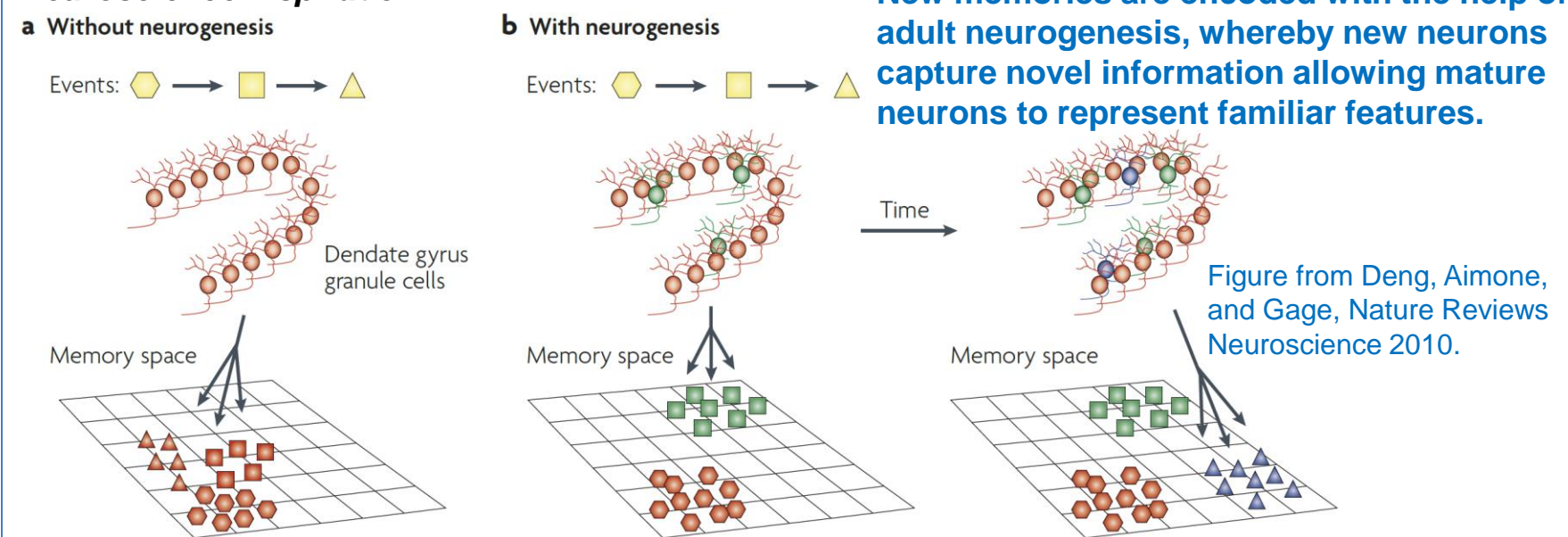




## PROBLEM: How to expand trained deep nets to accommodate new data classes?

In contrast to biological neural systems, capable of continuous learning, DNNs have a limited ability to incorporate new information in a trained network, so methods for continuous learning may be highly impactful in enabling the application of DNNs to dynamic data sets. Inspired by adult neurogenesis in the hippocampus, we explore the potential for adding new nodes to layers of DNNs to facilitate their acquisition of novel information while preserving previously trained representations. Results demonstrate that neurogenesis is well suited for addressing the stability-plasticity dilemma that has long challenged adaptive machine learning algorithms.

### Neuroscience Inspiration



### Neurogenesis algorithm as an efficient method for adapting DNNs

The value of a model to continuously adapt to changing data is challenging to quantify. Here, we quantify the value of a machine learning algorithm at a given time as follows.

$$\text{Utility} = \text{Benefit} - (\text{Cost of Model / Lifetime}) - \text{Cost of runtime}$$

**Extending the lifetime of a model by adapting in response to real-world data changes (e.g., via neurogenesis) mitigates the high initial training costs of DNNs.**

