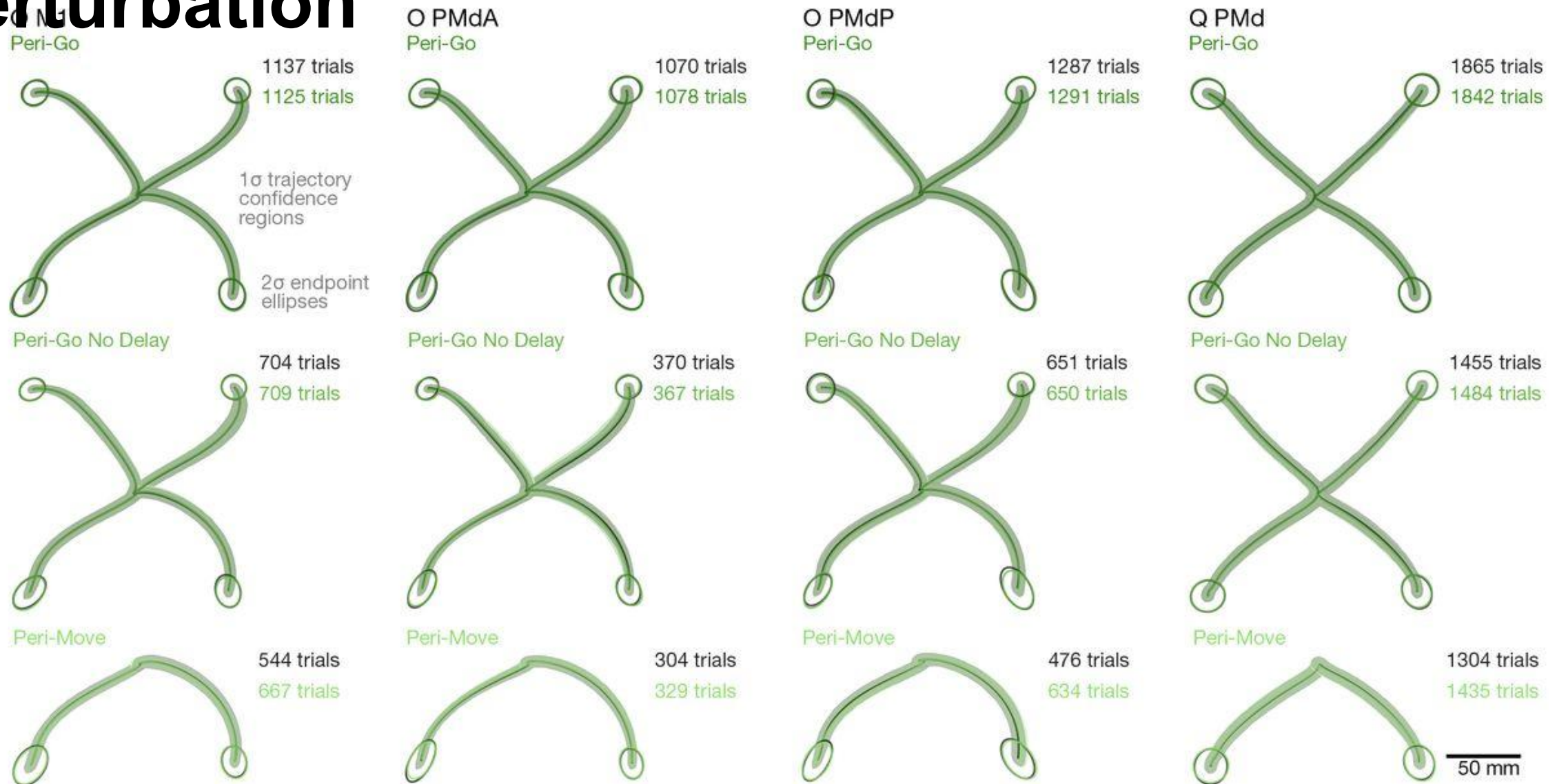
Neural Effects Of Optogenetic Perturbation



Neural Effects Of Optogenetic Perturbation

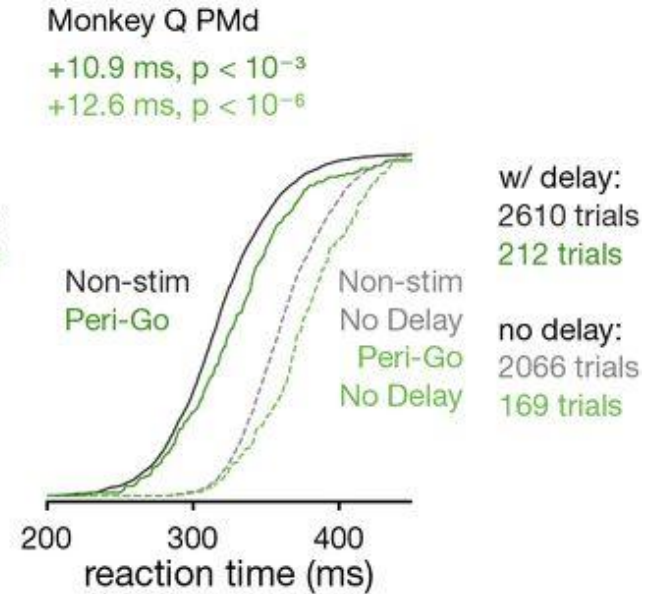
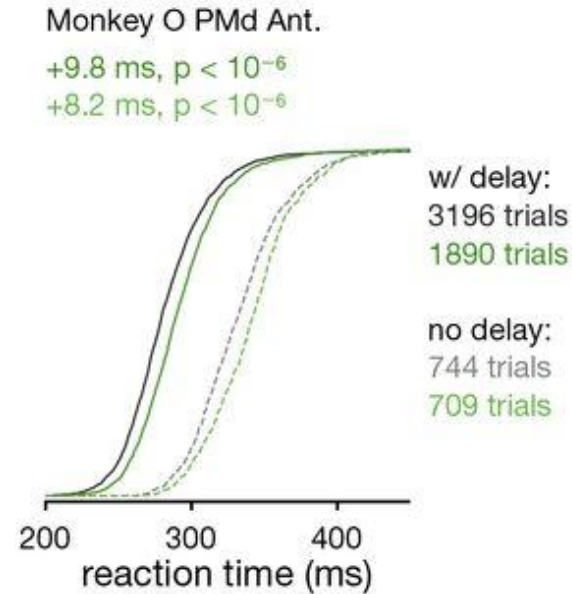
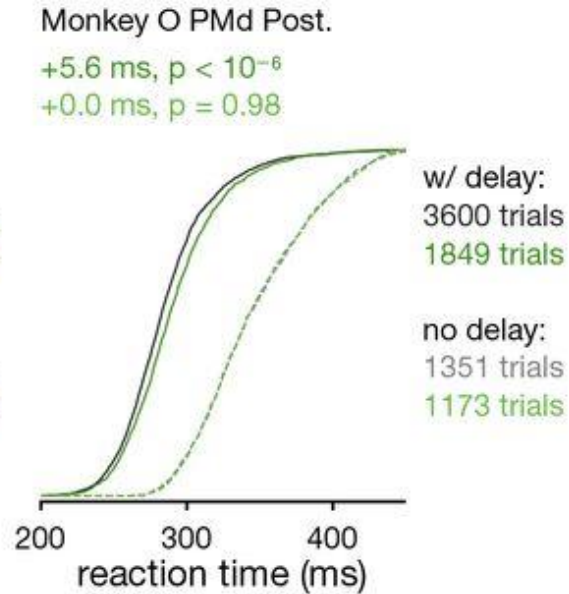
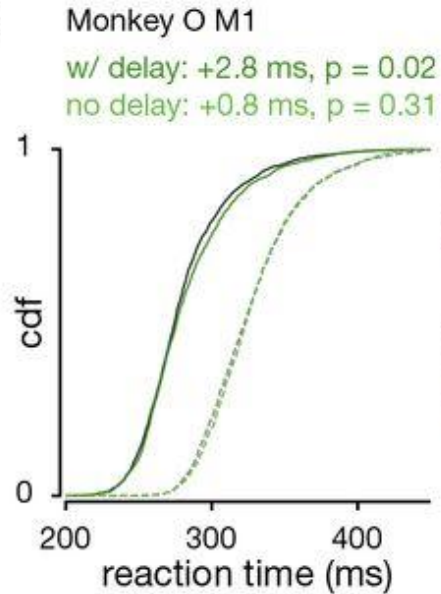


No Behavioral Effects Of Optogenetic Perturbation



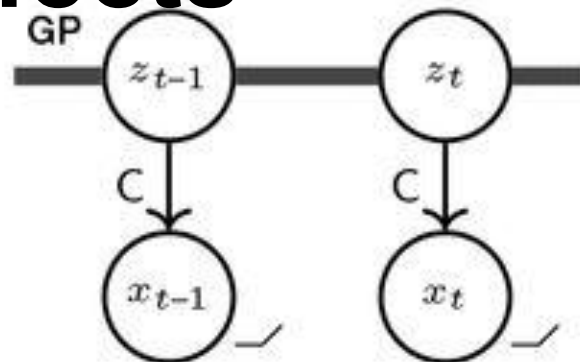
No Behavioral Effects Of Optogenetic Perturbation

c



Inferring Perturbation Effects

a non-stimulation neural data:



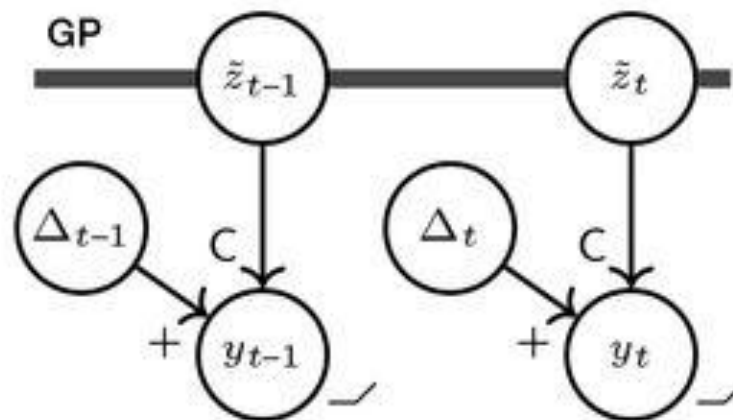
non-stim latents

$$z_t \in \mathbb{R}^K, \text{ with GP prior}$$

non-stim rates

$$x_t = g(Cz_t + d) \in \mathbb{R}^D$$

Stimulation neural data:



stim latents

$$\tilde{z}_t \in \mathbb{R}^K, \text{ with GP prior}$$

additive component

$$\Delta \in \mathbb{R}^{D \times T}, \text{ low-rank}$$

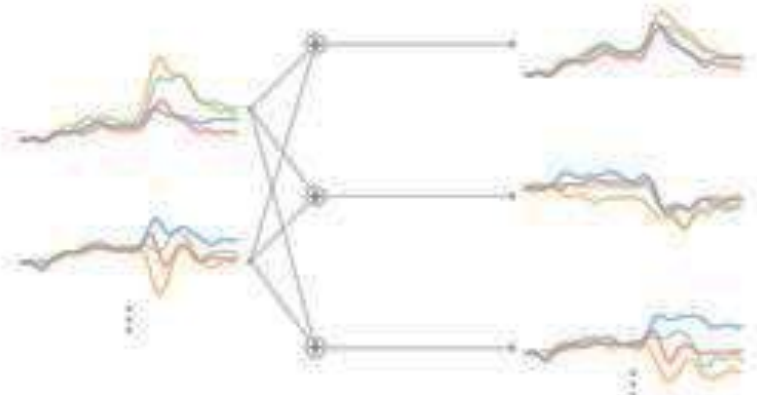
stim rates

$$y_t = g(C\tilde{z}_t + \Delta_t + d)$$

w/ rectifying non-linearity to correctly handle silencing

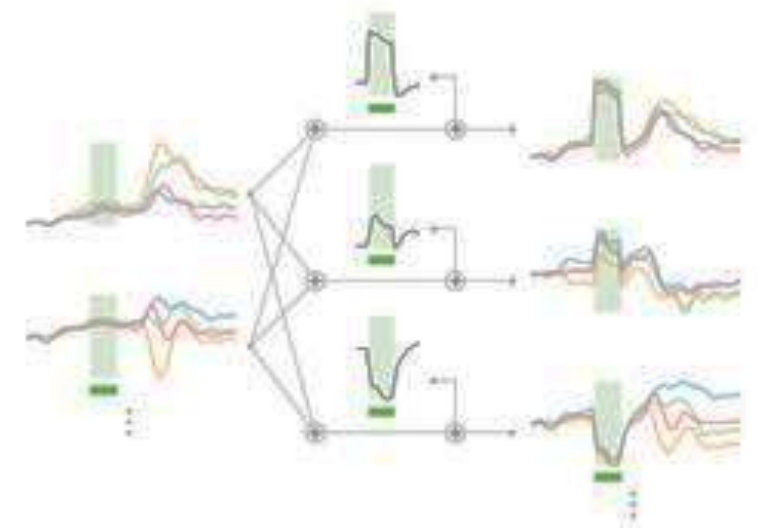
g Latent Variable Model

non-stim latents non-stim firing rates
latent dimensions neurons (PSTHs)



stim latents stim deltas stim rates

independent of non-stim condition-independent, low-rank neurons X time

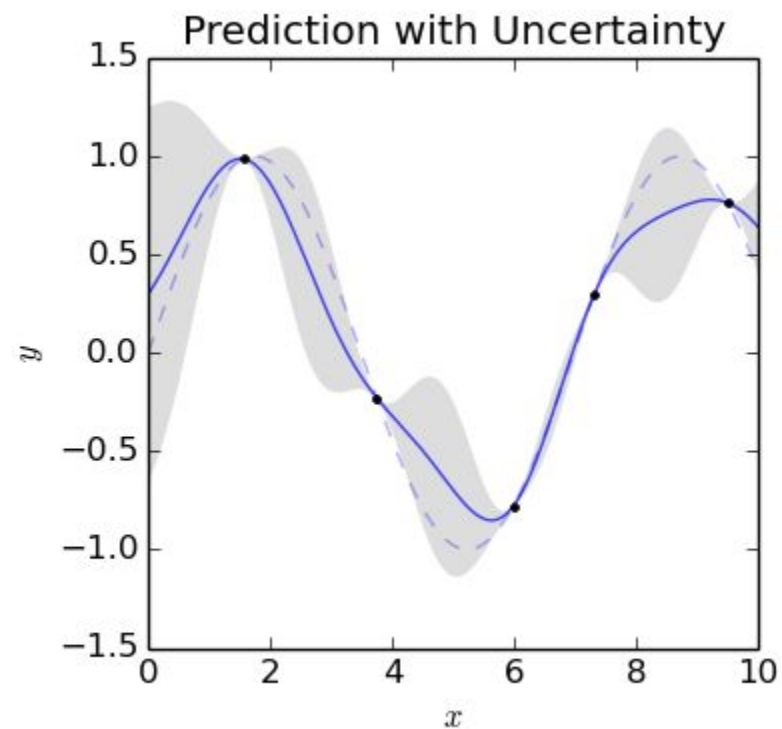
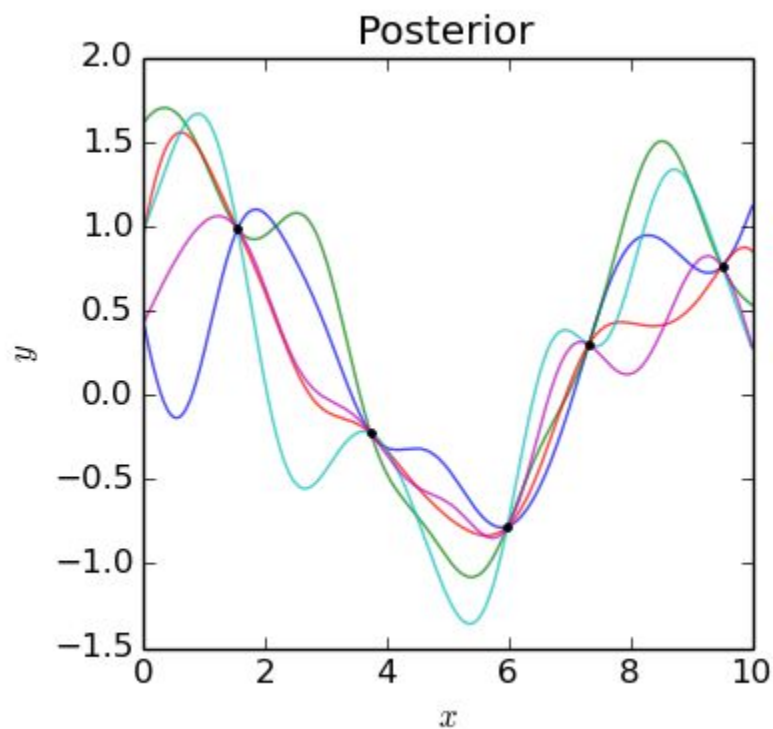
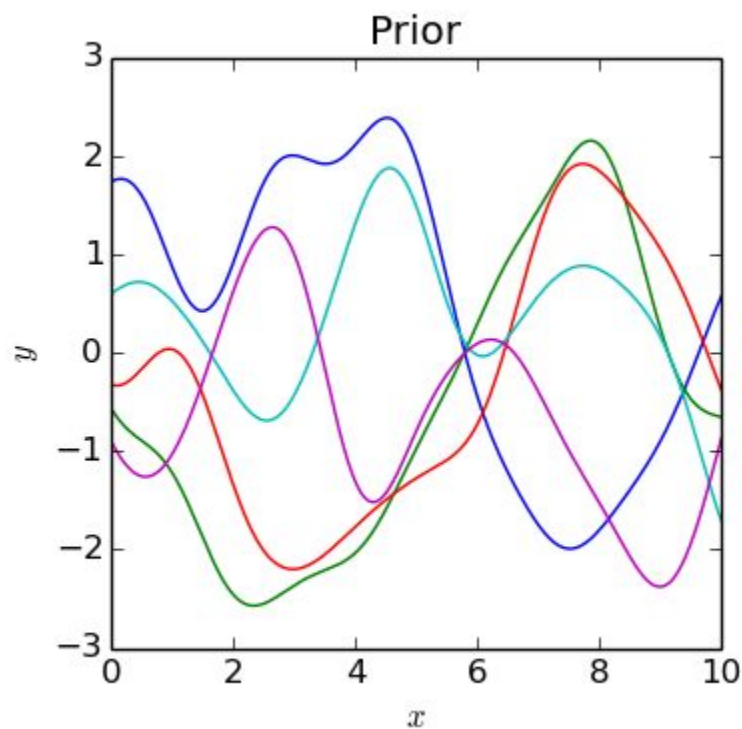


