# Reverse-engineering recurrent neural network solutions to a hierarchical inference task for mice
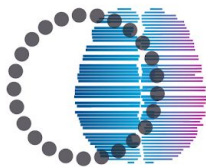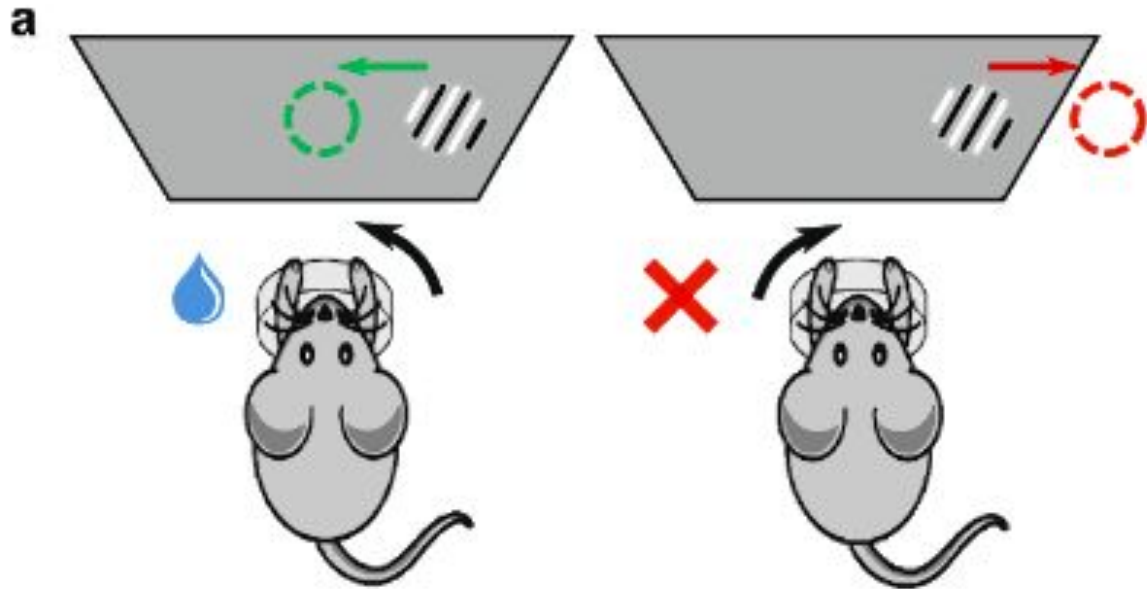
Rylan Schaeffer
NeurIPS 2020

Goal: reverse engineer how recurrent neural networks perform hierarchical inference
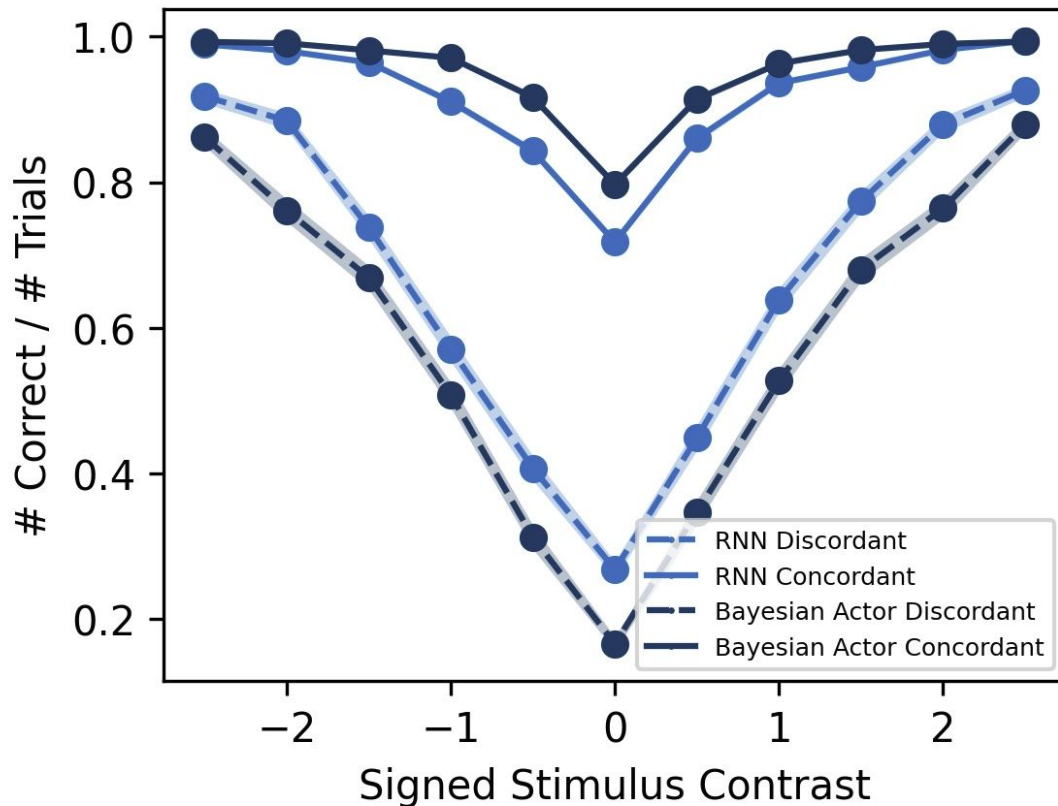
Questions