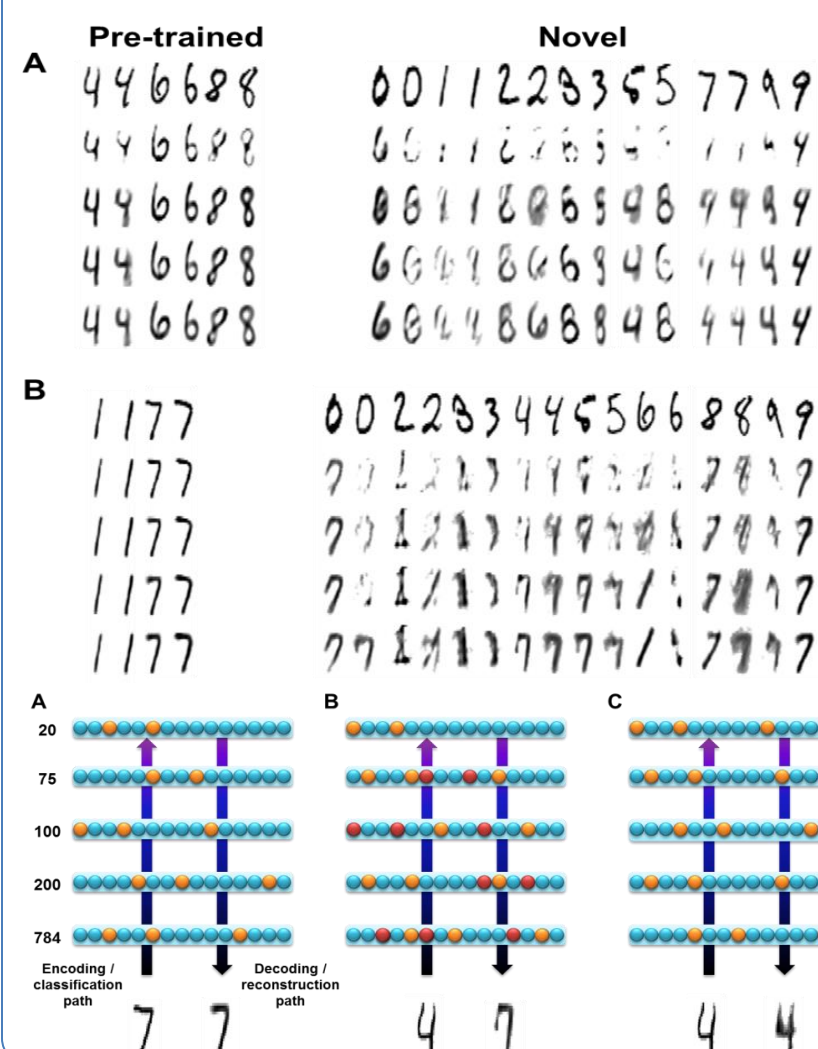
### Layer-wise Reconstruction Error as a measure of representation capability

RE is computed at internal layer  $L$  within an AE by encoding an input sample through  $L$  encode layers, then propagating through the corresponding  $L$  decode layers to the output. An AE parameterized with weights  $W$ , biases  $b$ , and activation function  $s$  is described from input,  $x$ , to output as  $N$  encode layers followed by  $N$  decode layers.

$$\text{Encoder: } f_{\theta_N} \circ f_{\theta_{N-1}} \cdots f_{\theta_2} \circ f_{\theta_1}(x) \text{ where } y = f_{\theta}(x) = s(Wx + b)$$

$$\text{Decoder: } g_{\theta'_N} \circ g_{\theta'_{N-1}} \cdots g_{\theta'_2} \circ g_{\theta'_1}(y) \text{ where } g_{\theta'}(y) = s(W'y + b')$$

$$\text{Then, the global RE at layer } L \text{ is } RE_{Global,L}(x) = (x - g_{\theta'_N} \circ \cdots g_{\theta'_{N-L}} \circ f_{\theta_L} \circ \cdots f_{\theta_1}(x))^2$$



### PROBLEM

Existing AE can't represent novel data

Networks initially trained on (A) 4, 6, and 8's and (B) 1 and 7's and not yet trained on any of the other MNIST digits reconstruct those novel digits using features biased by their original training data. The first row of each panel is the original data, with the 2<sup>nd</sup> - 5<sup>th</sup> rows the reconstruction through the 1<sup>st</sup> - 4<sup>th</sup> (full) encode and corresponding decode layers of the autoencoder, respectively.

### SOLUTION

Add new node(s) to each layer of an AE as needed to represent novel data

- Deep AE can faithfully reconstruct originally trained digit 7.
- AE fails at reconstructing novel digit 4.
- New nodes added to layers 2, 3, and 4 enable AE to be trained to reconstruct the 4.