Gaussian process autocorrelation timescale =

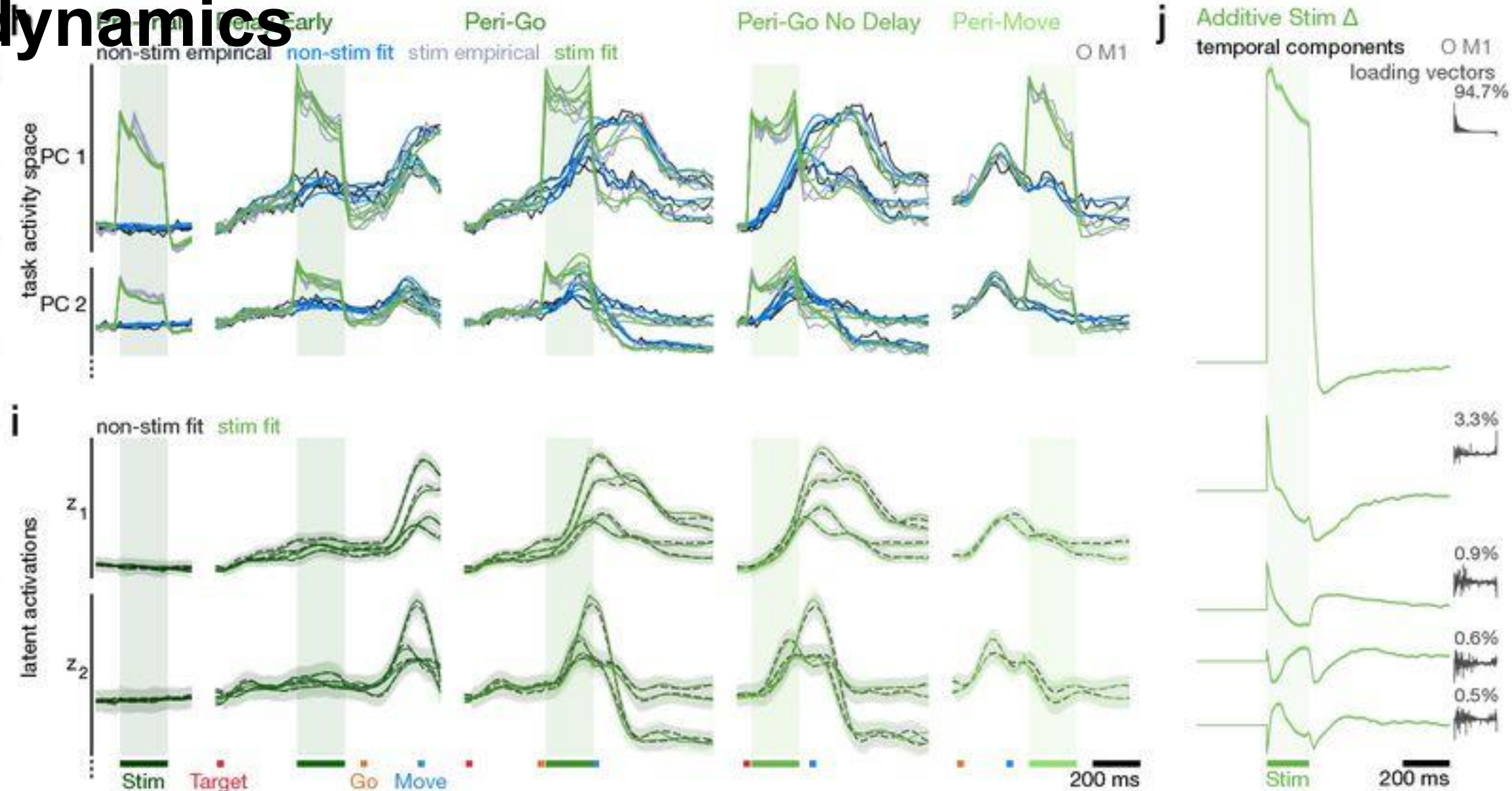
100ms

Additive component rank = 5

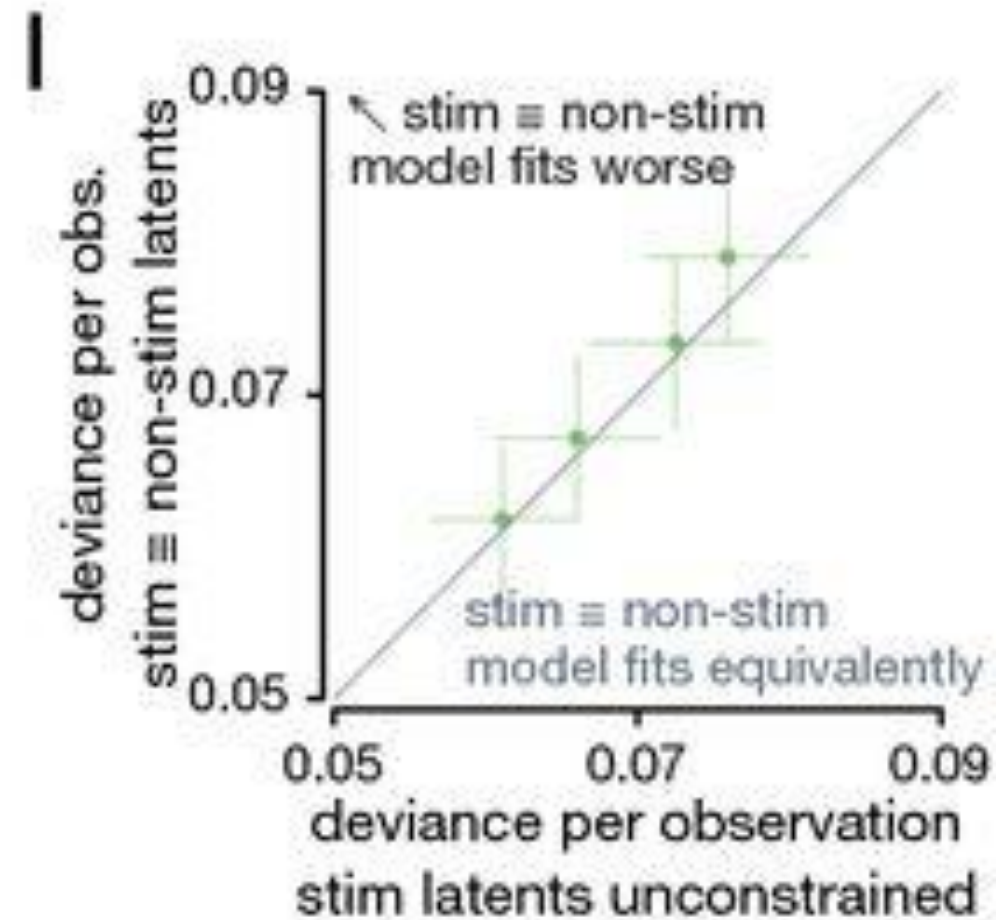
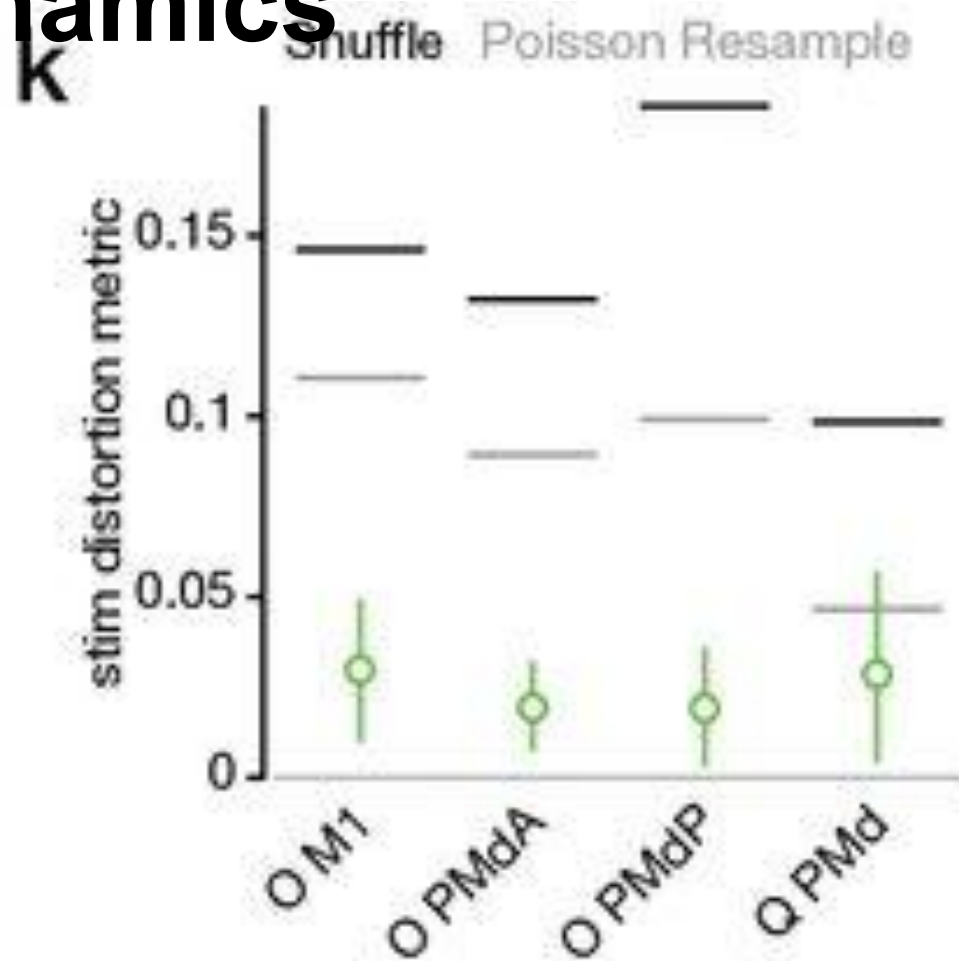
Gaussian what-now?



Perturbation doesn't affect underlying dynamics

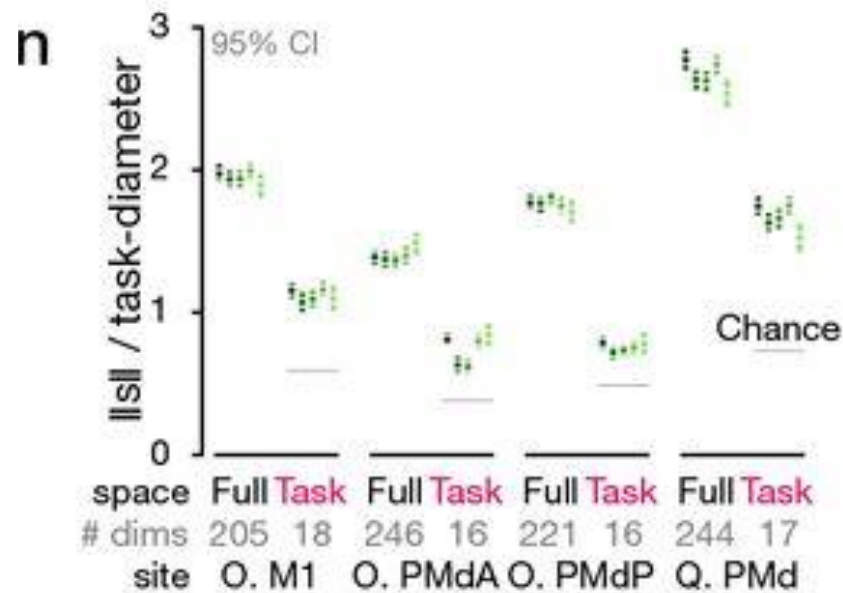
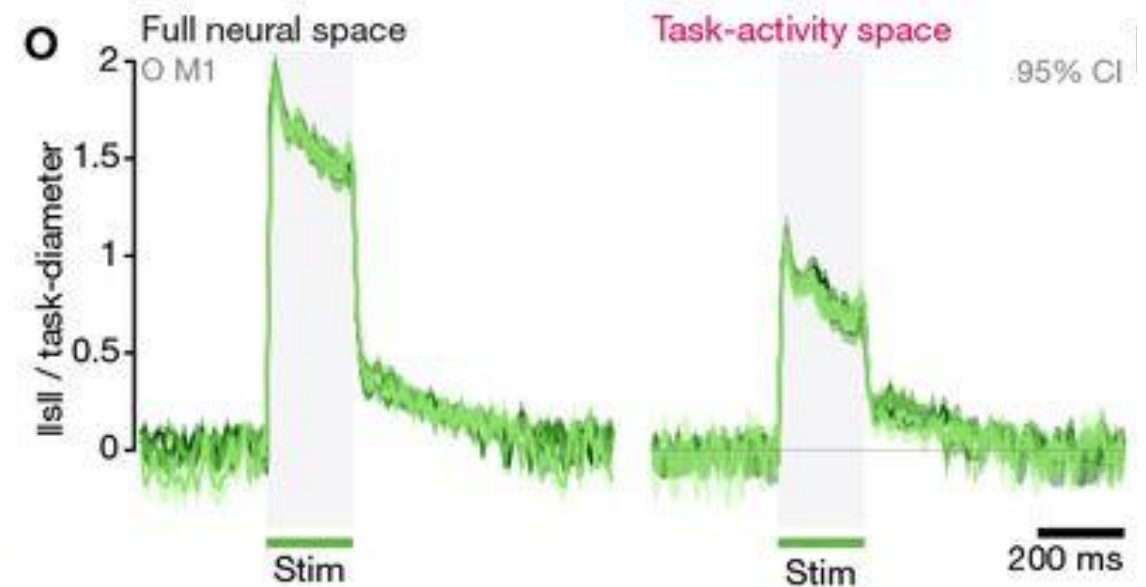
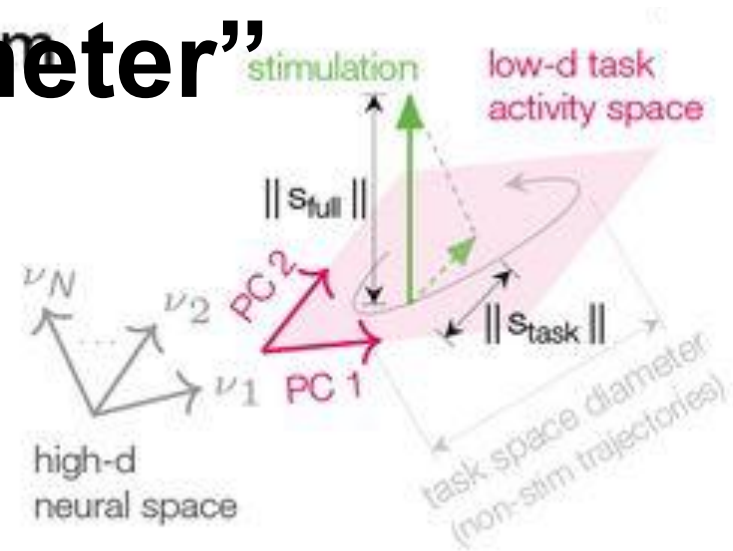


Perturbation doesn't affect underlying dynamics



Stim-distortion: Euclidean distance after alignment using rotation, reflection, translation (no scaling)

Introducing “Task-Space Diameter”



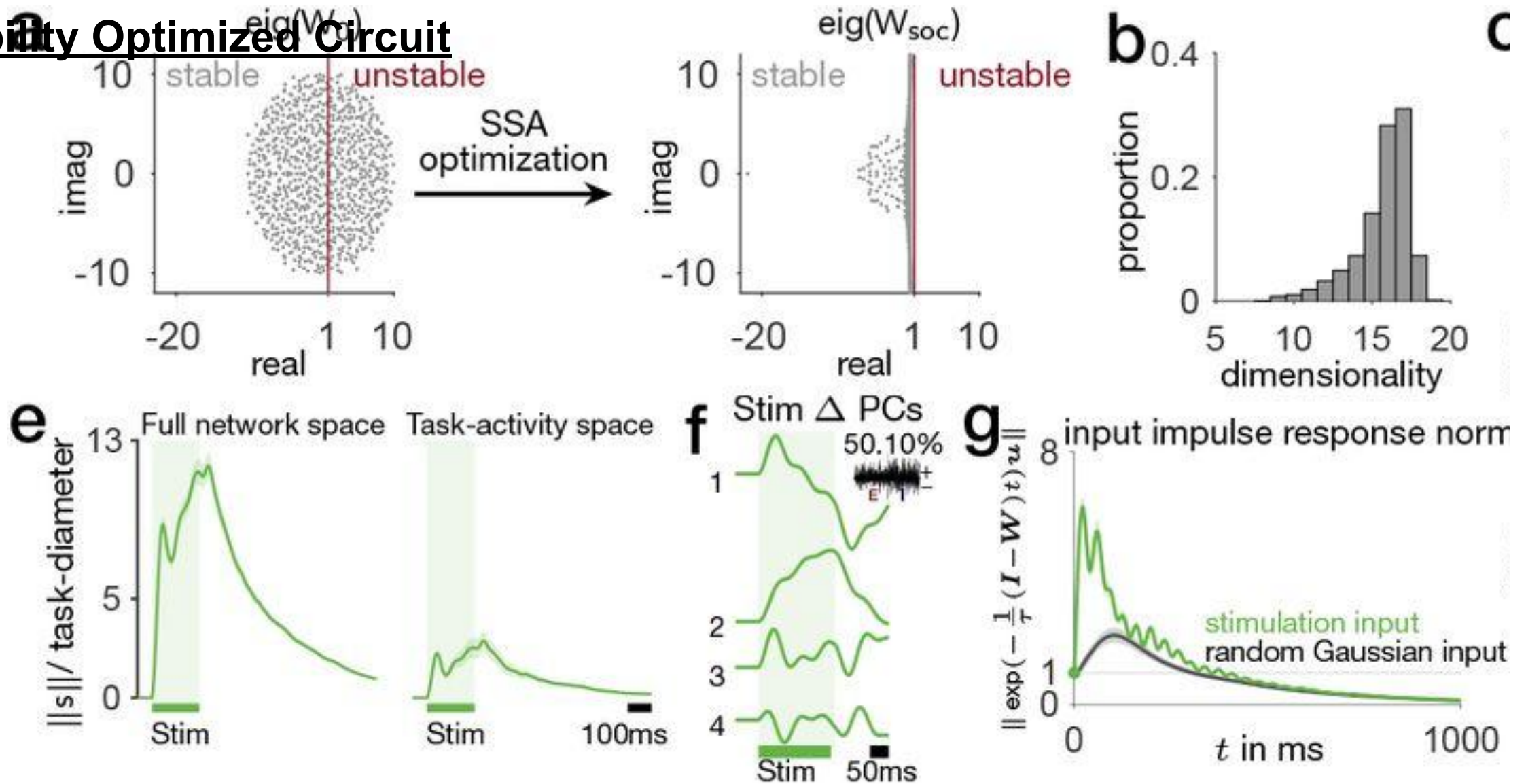
Magnitude of perturbation vector projected in “task subspace”

Task Activity Space: Dimensions with 95% of variance using principal components

Moving on to models...

Slow perturbation decay in a type of Reservoir Network

Stability Optimized Circuit



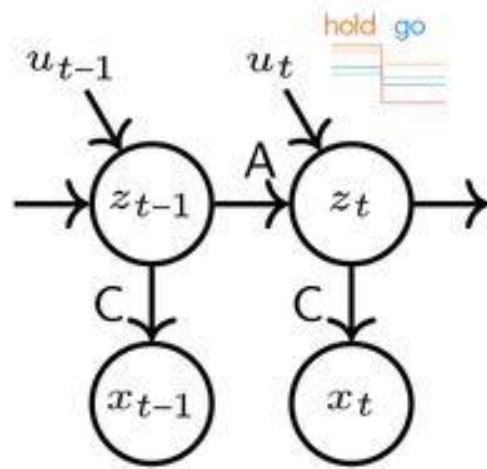
Done with H1?

- These predicted responses were qualitatively different from the observed neural responses arguing against the presence of high-dimensional reservoir dynamics in the motor cortex (**H1**).
- BUT...
- Following [14], we set $\tau = 200$ to match the time-scales observed in motor cortical activity patterns.

