

# Expressivity of Neural Networks with Random Weights and Learned Biases

Ezekiel Williams<sup>1</sup>, Alexandre Payeur<sup>1</sup>, Avery Hee-Woon Ryoo<sup>1</sup>, Thomas Jiralerspong<sup>1</sup>, Matthew Perich<sup>1</sup>, Luca Mazzucato<sup>2</sup>, Guillaume Lajoie<sup>1</sup> ✉ezekiel.williams@mila.quebec

## Summary

When all weights are randomly initialized and frozen and only biases are optimized, which we term *bias learning*, we provide math proofs and experiments showing that for wide enough hidden layers:

1. Single-hidden layer Multi-Layer Perceptrons (MLPs) can arbitrarily approximate continuous functions with high probability.
2. Recurrent Neural Networks (RNNs) can arbitrarily approximate finite trajectories from continuous dynamical systems with high probability.

## Motivation

**Neuro:** many biological network parameters are plastic beyond just synapses. In a rate model these are often modelled by bias terms (e.g. tonic input). Could the brain span a wide range of dynamics by adapting only such non-synaptic parameters?

**AI:** tuning weights and biases in a RNN (MLP) yields arbitrary approximation of many dynamical systems (functions) as hidden layer width increases. We don't know if such a result holds for bias learning, despite its relevance to multi-task and in-context learning.

## Theoretical Results for RNNs

Consider the dynamical system, with compact invariant set, given by:

$$\begin{aligned} z[t+1] &= F(z[t], x[t]) \\ y[t+1] &= Cz[t+1], \end{aligned}$$

where  $F$  is continuous, and the RNN:

$$\begin{aligned} r[t+1] &= r[t] + \phi(Wr[t] + Bx[t] + b) \\ \hat{y}[t+1] &= Ar[t+1], \end{aligned}$$

where  $r, b \in \mathbb{R}^m$ .

Let  $\theta \sim \text{unif}[-a, a]$ ,  $a > 0$  for  $\theta$  any element of  $A, B$  or  $W$ .

For finite  $T$ ,  $\varepsilon > 0$ ,  $\delta \in (0, 1)$ , we can find  $m$ ,  $b$ , and  $r_0(z_0)$  such that

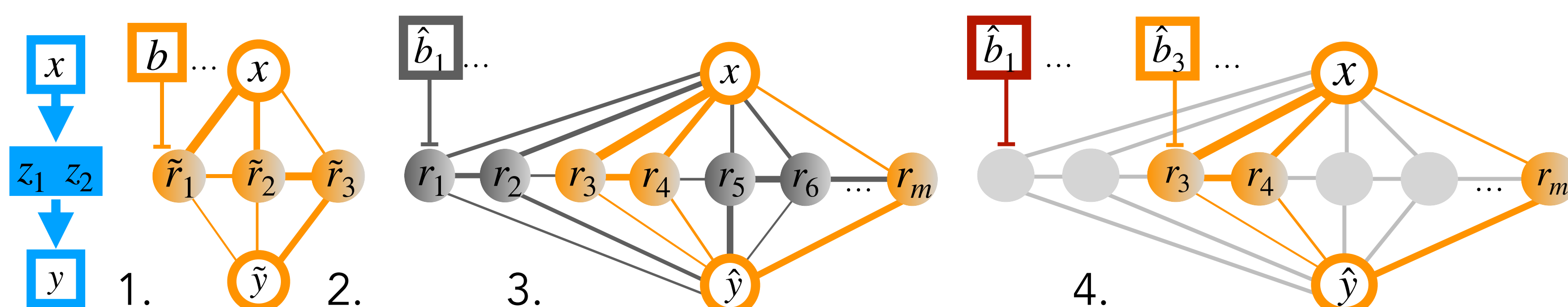
$$\sup_{z_0, x} \sum_{t=1}^T ||y_t(z_0, x) - \hat{y}_t(r_0, x)|| < \varepsilon,$$

with probability greater than  $1 - \delta$  for  $\phi$  ReLU.