## Neurogenic Deep Learning (NDL) Algorithm

The NDL algorithm adds and trains new nodes to a layer of an AE similar to layerwise pretraining when a critical number of input samples fail to achieve adequate representation.

- Plasticity** occurs as new nodes are added to represent novel data, then network
- Stability** occurs by leveraging both new data and replayed samples from previously seen classes. Samples from old classes are created using the representation capability of the AE in a process we call "intrinsic replay" (see below).

**Input:** 2N-layer autoencoder AE trained on data classes  $\{D_1, D_2, \dots, D_{U-1}\}$ , new class of data  $D_U$ , vector of per-layer RE thresholds  $Th$ , vector of per-layer maximum nodes allowed to add  $MaxNodes$ , maximum number of samples allowed to have  $RE_L > Th_L$ ,  $MaxOutliers$ , Learning Rate LR

**Output:** Autoencoder AE capable of representing data classes  $\{D_1, D_2, \dots, D_U\}$

// Combine samples from the new class of data with replayed samples of old data  
 $AE\_TrainingSamples \leftarrow \{D_U \cup \text{IntrinsicReplay}(D_1, D_2, \dots, D_{U-1})\}$

// Perform neurogenesis layer by layer

$numOutliers \leftarrow |D_U|$

for Layer  $L \leftarrow 1$  to  $N$

$NewNodes \leftarrow 0$

$Outliers \leftarrow \{d \in D_U \mid RE_{Global,L}(d) > Th_L\}$

$N_{Out} \leftarrow |Outliers|$

$N_{Out}^{Prev} \leftarrow N_{Out} + 1$

// Add and train new nodes to layer  $L$

while  $N_{Out} > MaxOutliers$  and  $NewNodes < MaxNodes_L$  and  $N_{Out} < N_{Out}^{Prev}$

$AE_L \leftarrow (W_L, b_L; W'_{N+1-L}, b'_{N+1-L})$  from AE

**Plasticity:** Add a node with random weights to  $AE_L$  and train on  $Outliers$

Use LR to update encoder weights connected into new node only

Use LR/100 to update decoder weights

**Stability:** Train  $AE_L$  on  $AE\_TrainingSamples$

Using LR to update all weights

$(W_L, b_L; W'_{N+1-L}, b'_{N+1-L}) \leftarrow AE_L$

$Outliers \leftarrow \{d \in D_U \mid RE_{Global,L}(d) > Th_L\}$

$N_{Out}^{Prev} \leftarrow N_{Out}$

$N_{Out} \leftarrow |Outliers|$

$NewNodes \leftarrow NewNodes + 1$

// Add connections from new nodes in layer  $L$  to existing nodes in layer  $L+1$  and train

If  $NewNodes > 0$  &  $L < N$

**Plasticity:** Add weights initialized to zero to  $AE_{L+1}$  connected to new nodes from layer  $L$

**Stability:** Train  $AE_{L+1}$  on  $AE\_TrainingSamples$

$(W_L, b_L; W'_{N+1-L}, b'_{N+1-L}) \leftarrow AE_L$

## Intrinsic Replay (IR) Algorithm

The hippocampus region of the brain is known to "replay" experienced activity to aid the consolidation of newly acquired information into memory.