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x}(t) + W\mathbf{x}(t) + \mathbf{I}(t)$$

Modelling perturbation of low-D dynamics

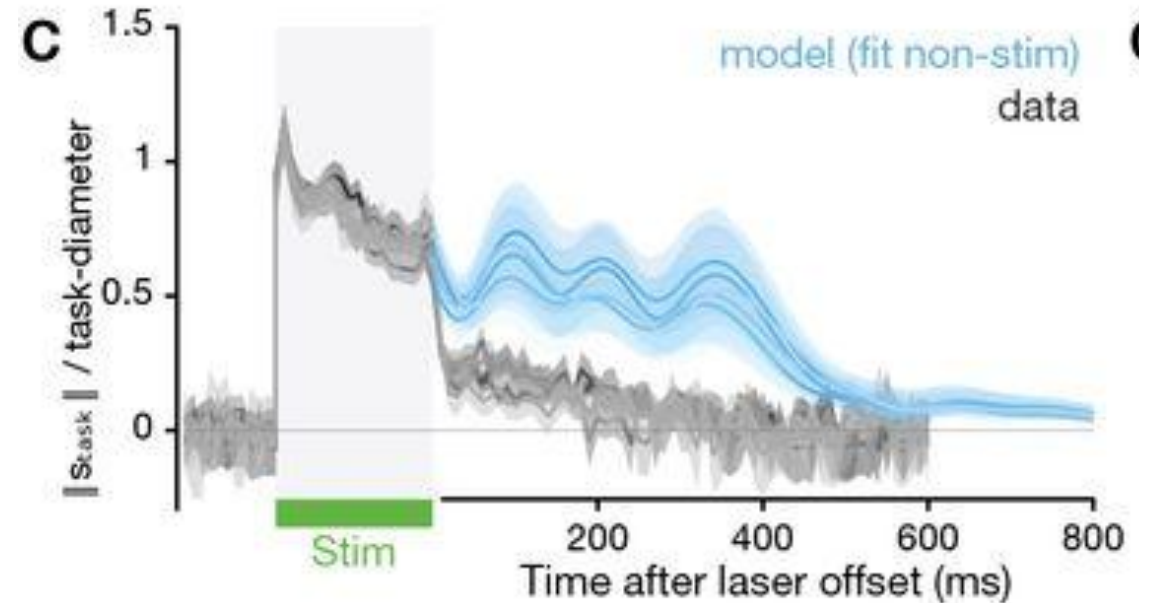
a Latent LDS model



condition-specific input
piecewise-constant, $u_t \in \mathbb{R}^K$

low-d latent dynamics
 $z_t = Az_{t-1} + u_t \in \mathbb{R}^K$

neuron rates
 $x_t = Cz_t + d \in \mathbb{R}^D$



Rationalizing LDS perturbation results

- Task activity space = Task dynamics space?
- Might not be perturbing the “task dynamics space”
- Can’t selectively perturb LDS, therefore...
- Fit model where we know “task dynamics space” and perturb that model.

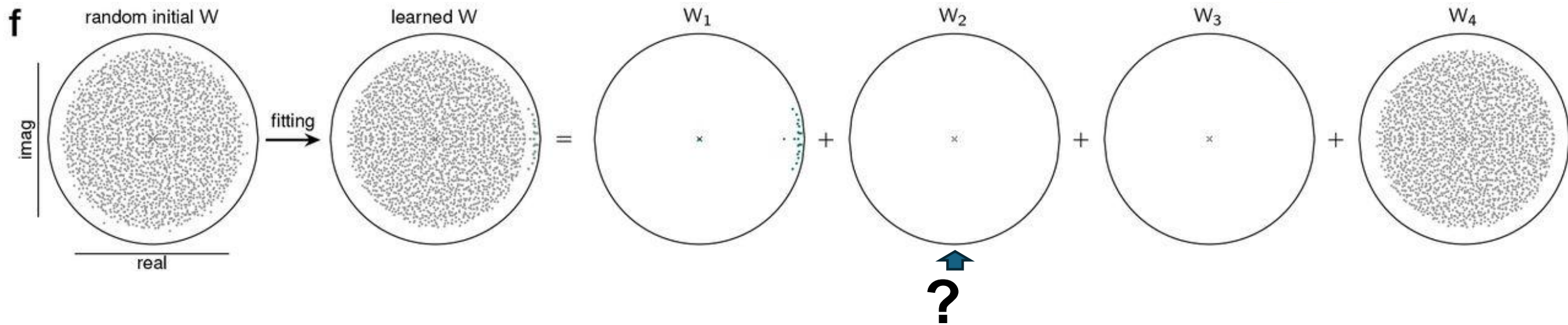
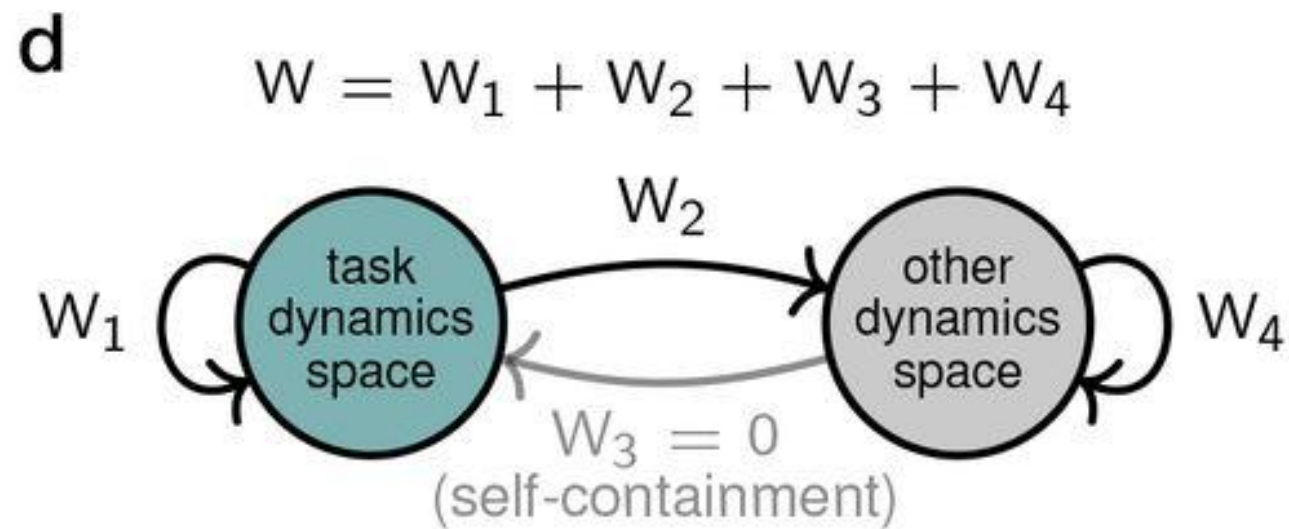
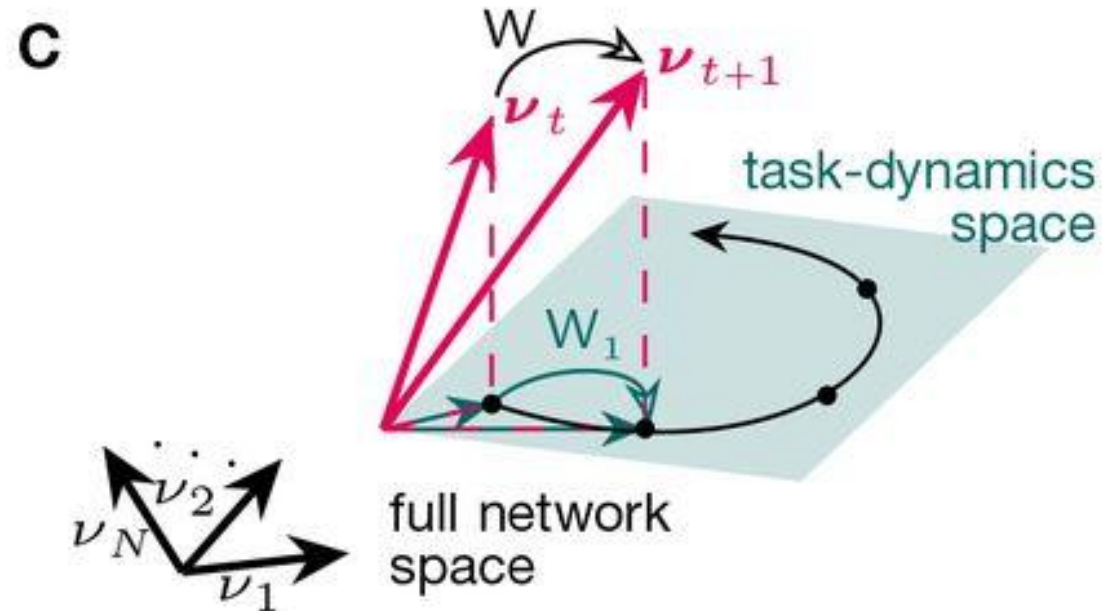
Mechanistic Explanation of perturbation results

More realistic model of subspace-structured

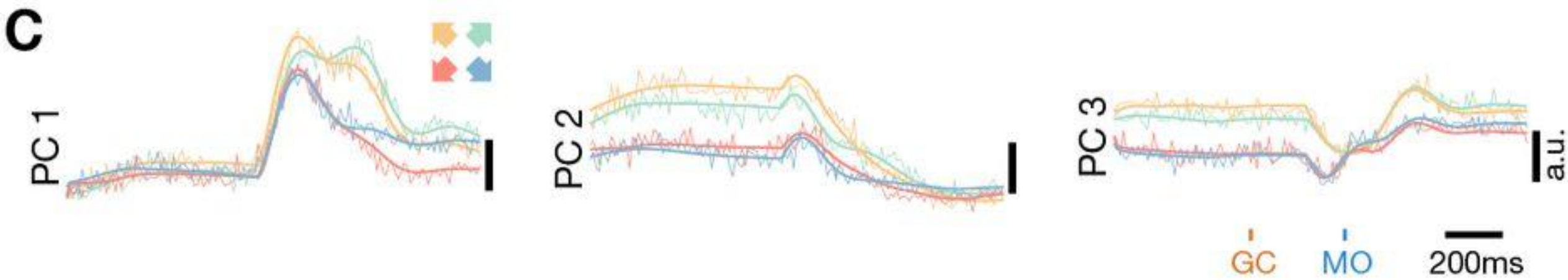
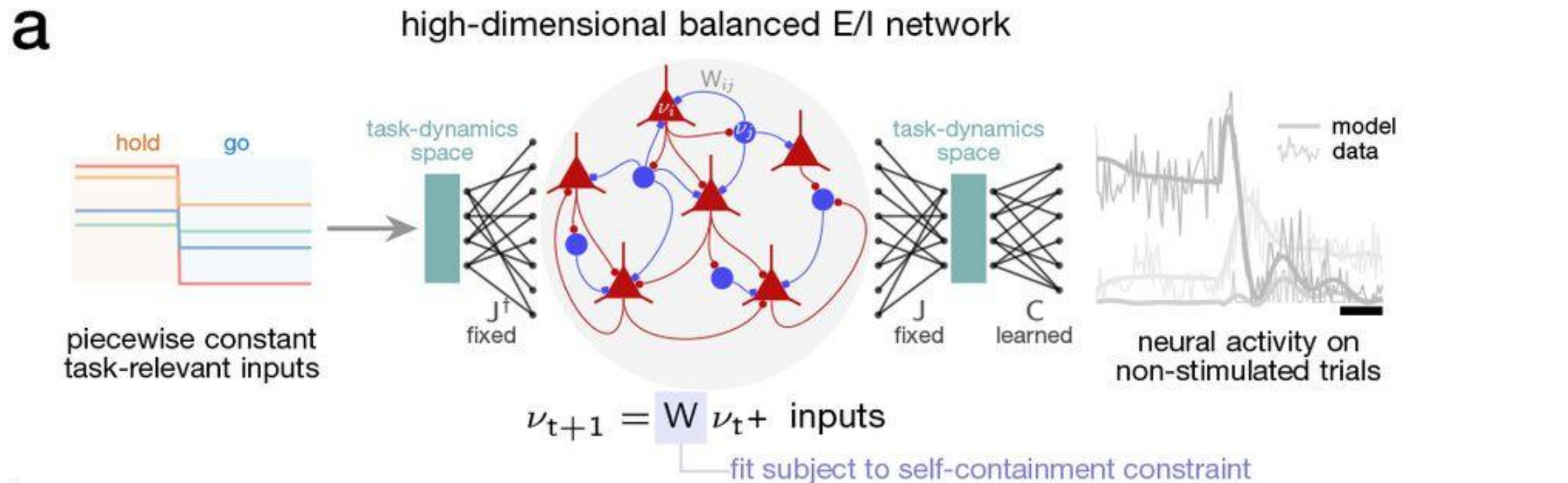
dynamics

- Connectivity satisfies: Dale's law⁵³ and E/I balance.⁵³⁻⁵⁵
- high-dimensional ($N=1500-2500$ units), linear network
- low-dimensional subspace of size $K=N/100$ driving response
- We then optimized the E/I network connectivity matrix to maximize the likelihood that a linear readout from this task dynamics space matched the measured neural responses.
- **The E/I model was fit exclusively to data from trials without stimulation.**

More realistic model of subspace-structured dynamics

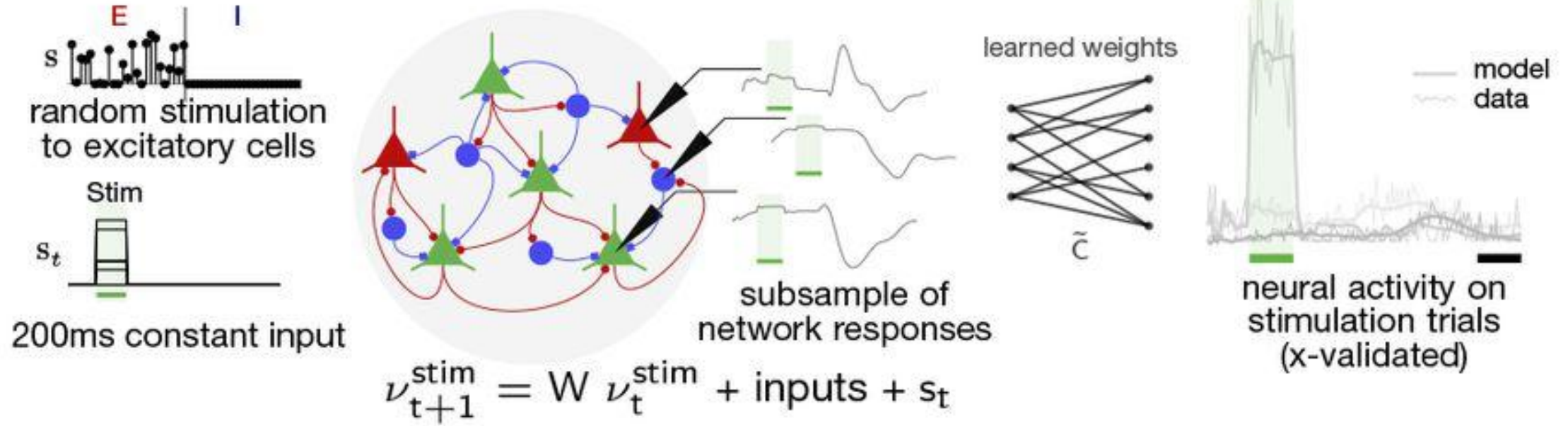


More realistic model of subspace-structured dynamics

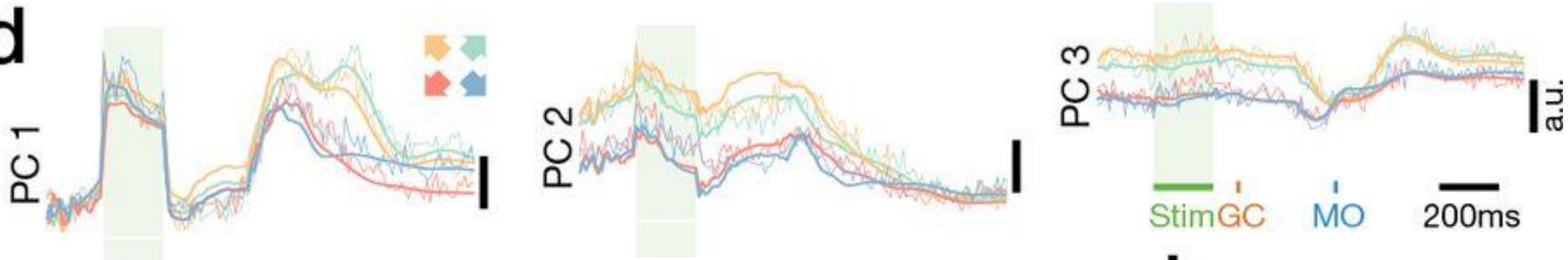


More realistic model of subspace-structured dynamics

b

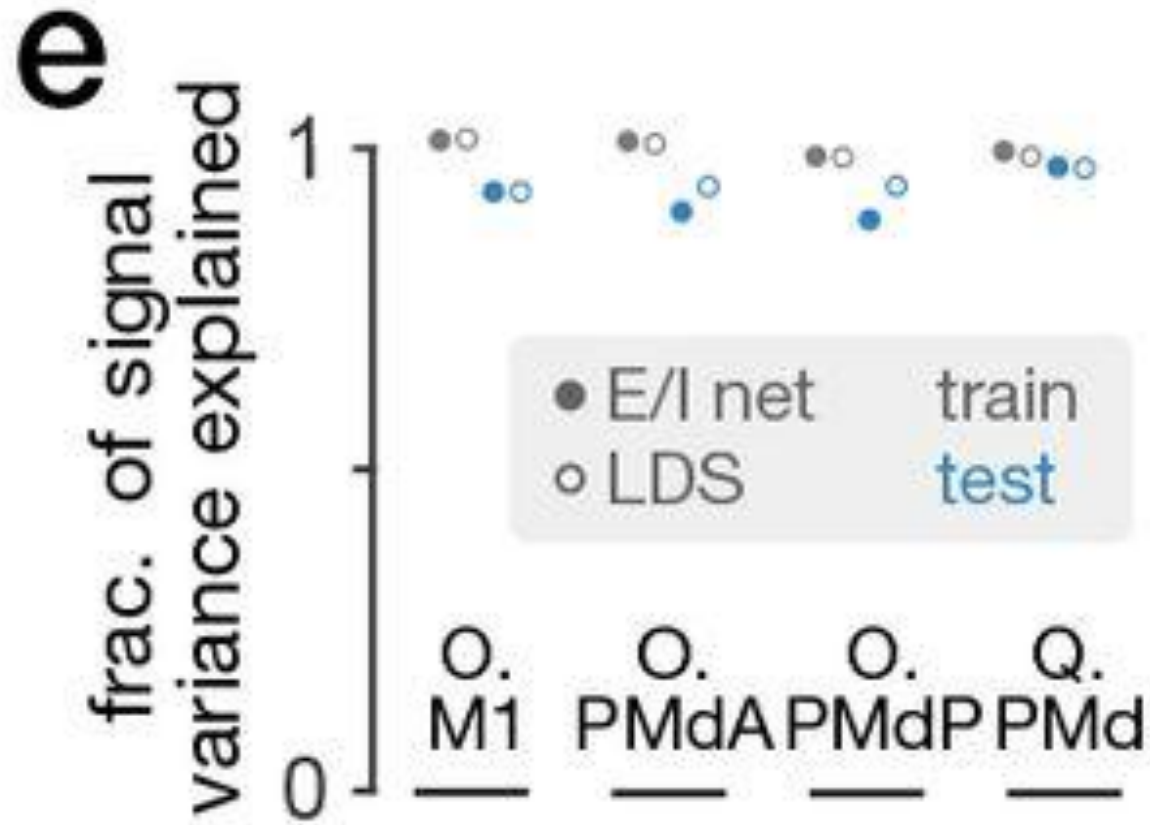


d



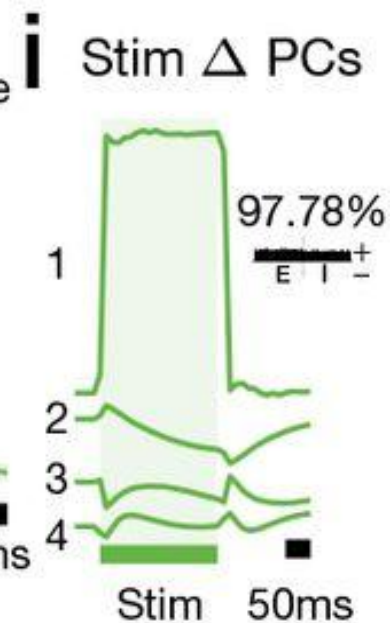
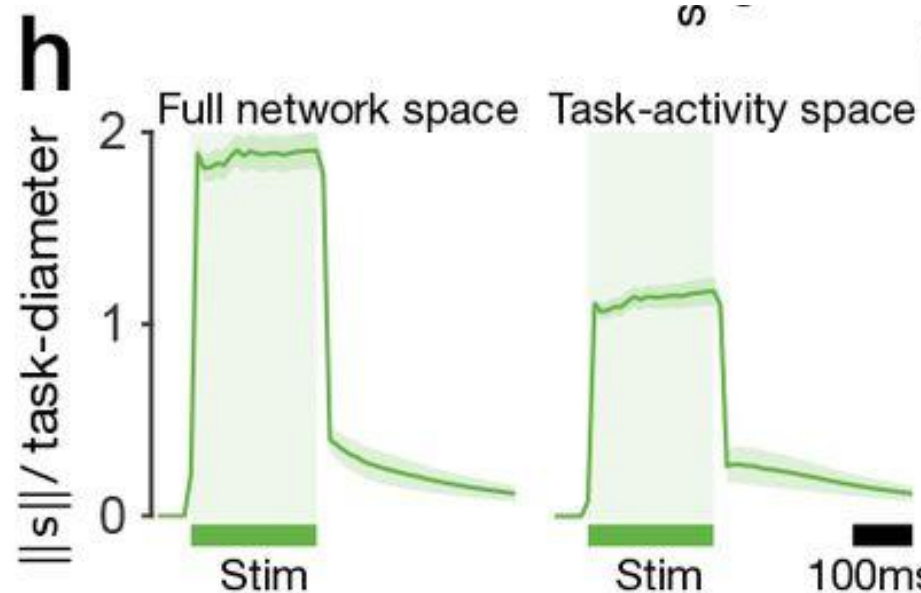
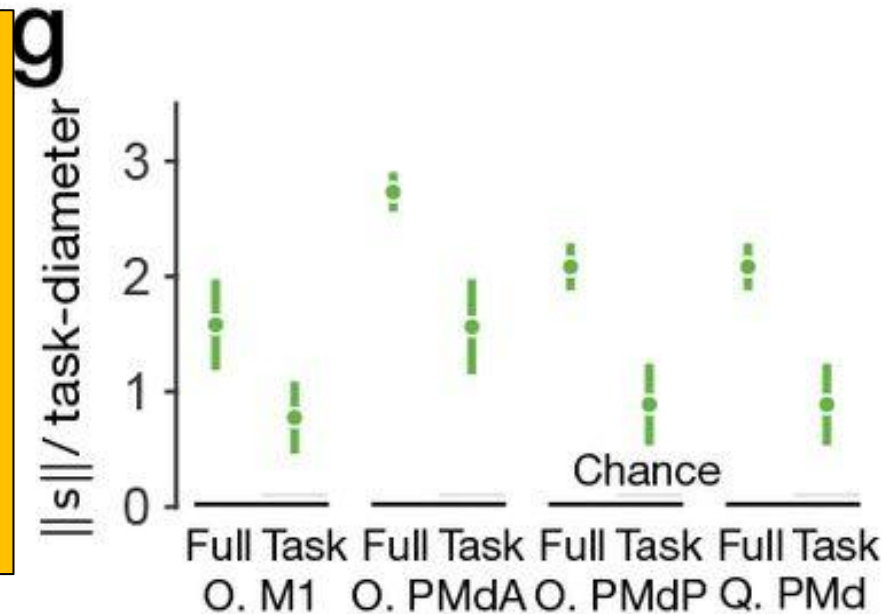
optogenetic stimulation was modeled with additive noisy, positive input patterns targeting random subsets of excitatory cells in the E/I network