1. How well do RNNs compare against normative Bayesian baselines?

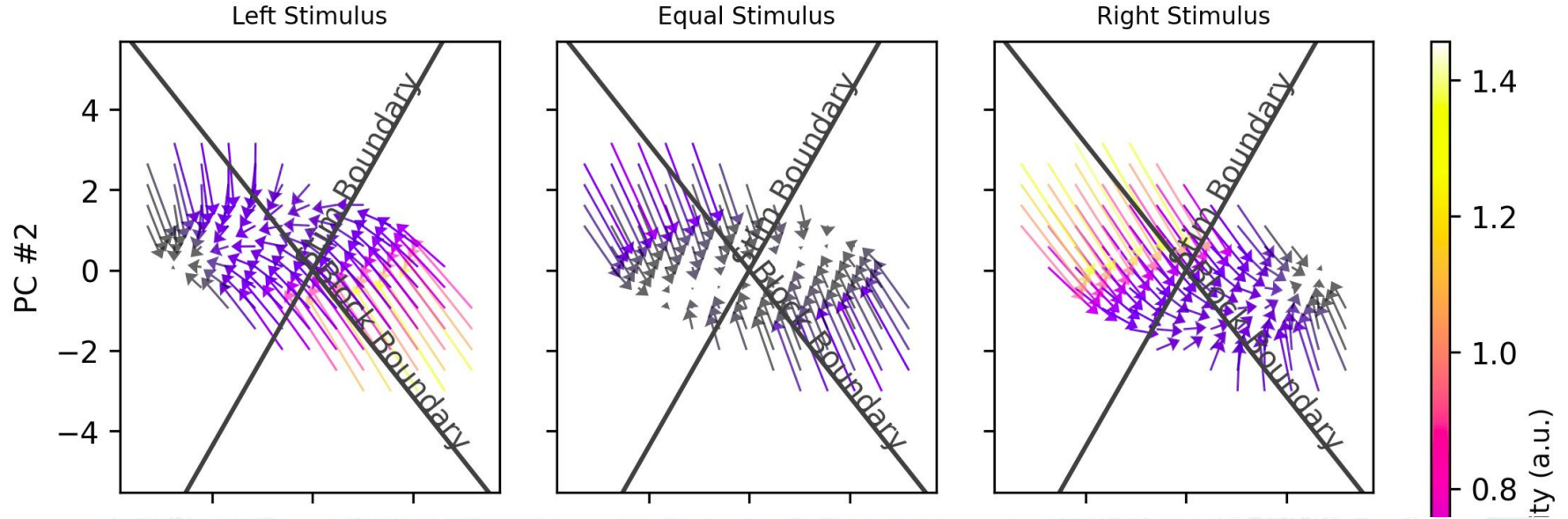2. What are the representations, dynamics and mechanisms RNNs employ to perform inference?

# RNN State Space Displays Two Kinds of Dynamical Behavior

$$\hat{z}_{n,t} = \begin{bmatrix} \text{Stimulus Belief}_{n,t} \\ \text{Block Belief}_{n,t} \end{bmatrix}$$
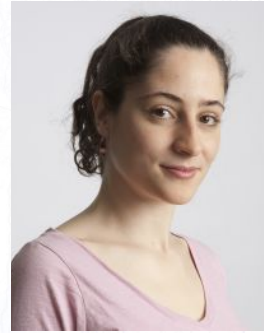
$$\hat{z}_{n,t} = \tanh\left( \begin{bmatrix} 0.54 & 0.31 \\ 0.19 & 0.84 \end{bmatrix} \hat{z}_{n,t-1} + \begin{bmatrix} -0.20 & 0.20 & 0.005 \\ -0.04 & 0.04 & 0.021 \end{bmatrix} \begin{bmatrix} o_{n,t}^{L} \\ o_{n,t}^{R} \\ r_{n,t} \end{bmatrix} \right)$$

Rylan Schaeffer
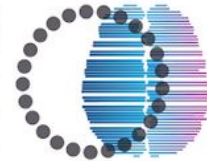Mikail Khona
Leenoy Meshulam
Ila Fiete
& IBL Theory Working Group