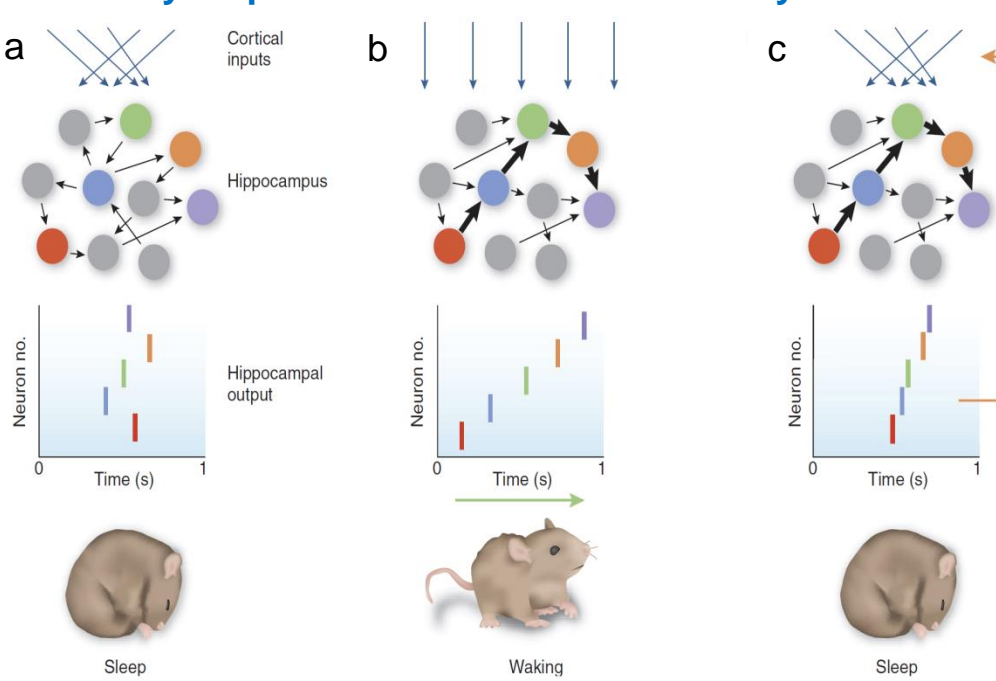
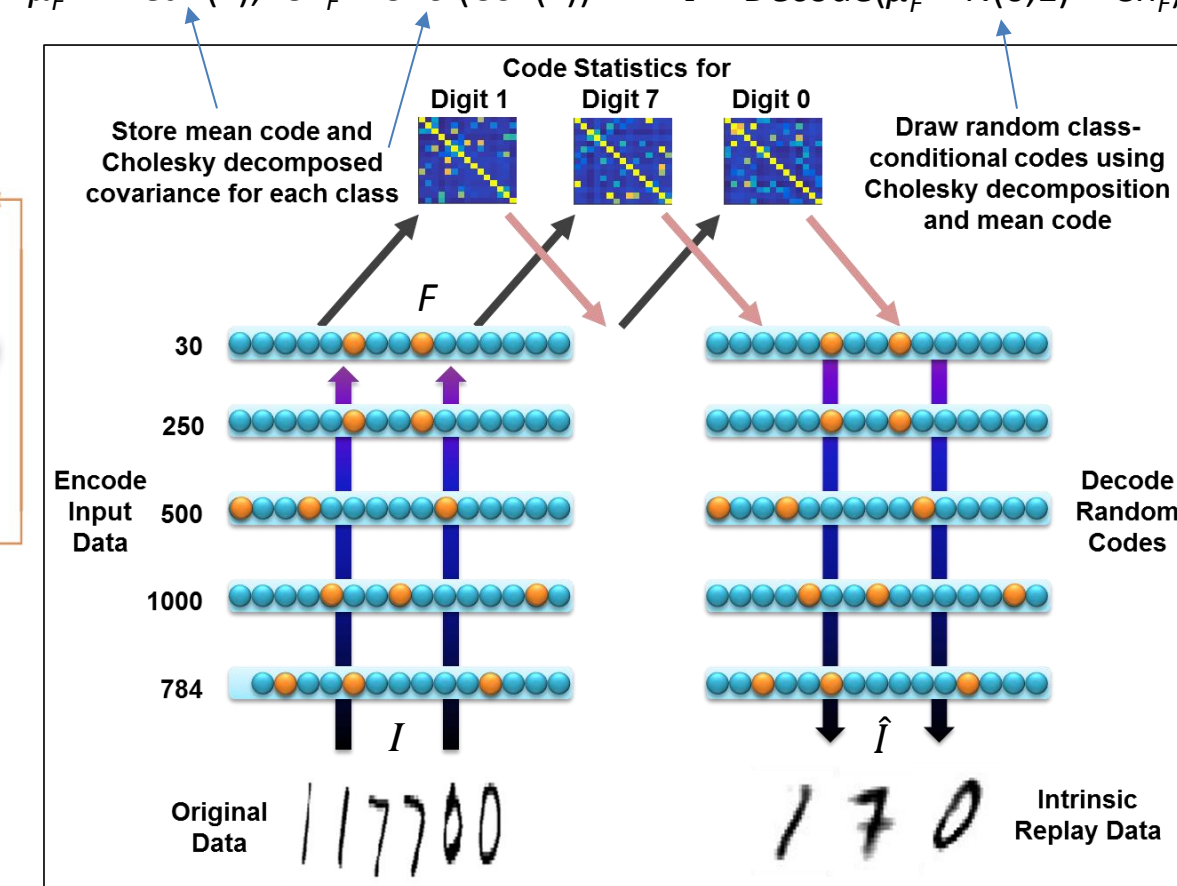