More realistic model of subspace-structured dynamics

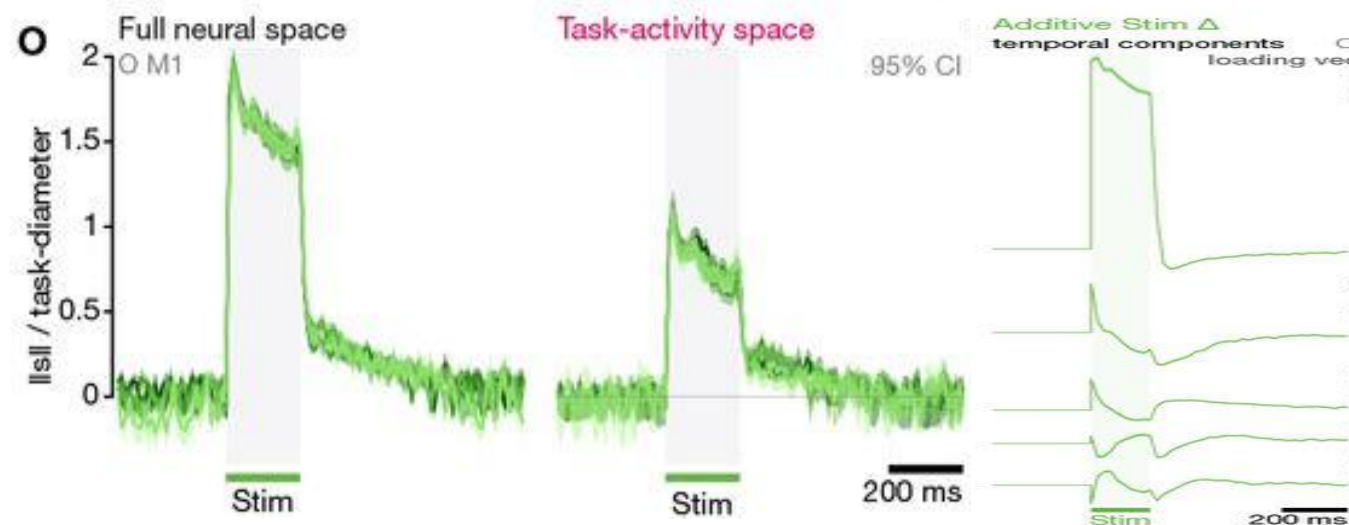
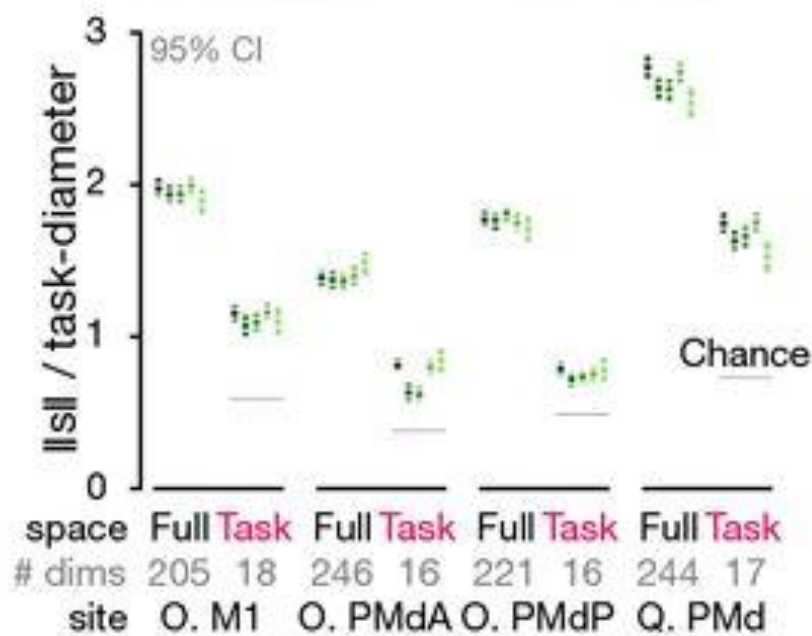


More realistic model of subspace-structured dynamics

Model

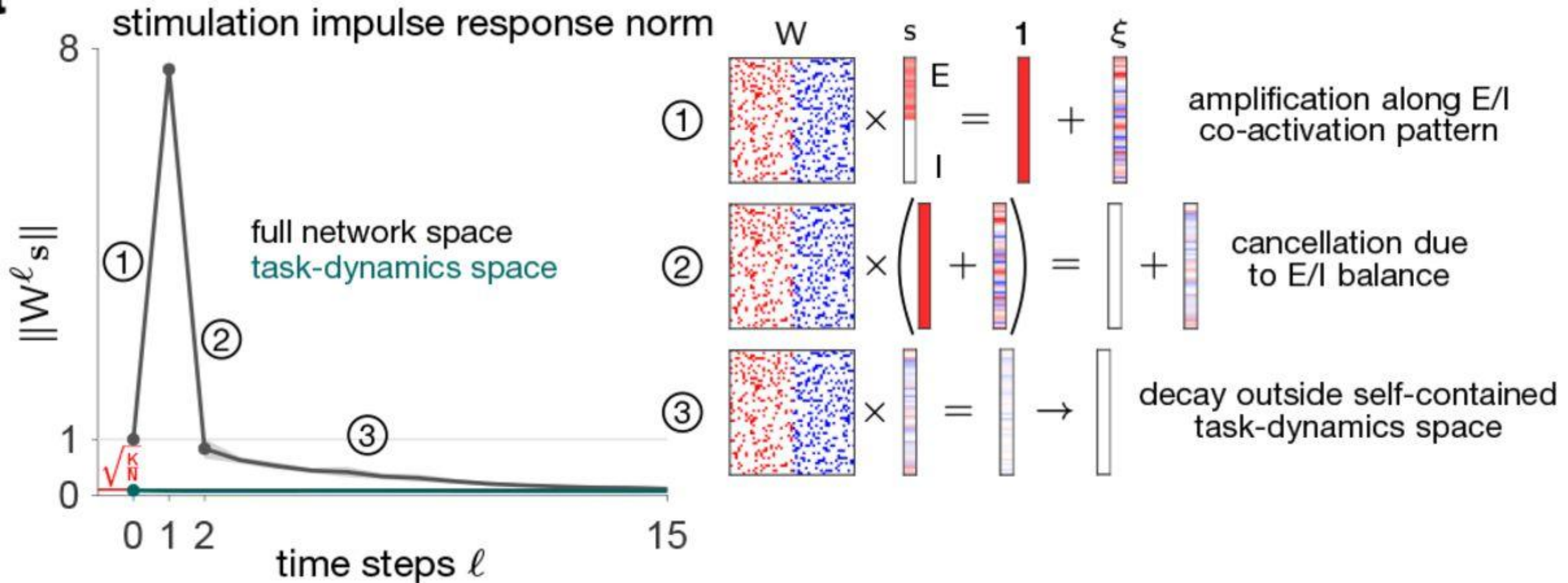


Data

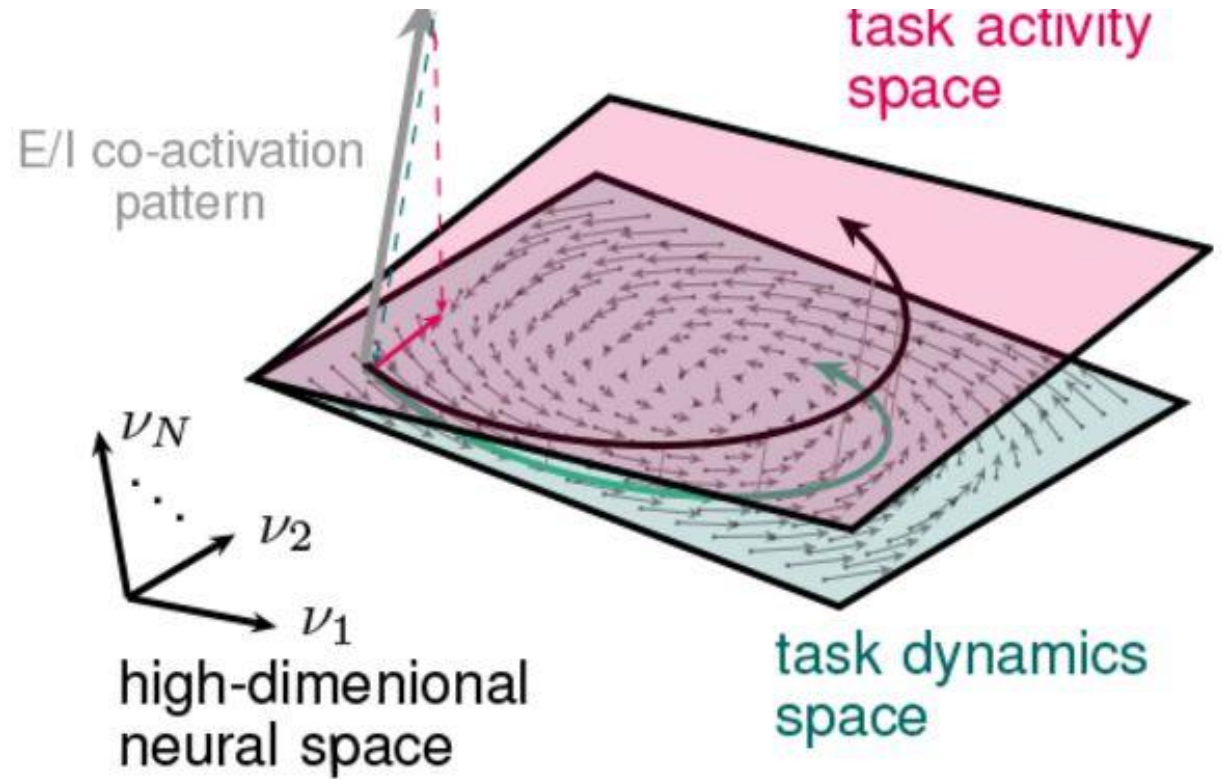
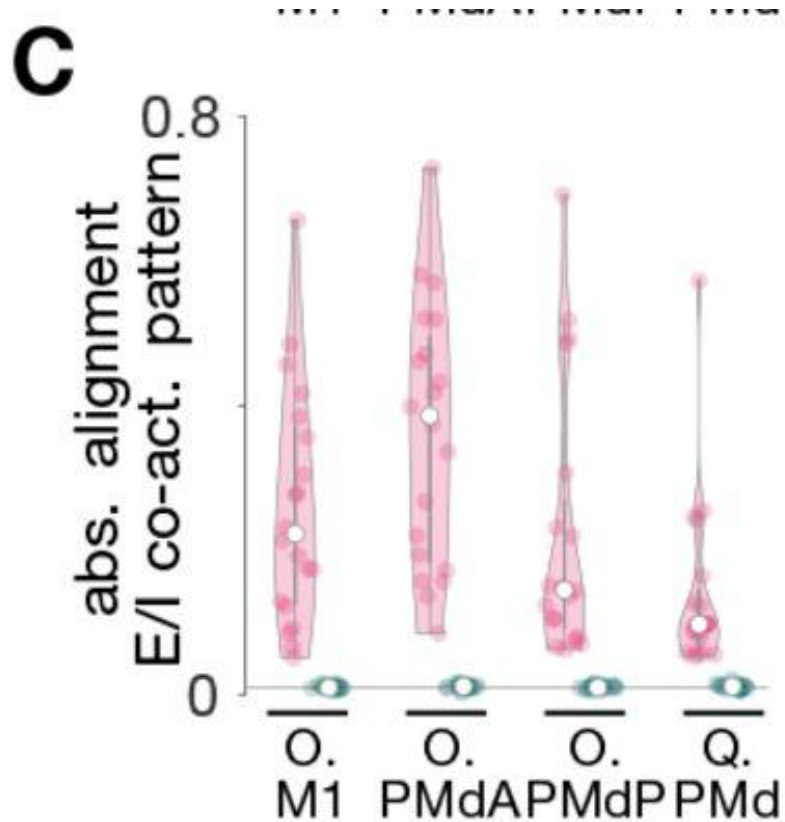


Mechanistic explanation for perturbation results

a



Mechanistic explanation for perturbation results



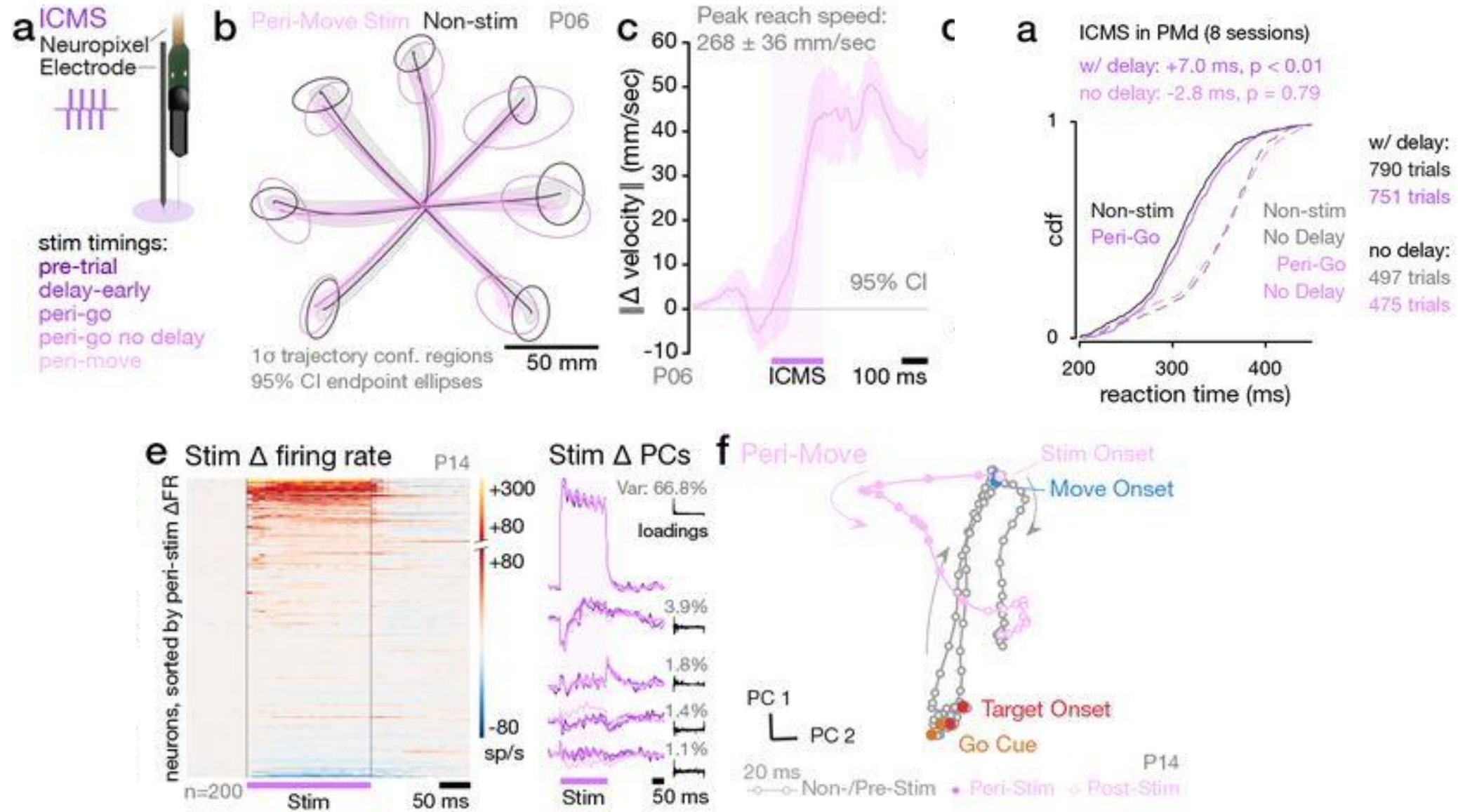
Checkpoint

- Path-following dynamics (H3): all local perturbations should decay rapidly back towards the externally imposed trajectory
- Subspace-structured dynamics (H2) : If perturbation doesn't reach task-dynamics subspace...perturbations should also decay rapidly.
- Therefore, **H2** and **H3** cannot be distinguished based on the results from these optogenetic perturbations alone.