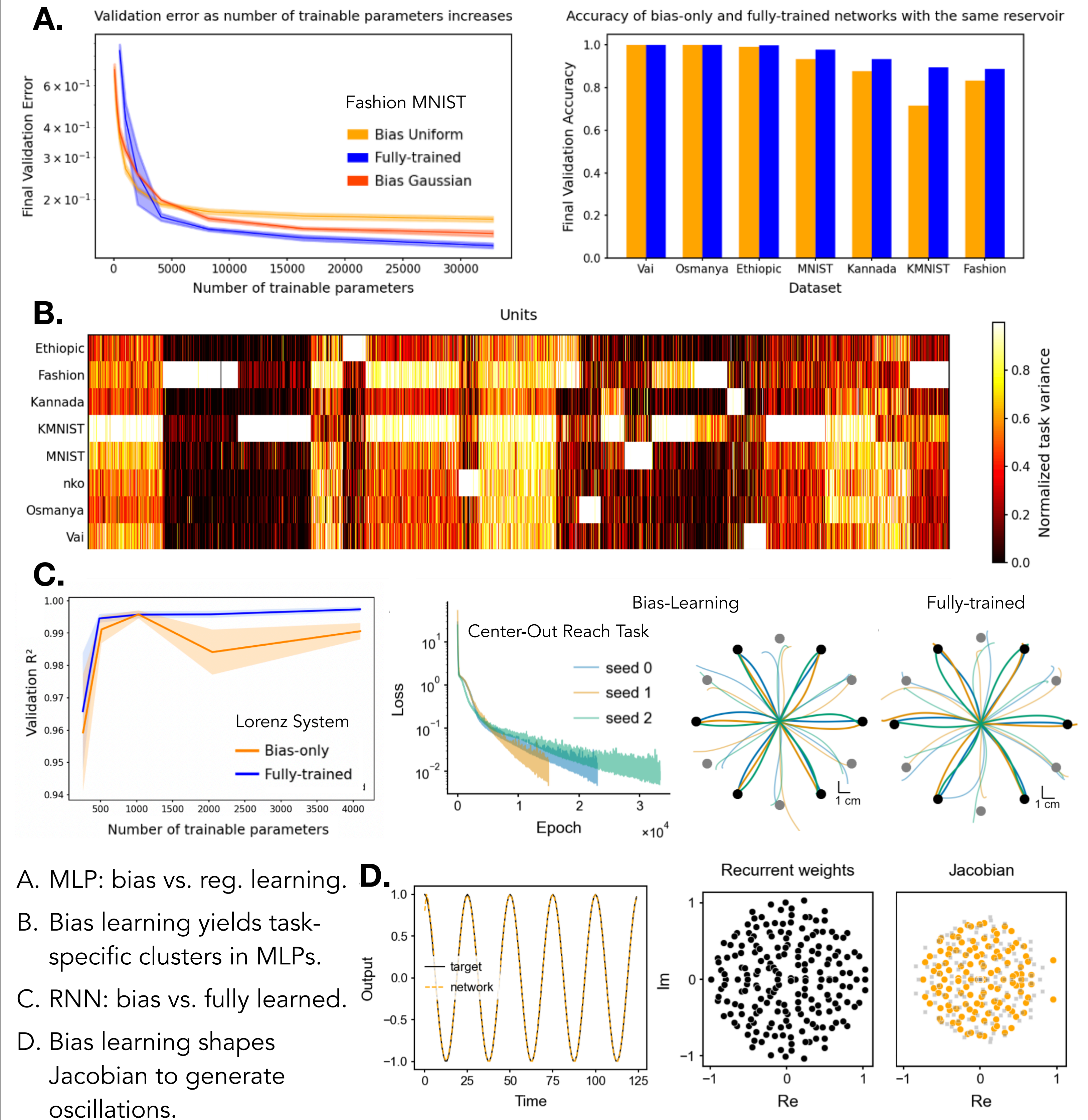
Result for MLPs in preprint

## Theory Intuition

1. Consider a continuous, partially observed Dynamical System (DS).
2. Let  $0 < t \leq T$ . Approximate DS using RNN with trained weights and biases.
3. Randomly initialize a larger RNN. Check for sub-network with weights matching RNN.
4. Adjust biases outside (inside) sub-network to be very negative (match those of RNN).



## Empirical Results



## Conclusions and Future Directions

Synaptic plasticity is not necessarily required for learning task-relevant dynamics in sufficiently large neural networks.

Characterization of network size needed for given degree of error?

What are better distributions for initializing synaptic weights?

Might tonic inputs be adjusted for rapid learning?

Could bias learning yield simpler temporal credit assignment?

## Link to Paper