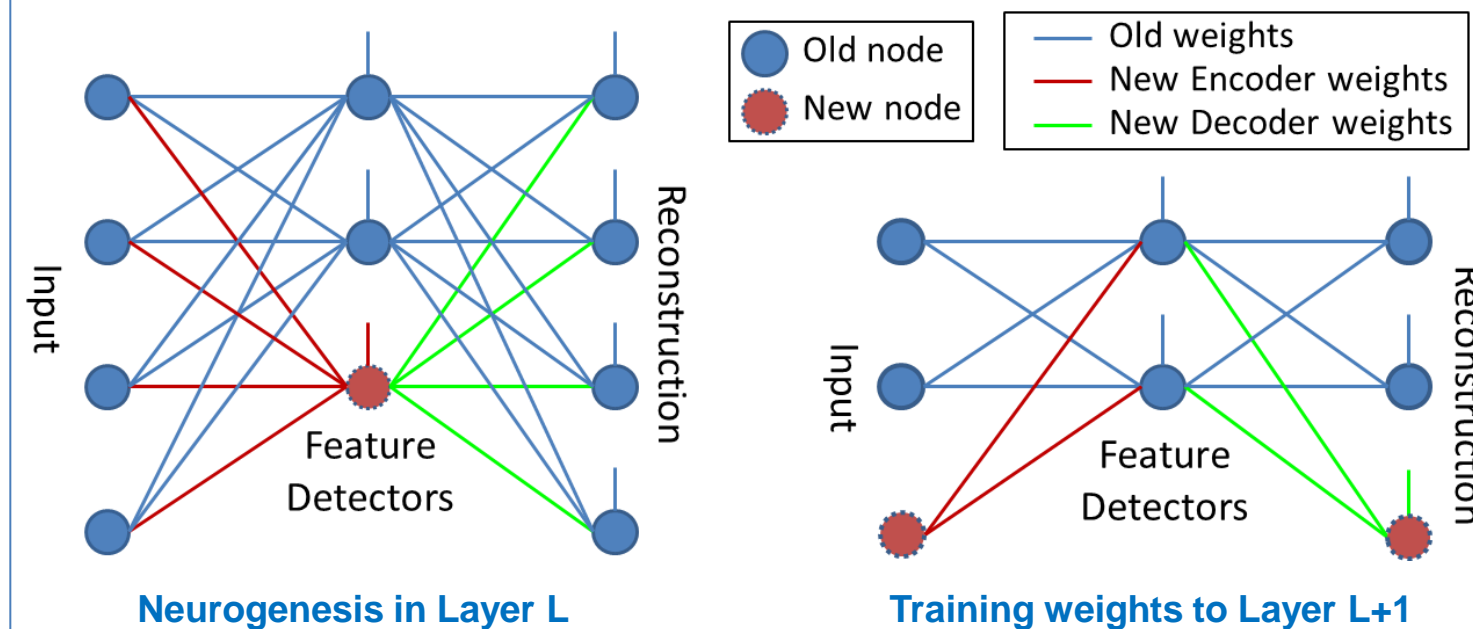
Figure from Mehta et al., Nature Neuroscience 2007

$$\mu_F = \text{Mean}(F), \quad Ch_F = \text{Chol}(\text{Cov}(F)) \quad \hat{I} = \text{Decode}(\mu_F + N(0,1) * Ch_F)$$