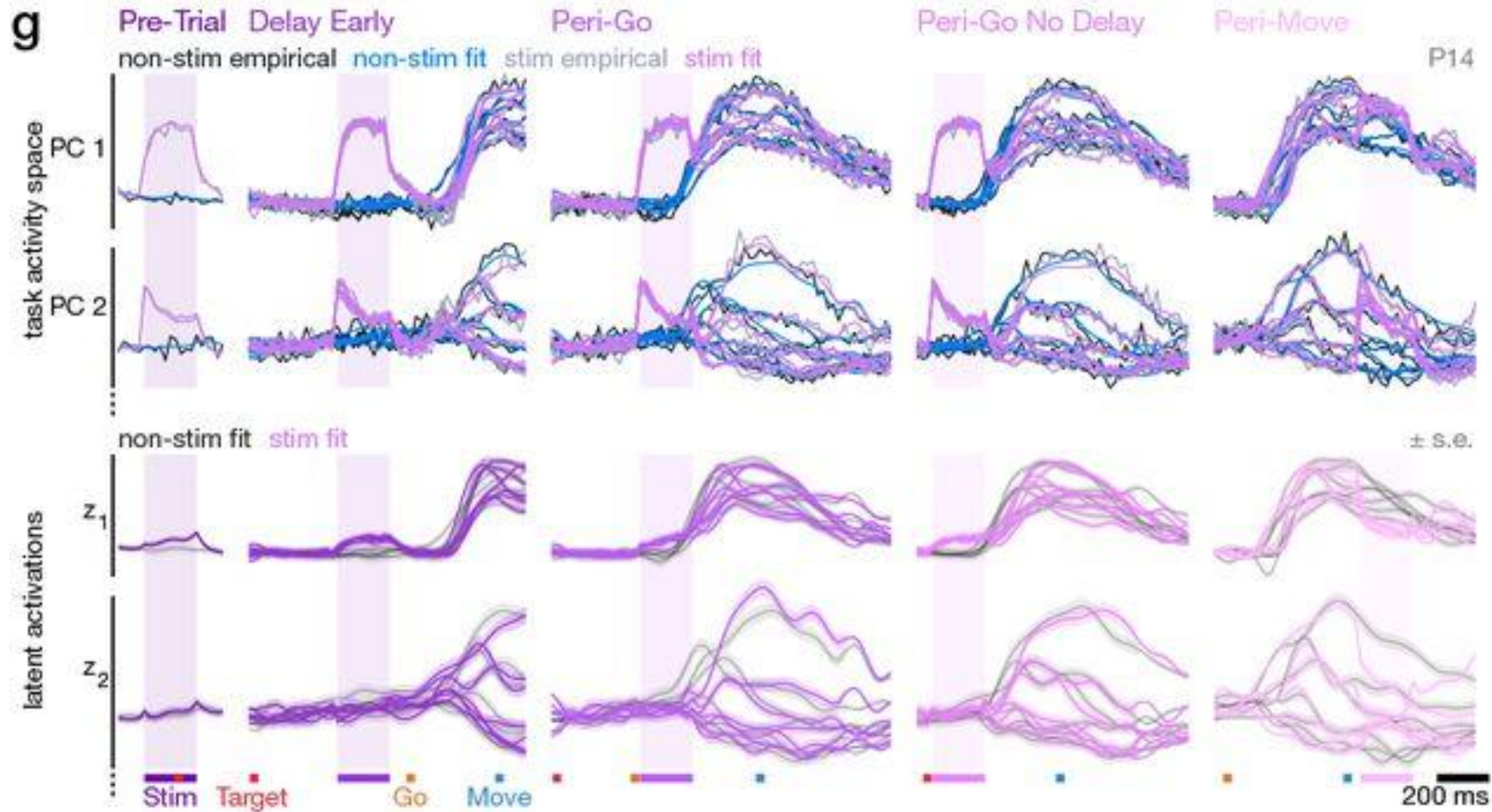
- Intro + Background
 - Hypotheses (total of 3)
 - Experimental Setup
- Results (knock out 1 hypothesis with each section):
 - Optogenetics
 - Intra-Cranial Microstimulation

- We reasoned that the ability to influence task-relevant behavioral outcomes might indicate that an experimental intervention is able to induce perturbations within the task dynamics space.
- Intracortical electrical microstimulation (ICMS), which is known to disrupt motor preparation⁶⁰ and to evoke movements readily when delivered to motor cortex.

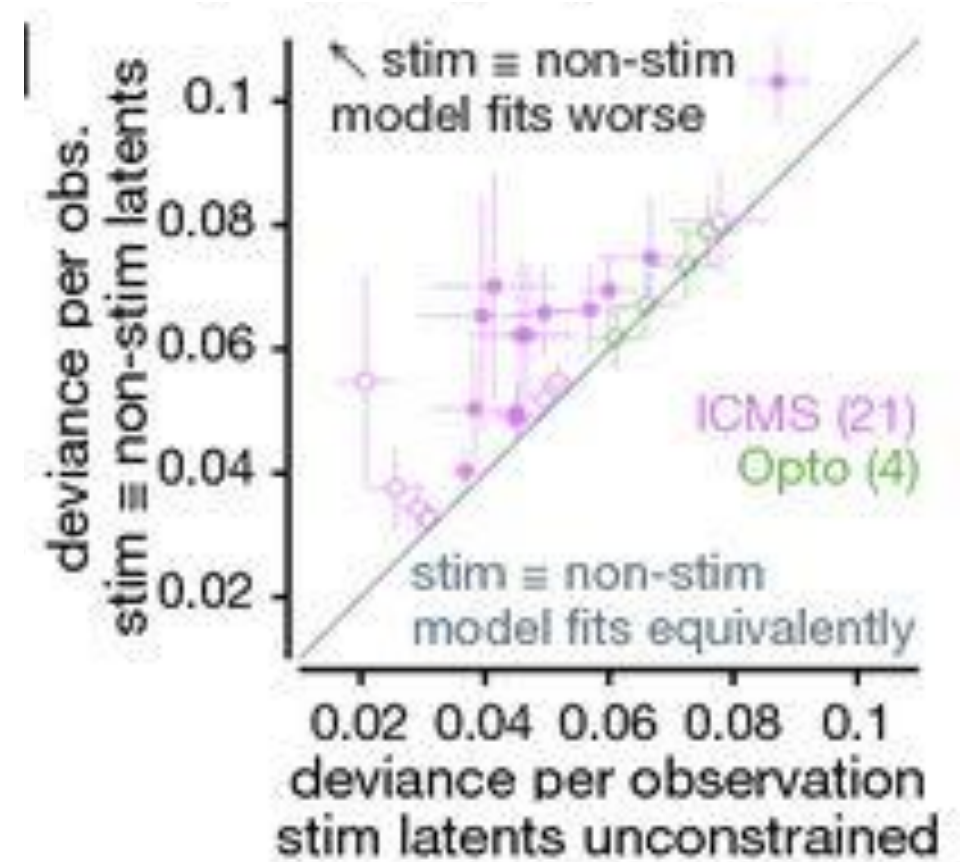
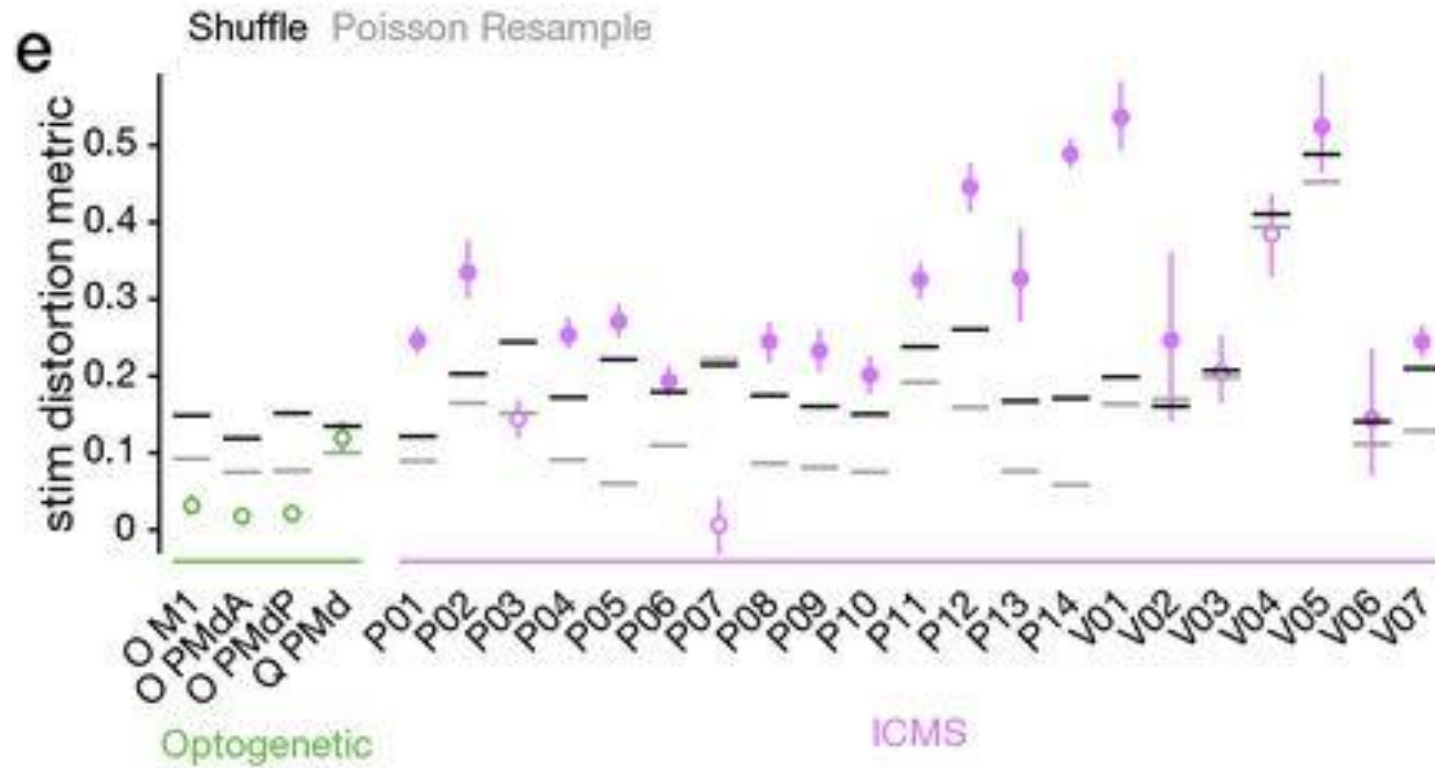
ICMS Effects



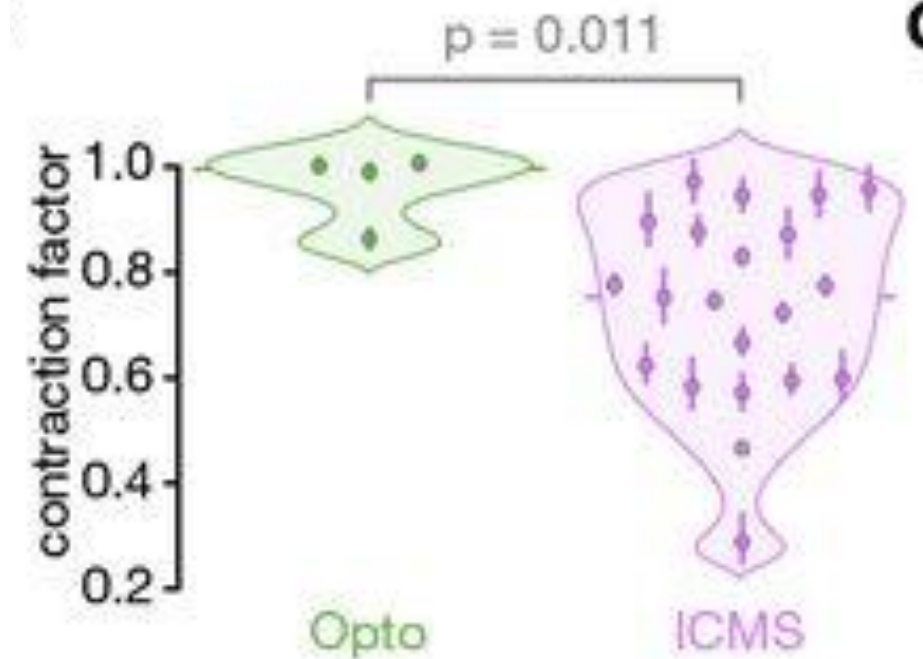
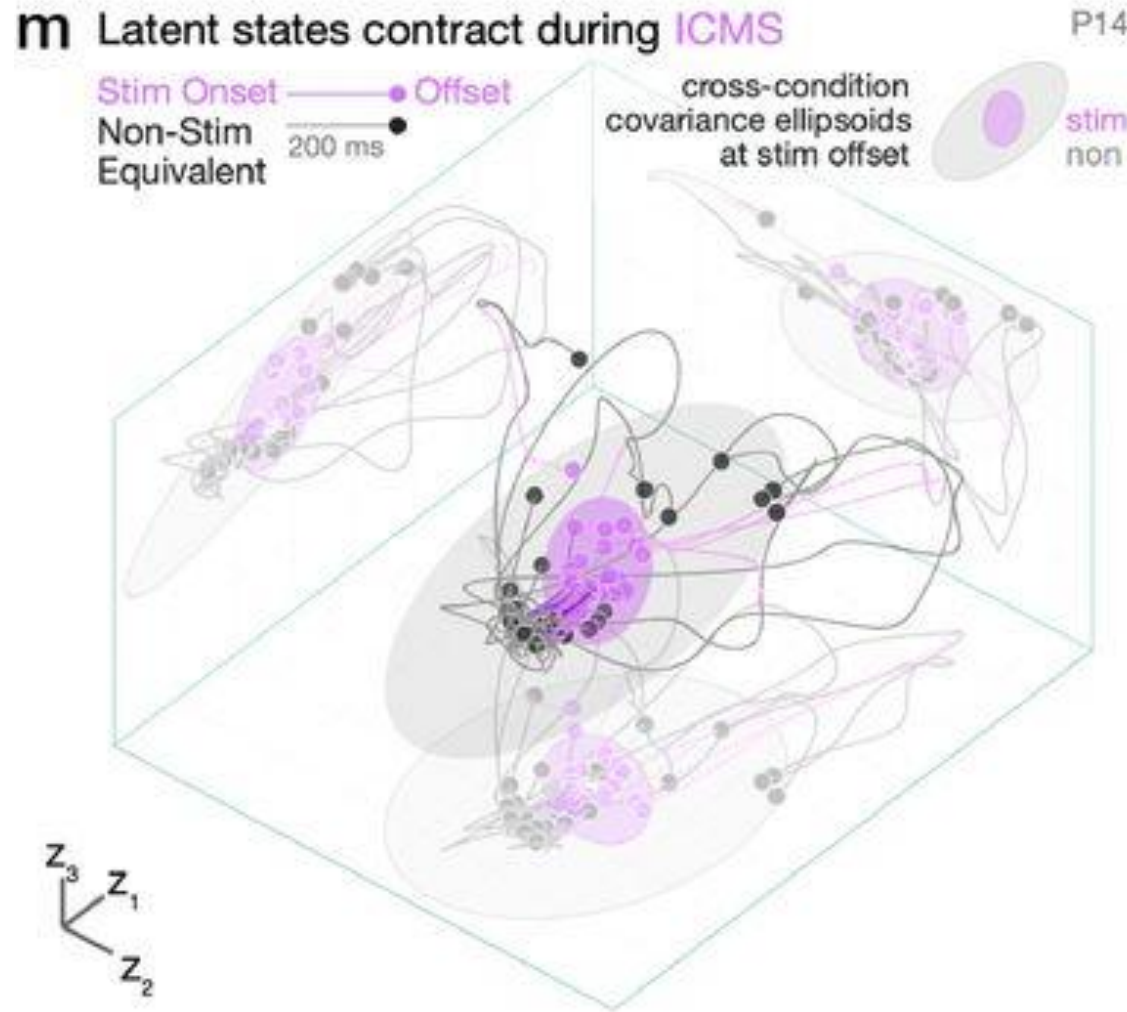
ICMS Effects



ICMS has stronger (actually significant) effect on neural dynamics



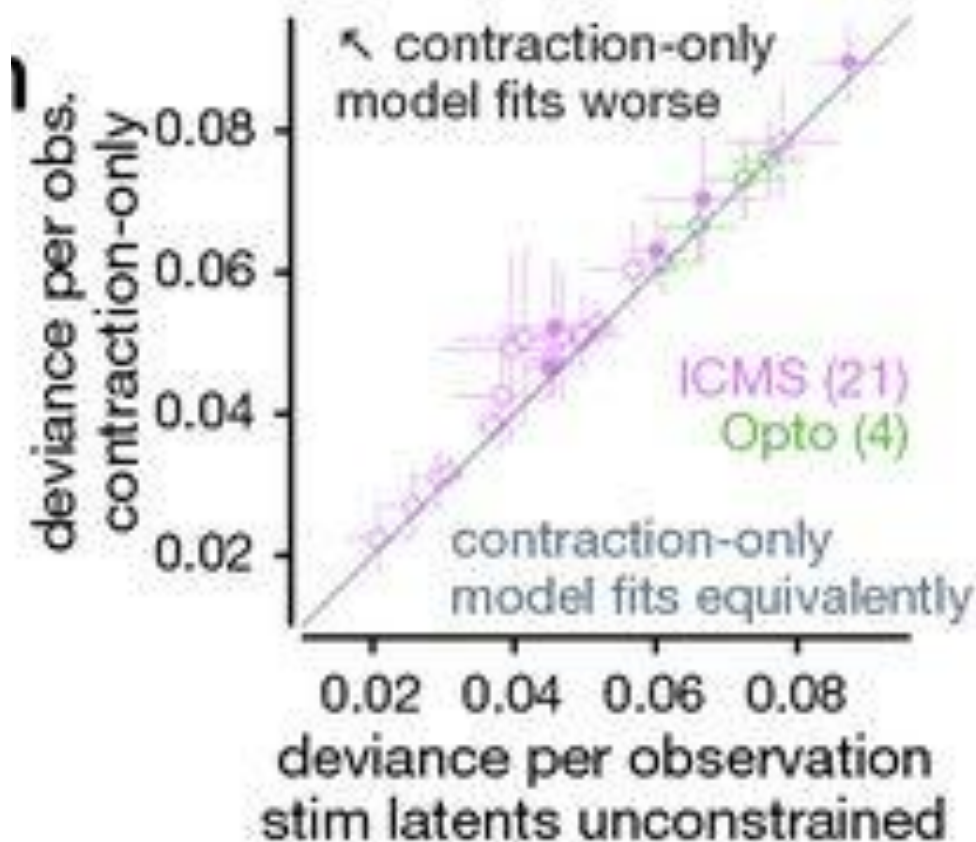
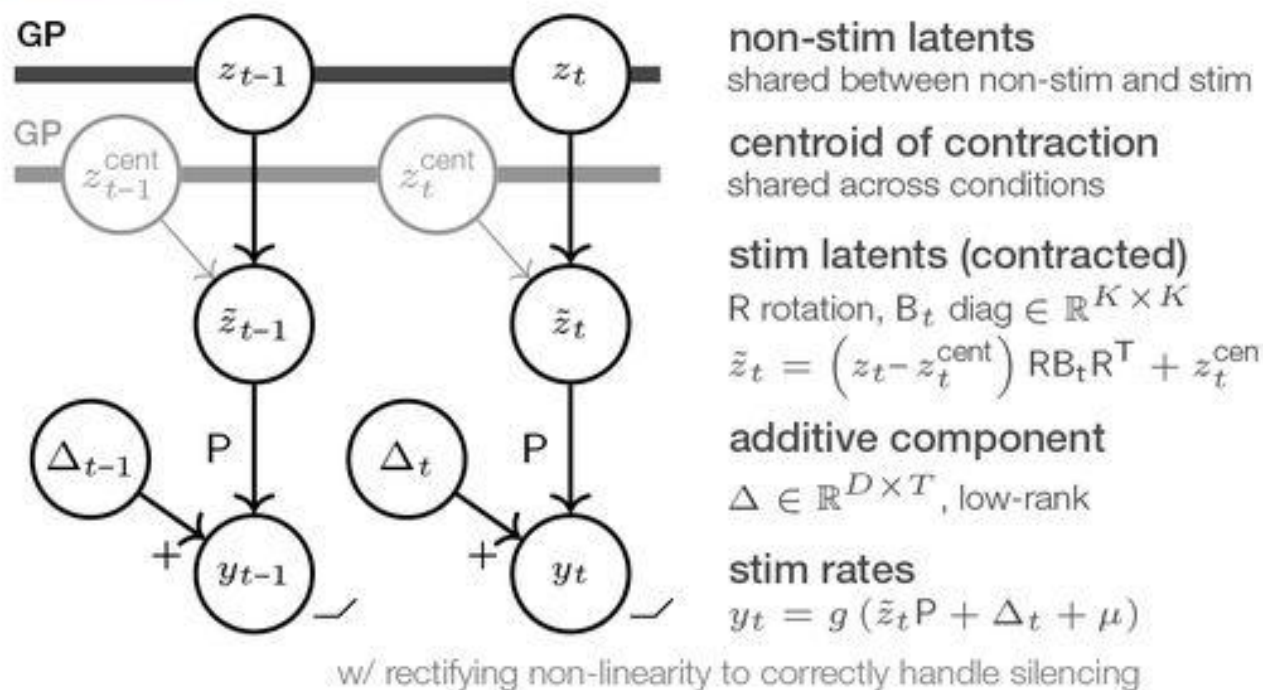
ICMS causes “contraction” of neural trajectory



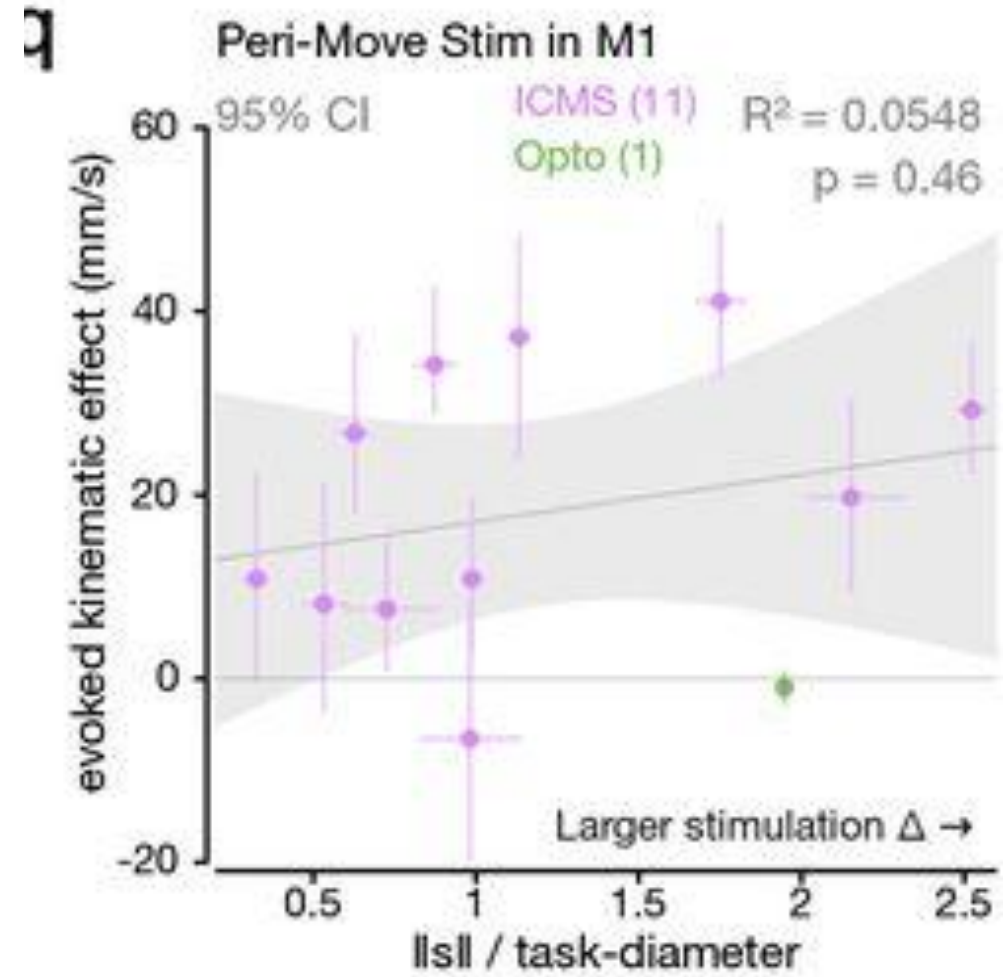
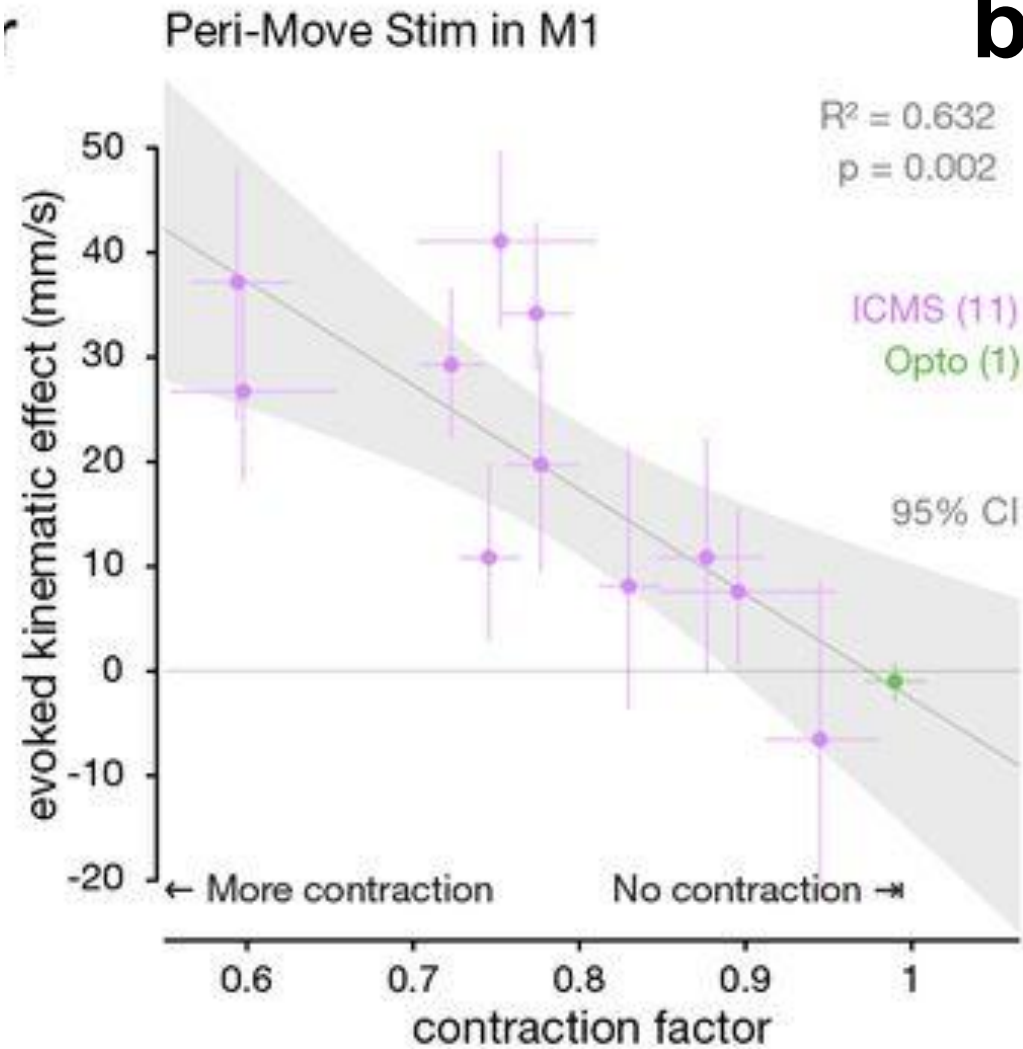
For each session, we computed a contraction factor, defined as the ratio of variance in the stimulated vs. non-stimulated neural states, with one corresponding to no contraction, and zero indicating complete contraction to a point.

Modelling ICMS-caused contraction

k Stimulation neural data, contraction-only model:

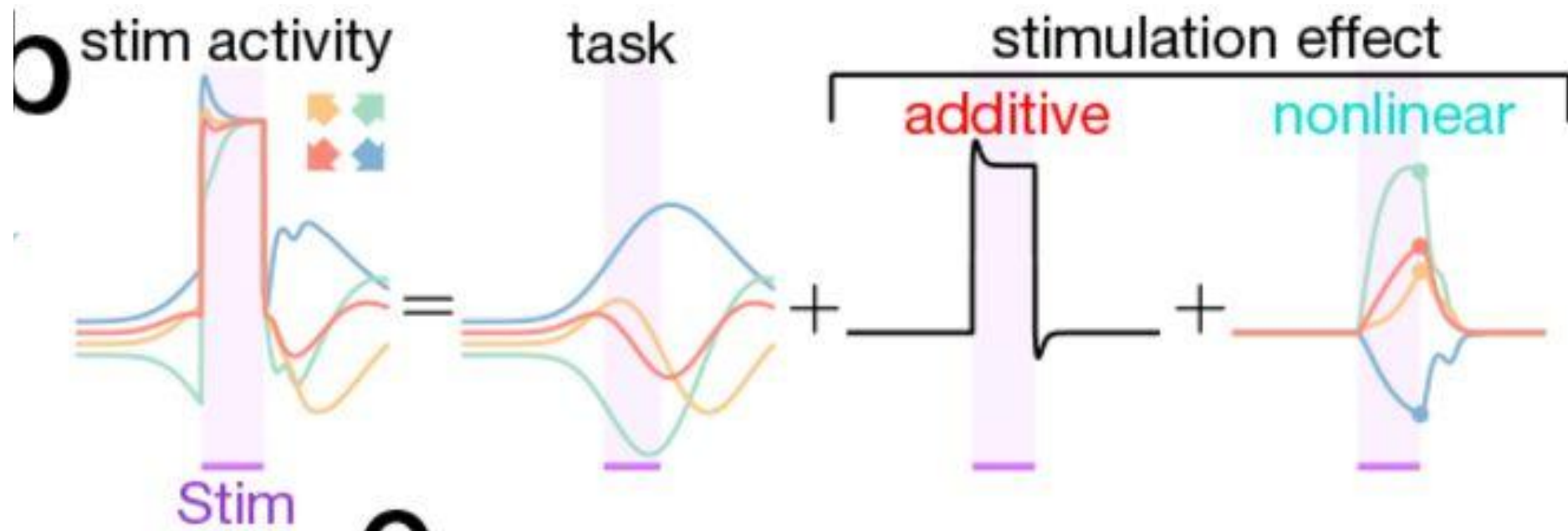
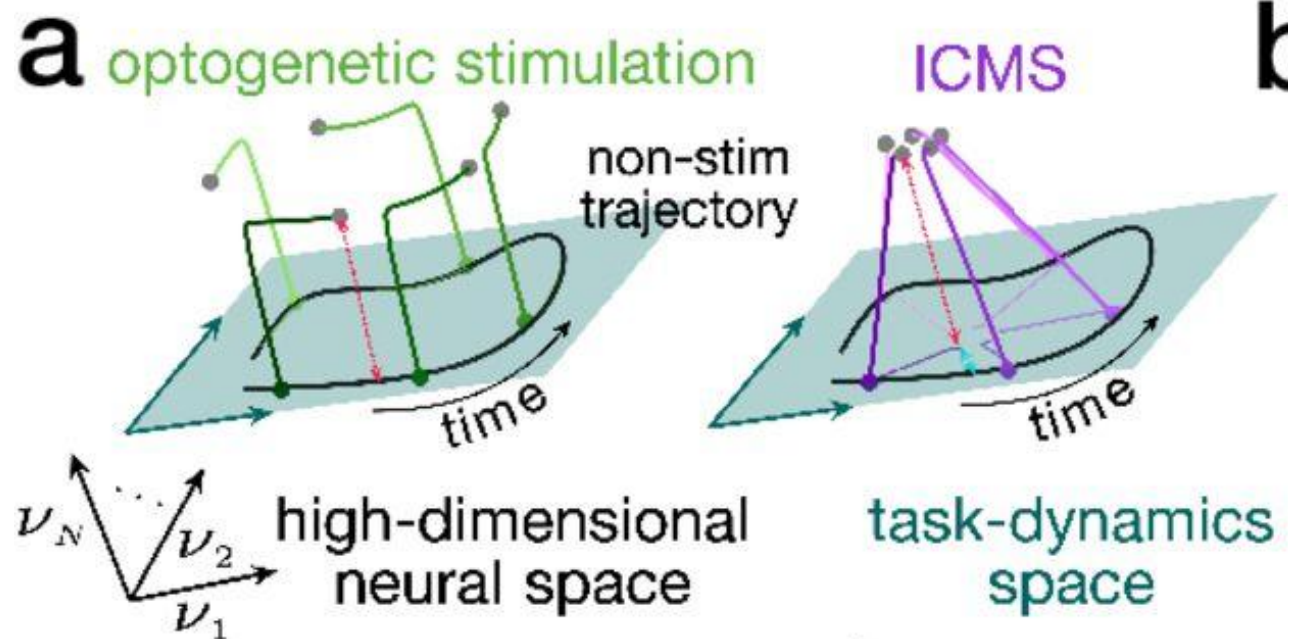


ICMS-caused dynamics contraction is related to behavior



Mechanistic model of ICMS-caused trajectory contraction

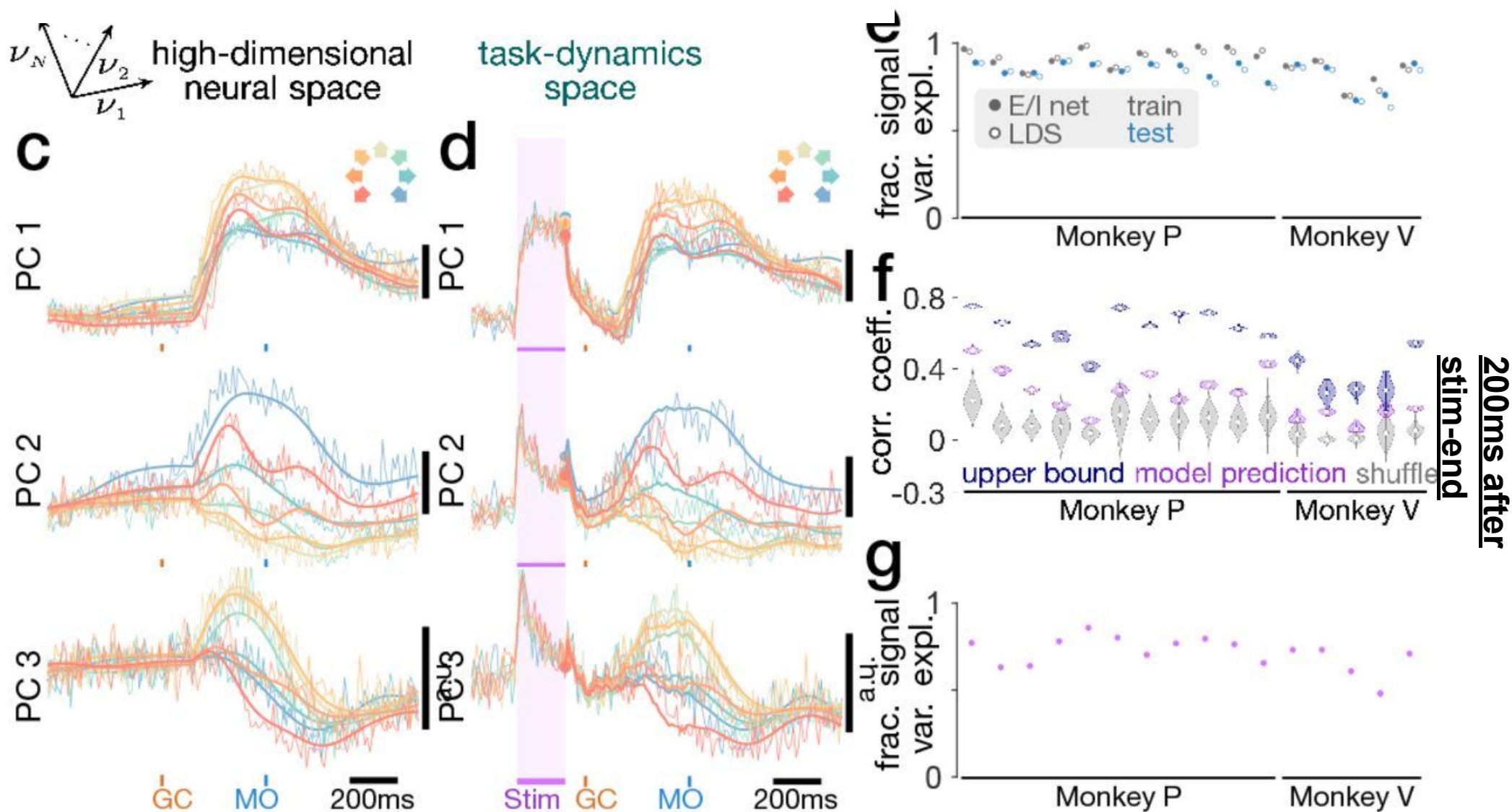
- The E/I network model was again fit exclusively to neural activity recorded in the absence of stimulation.
- The additive component of the stimulation effect was modeled as the fitted E/I network's response to a constant, additive input vector. The elements of this vector were drawn at random from a Gaussian distribution, since ICMS was not targeted to specific cell-types.
- The network responses to the additive perturbation were unable to capture the condition-dependent, nonlinear component of the stimulation responses.



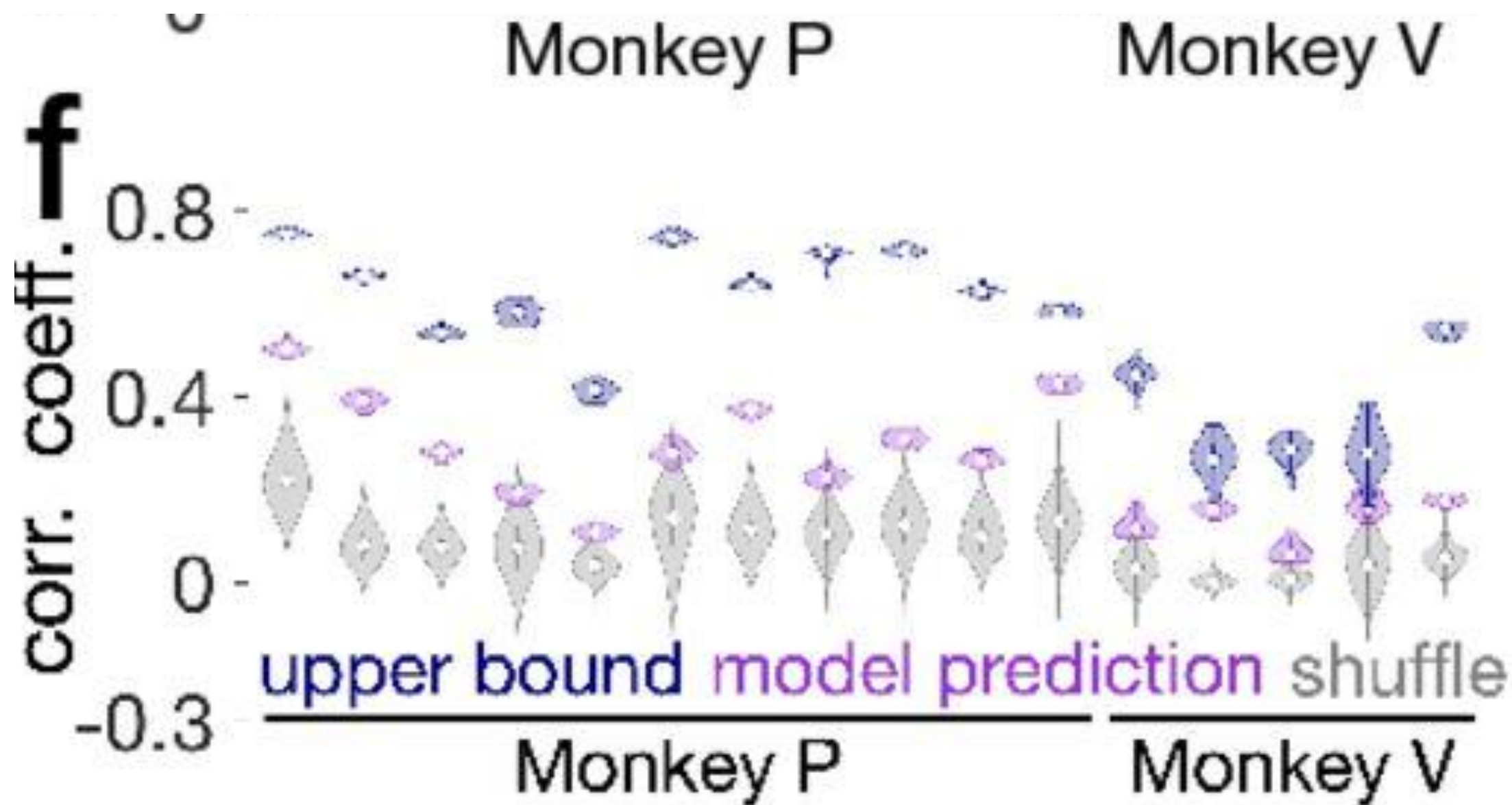
Mechanistic model of ICMS-caused trajectory contraction

- Thus, the nonlinear stimulation effects could be estimated by subtracting estimates of both task-related activity and additive effects from the recorded neural firing rates (**Fig. 6b**).
- We projected the state of the nonlinear stimulation effect at the last time-point of ICMS into the task dynamics space of our E/I model.
- Taking this projected state as the initial condition, we then predicted forward in time using the learned dynamics within the task dynamics space and evaluated Pearson's correlation coefficient between the model prediction and empirical nonlinear stimulation effect throughout the 200ms following the end of stimulation.