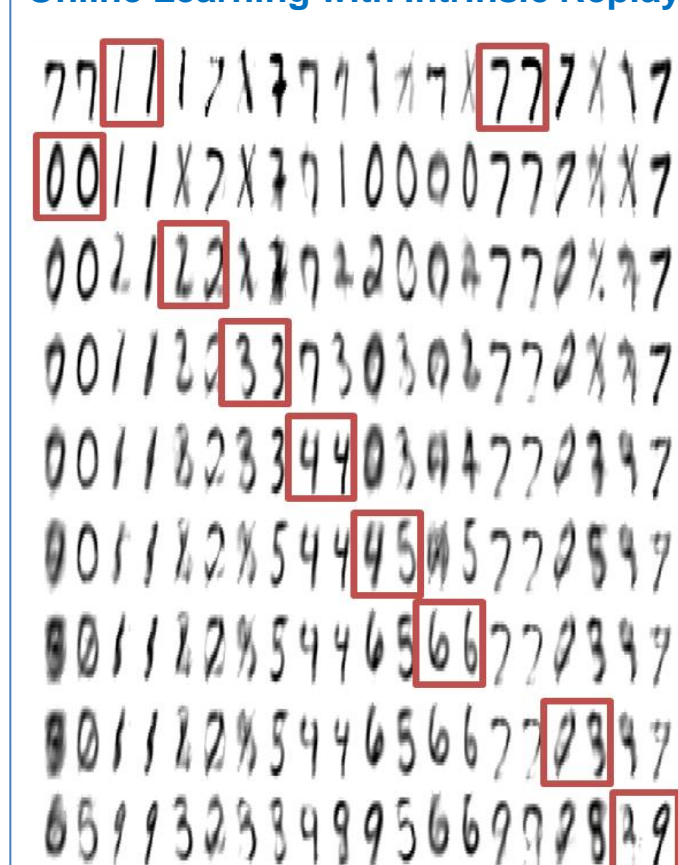


## Greedy Layerwise Neurogenesis on an Autoencoder.

The goal is to learn new feature detectors for novel data.

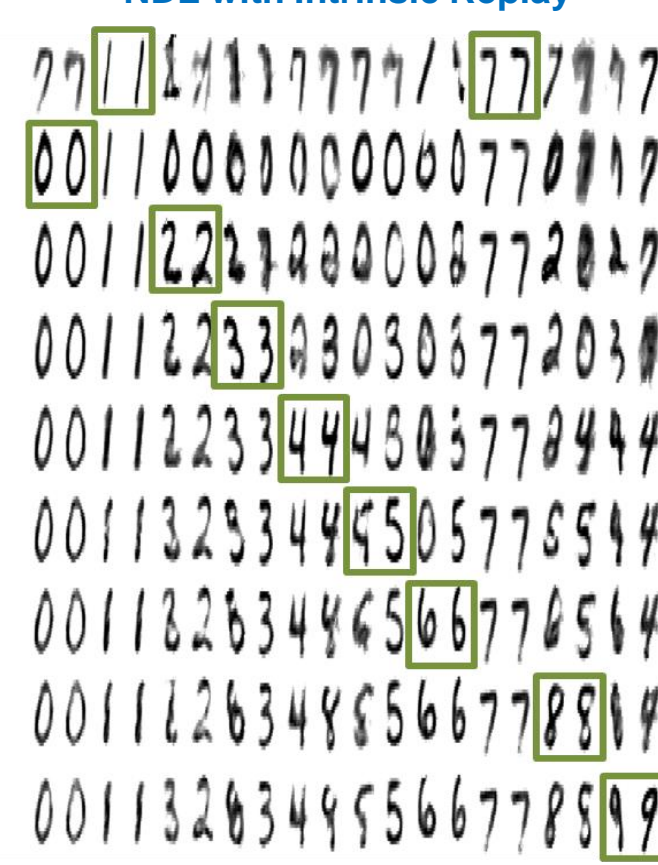


## Online Learning with Intrinsic Replay



Reconstructions of all digits by AE pre-trained on digits 1 and 7' networks and progressive learning of new classes.