E/I model does(n't) explain non-linear effects of ICMS



- **Strikingly, the model predictions showed a strong correlation with the empirical effect in the majority of datasets (Fig. 6f).**
- **This demonstrates that the stimulation indeed engaged the task dynamics space in a way that is predictable based on normal task-related dynamical structure, without any further adjustment of model parameters.**



- The observation of long-lasting, predictable effects following ICMS constitutes strong evidence against path-following dynamics (**H3**). Instead, neural population responses to both optogenetic stimulation and ICMS were indicative of subspace structured dynamics embedded within an E/I network (**H2**).

Abu's Conclusions

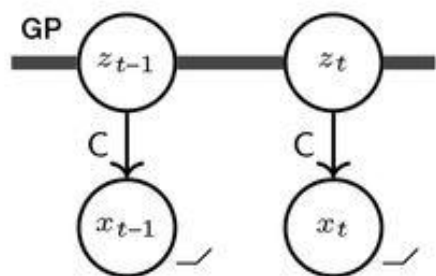
- Optogenetics:
 - Not hitting enough of motor cortex, or
 - Motor cortex is not involved in underlying dynamics
- ICMS
 - Finally hitting enough of motor cortex, or
 - Finally reaching enough circuitry are involved in underlying dynamics

Concerns

- No literature support for prediction of perturbation effects.
- Subspace-driven dynamics and Reservoir dynamics may be equivalent if only a few modes in the reservoir are driving the output?
- Their E/I model is a reservoir model?
 - We then optimized the E/I network connectivity matrix to maximize the likelihood that a linear readout from this task dynamics space matched the measured neural responses.

Supplement

a Non-stimulation neural data:



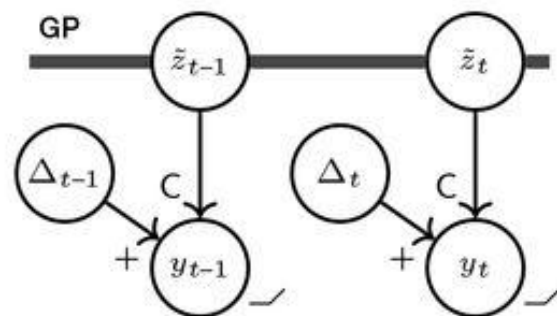
non-stim latents

$z_t \in \mathbb{R}^K$, with GP prior

non-stim rates

$x_t = g(Cz_t + d) \in \mathbb{R}^D$

Stimulation neural data:



stim latents

$\tilde{z}_t \in \mathbb{R}^K$, with GP prior

additive component

$\Delta \in \mathbb{R}^{D \times T}$, low-rank

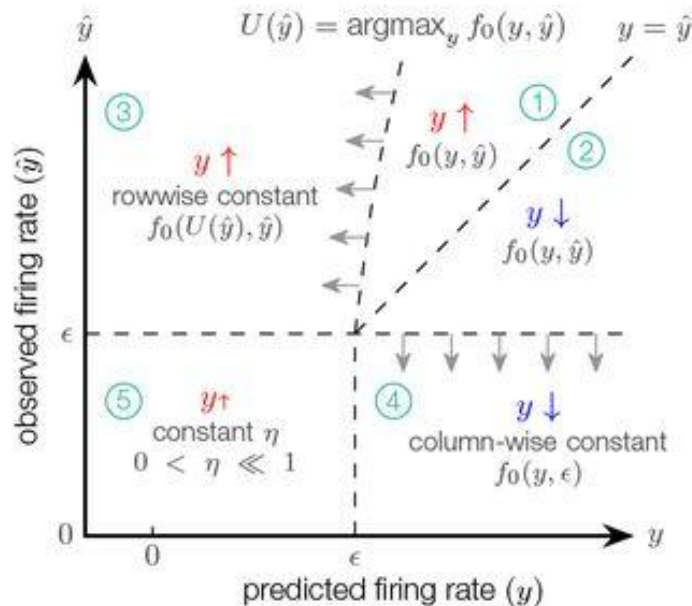
stim rates

$y_t = g(C\tilde{z}_t + \Delta_t + d)$

w/ rectifying non-linearity to correctly handle silencing

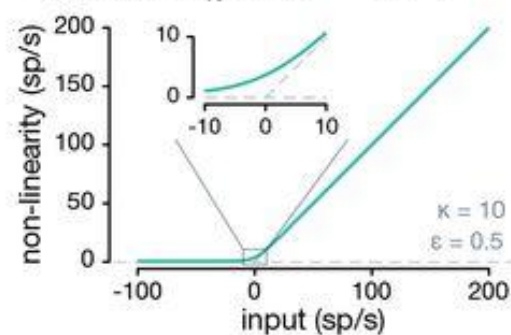
b Soft-plus rectified Poisson observation model:

Schematic of $f(y, \hat{y}) = \frac{\partial \ell}{\partial y} \Big|_{y, \hat{y}}$ defined piecewise via f_0 for Poisson model with soft-plus $g(\cdot)$ applied.



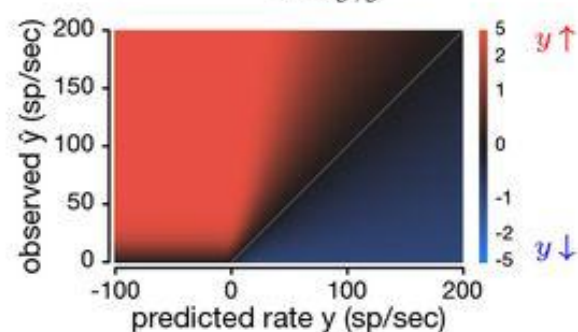
Soft-plus rectification of predicted rates

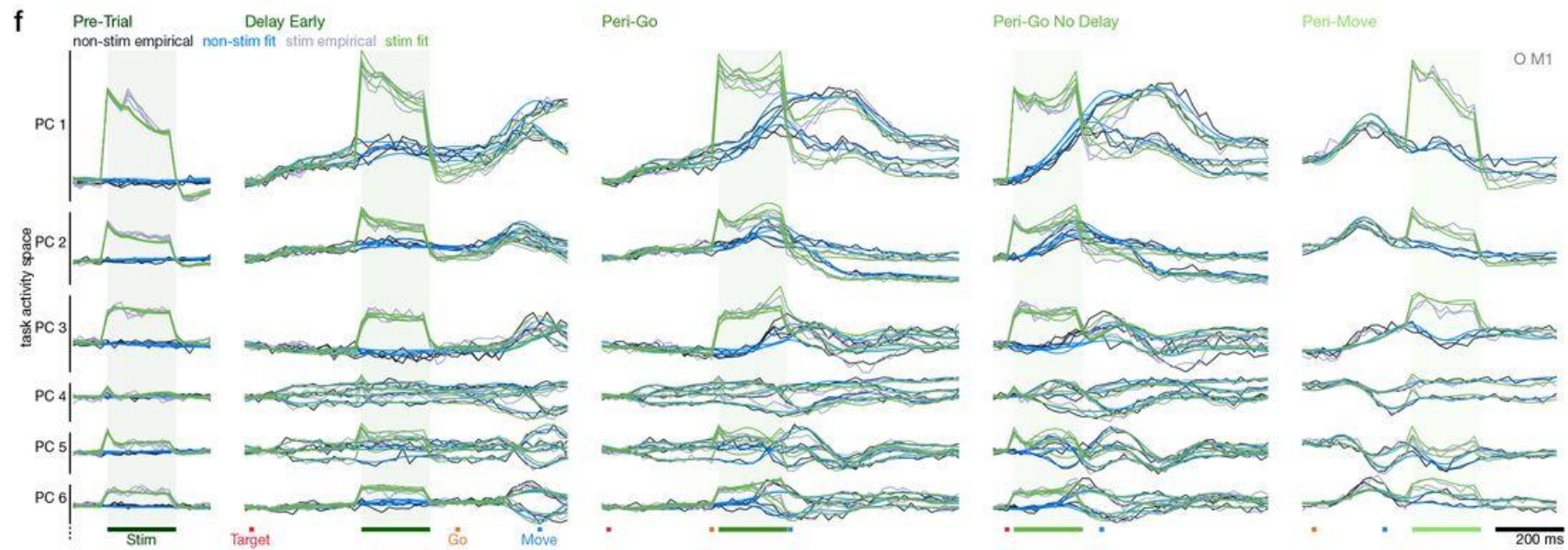
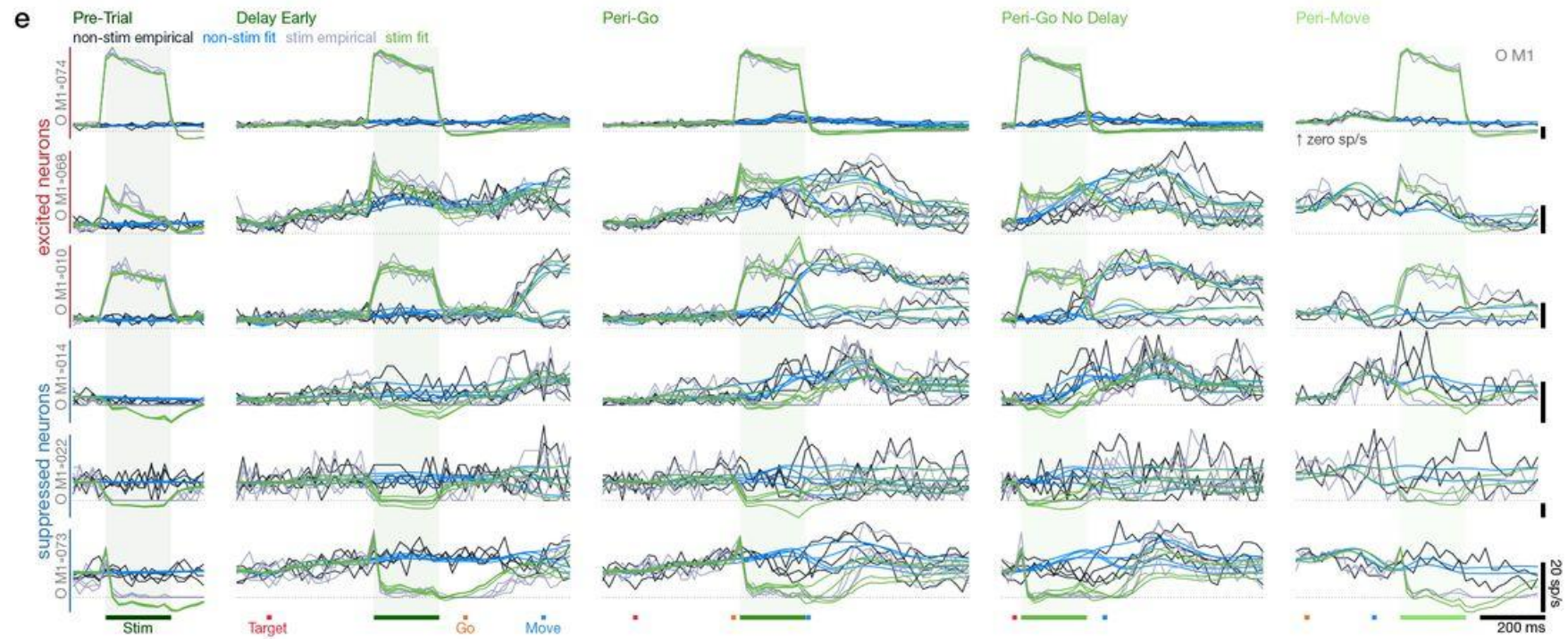
$$g_{\kappa, \epsilon}(y) = \frac{1}{\kappa} \log(e^{\kappa y} + e^{\kappa \epsilon})$$



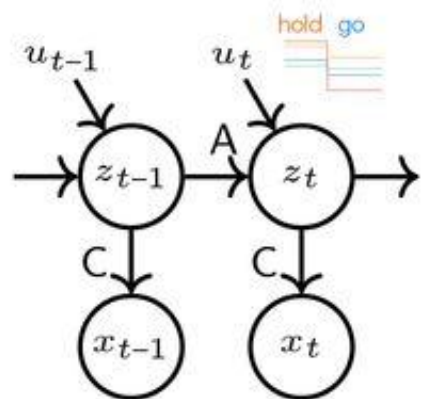
d Gradient of predicted rate

$$f_{\kappa, \epsilon, \eta}(y, \hat{y}) = \frac{\partial \ell}{\partial y} \Big|_{y, \hat{y}}$$





a Latent LDS model



condition-specific input

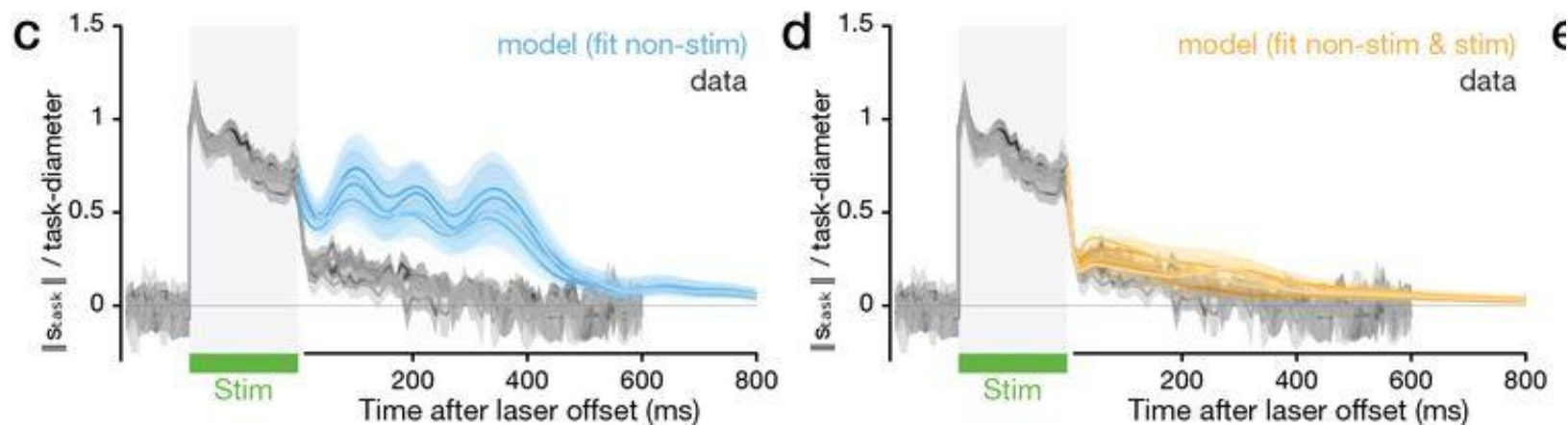
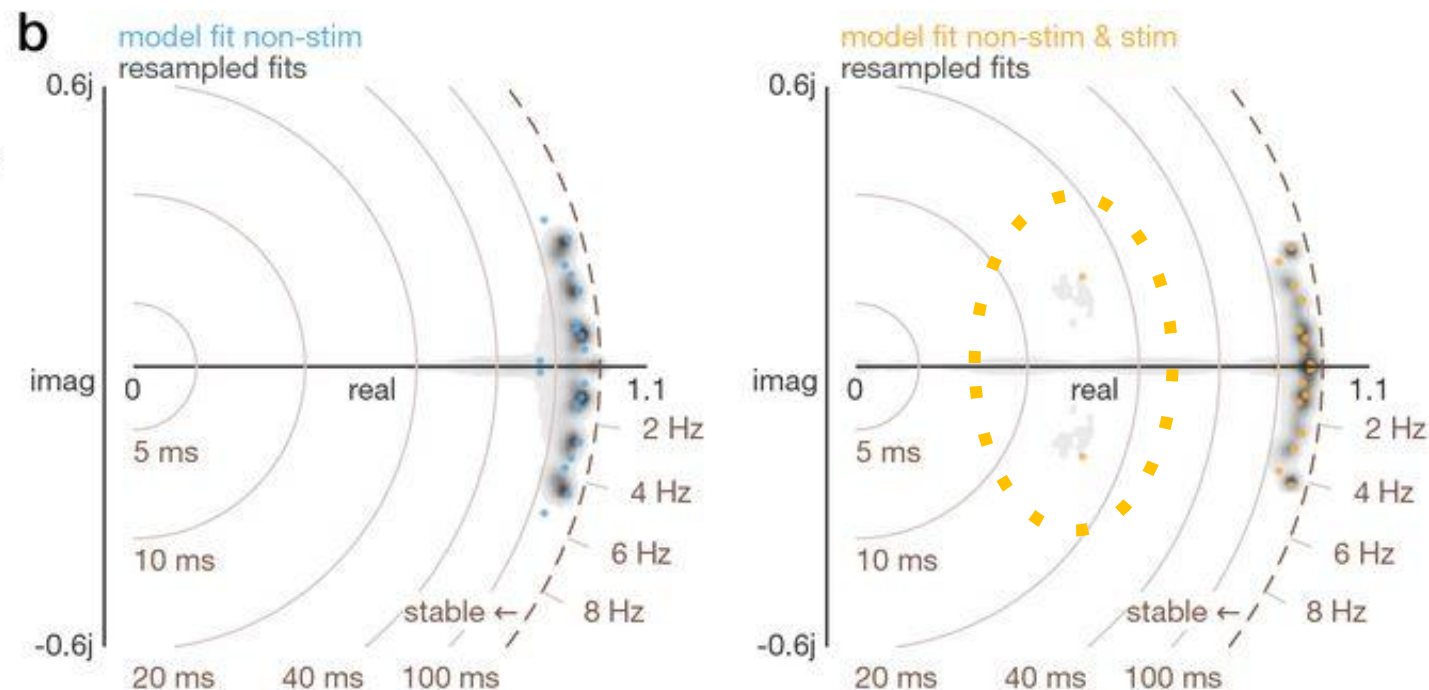
piecewise-constant, $u_t \in \mathbb{R}^K$