## NDL with Intrinsic Replay



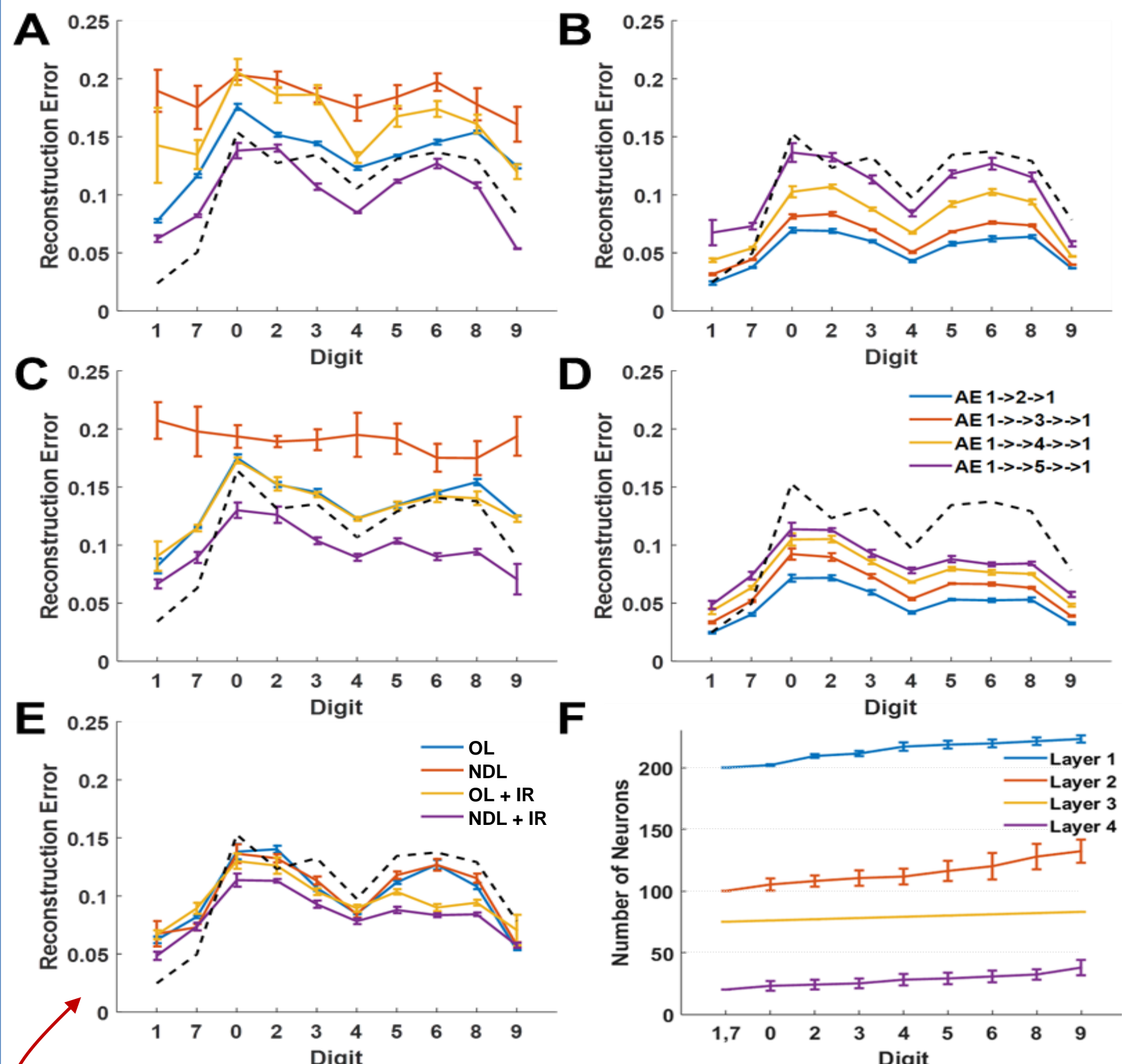
## RESULTS

We evaluate NDL on MNIST data, where a Deep AE is initially trained with two digits (1, 7), then learning a new task is simulated by progressively expanding the number of encountered classes (0, 2, 3, 4, 5, 6, 8, 9), one at a time. For each experiment, all training samples in a class are presented at once. For OL networks (A & C), the entire AE is retrained as each new class of data is presented.

- NDL with IR (NDL+IR) - Figure D: Starting network size of 784-200-100-75-20-75-100-200-784
- OL with IR (OL+IR) - Figure C below: Starting network size = NDL+IR generated network
- NDL without IR (NDL) - Figure B below: Starting network size of 784-200-100-75-20-75-100-200-784
- OL without IR (OL) - Figure A below: Starting network size = NDL generated network

Our results show that NDL with IR enables training of new digits while minimally impairing original representations.

- NG+IR outperforms OL not only overall, but in both the ability to represent the new data as well as preserving the ability to represent previously trained digits.
- OL+IR performs well on new digits, but poorly on retaining original old digits, whereas the NG+IR process does well on all digits.



A-E: Average REs of trained AEs after exposure to all 10 digits (legend in Plot D applies to Plots A, B, & C; dotted line shows REs of the original AE trained just on 1 and 7).

E: RE Comparison of the output of each full AE showing the superior performance of NDL+IR.

F: Neurogenesis contribution to network size in NDL+IR networks

## FUTURE WORK

- Apply Neurogenic Deep Learning to Convolutional Neural Nets
- Add new filters via neurogenesis
- Test classification performance
- Apply to streaming data to address concept drift