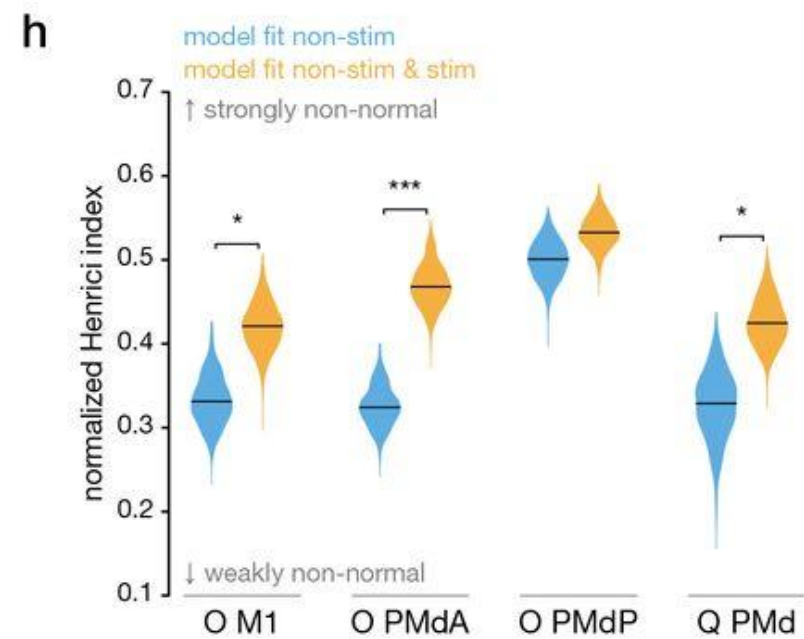
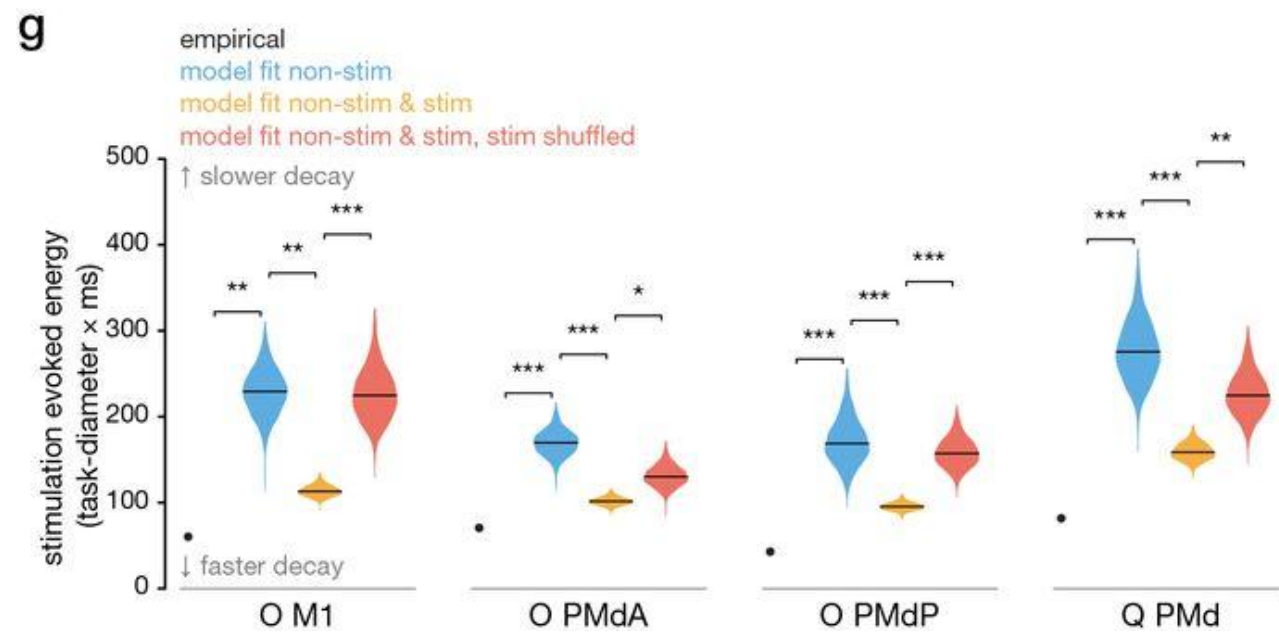
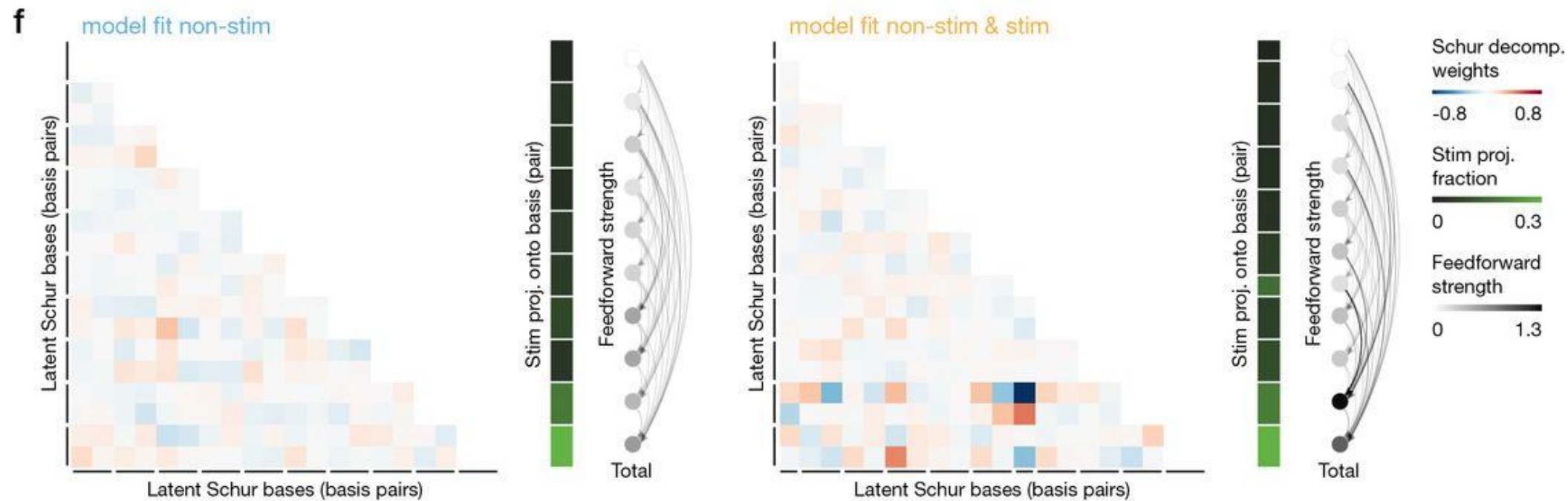
low-d latent dynamics

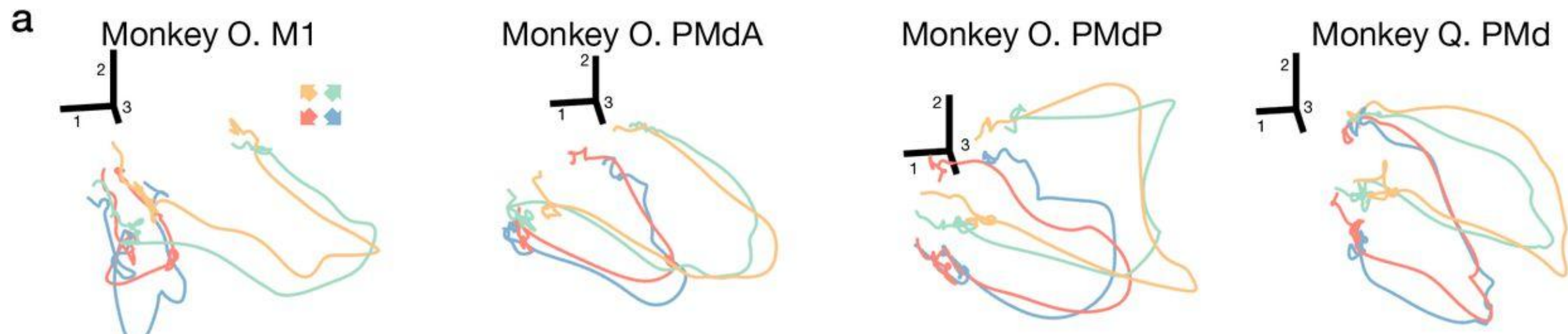
$$z_t = Az_{t-1} + u_t \in \mathbb{R}^K$$

neuron rates

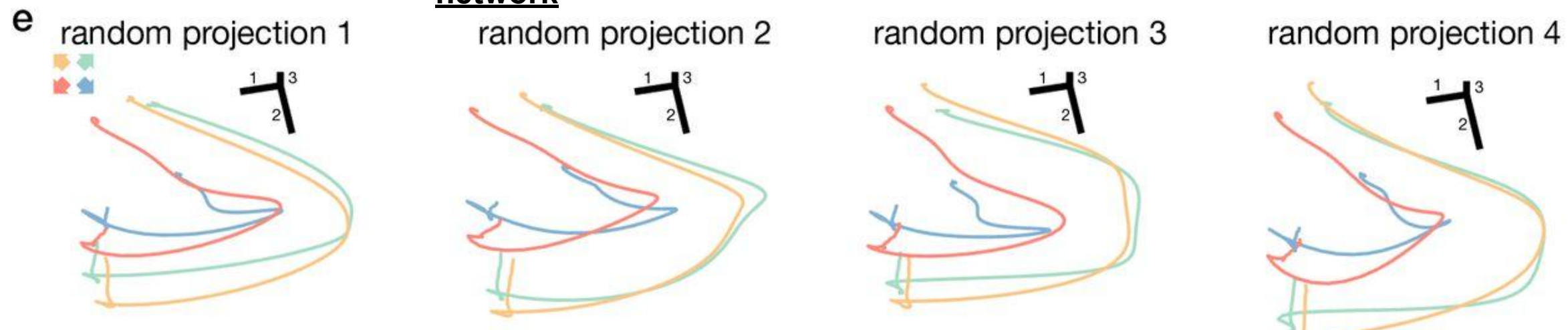
$$x_t = Cz_t + d \in \mathbb{R}^D$$





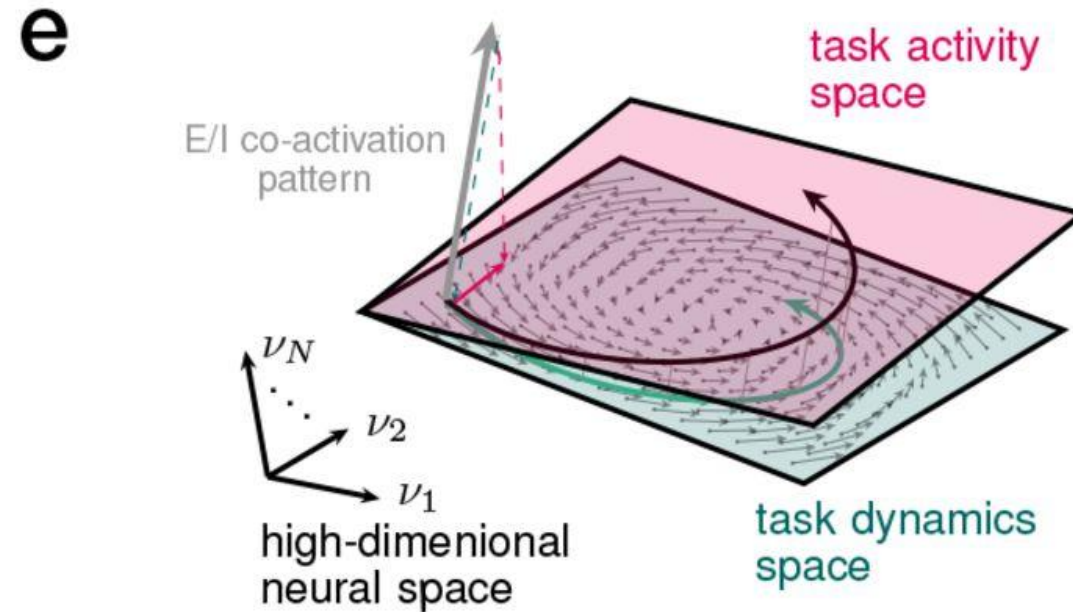
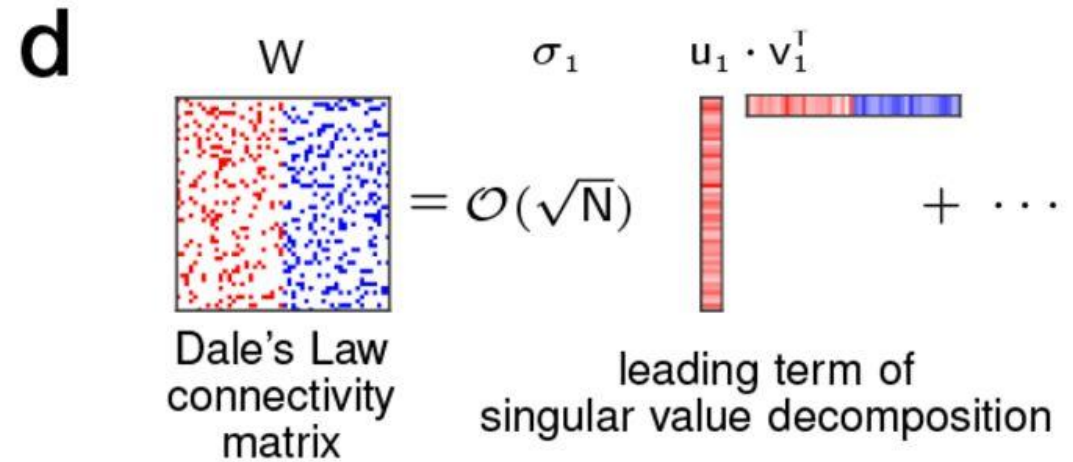
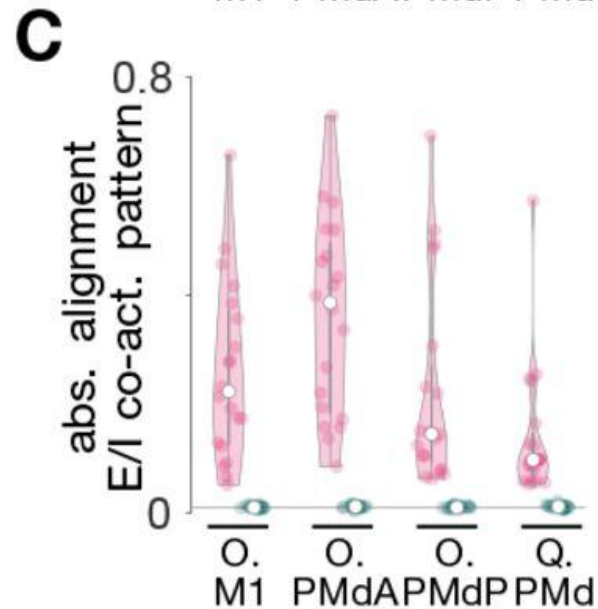
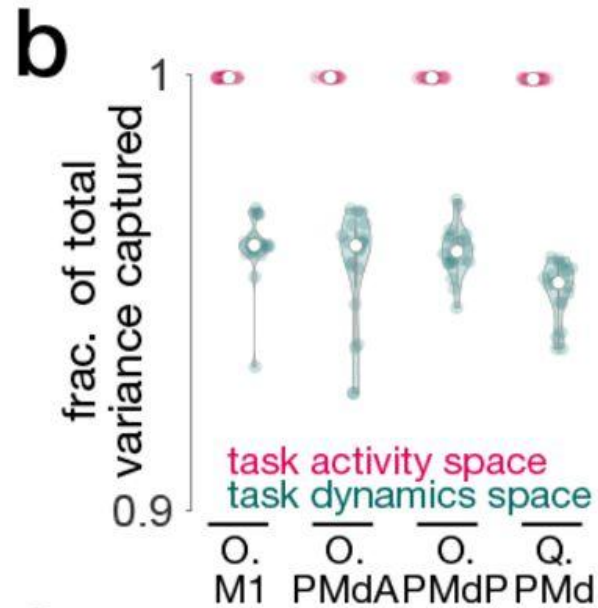


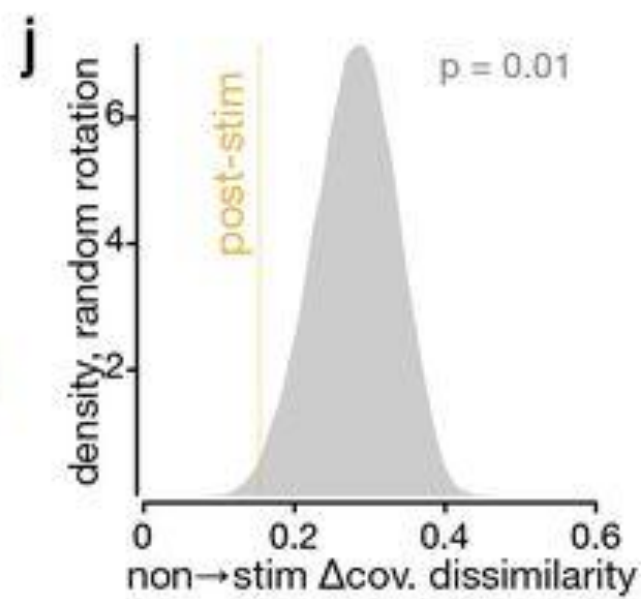
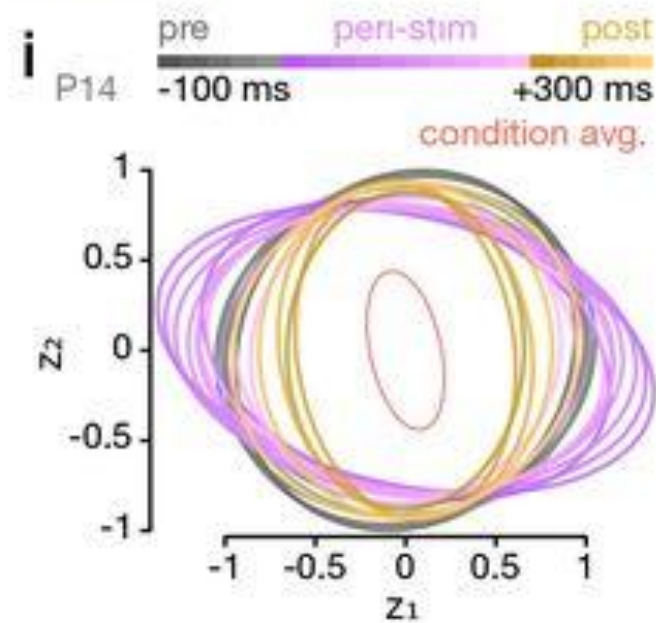
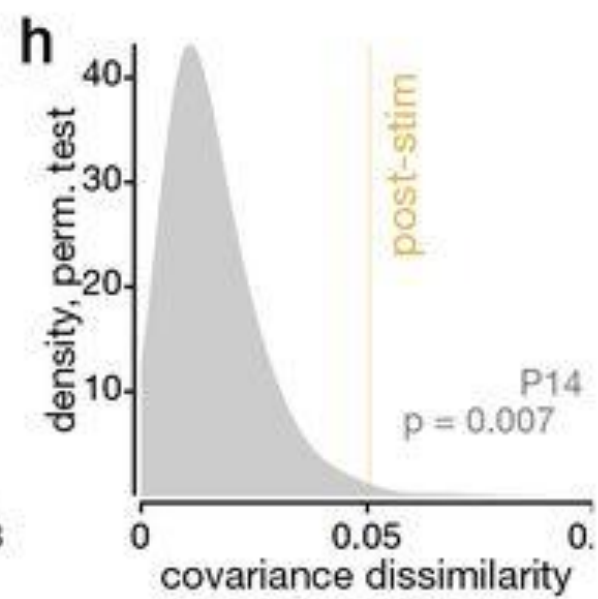
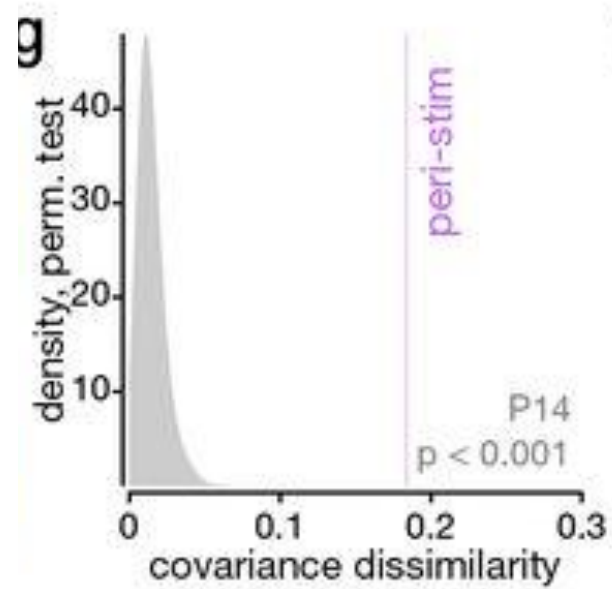
Random Projections of Activity from a trained network



Mechanistic explanation for perturbation results

The task-activity space dimensionality is chosen to match that of the task-dynamics space for each dataset.





k Stimulation neural data, contraction-only model:

