

---

# Efficient rescue of damaged neural networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural networks in the brain and in neuromorphic chips confer systems with the ability to perform multiple cognitive tasks. However,  
2 both kinds of networks experience a wide range of physical perturbations, ranging from damage to edges of the network to complete  
3 node deletions, that ultimately could lead to network failure. A critical question is to understand how the computational properties of  
4 neural networks change in response to node-damage and whether there exist strategies to repair these networks in order to compensate  
5 for performance degradation. Here, we study the damage-response characteristics of two classes of neural networks, namely multilayer  
6 perceptrons (MLPs) and convolutional neural networks (CNNs) trained to classify images from MNIST and CIFAR-10 datasets  
7 respectively. We also propose a new framework to discover efficient repair strategies to rescue damaged neural networks. The  
8 framework involves defining damage and repair operators for dynamically traversing the neural networks loss landscape, with the  
9 goal of mapping its salient geometric features. Using this strategy, we discover features that resemble path-connected attractor sets in  
10 the loss landscape. We also identify that a dynamic recovery scheme, where networks are constantly damaged and repaired, produces  
11 a group of networks resilient to damage as it can be quickly rescued. Broadly, our work shows that we can design fault-tolerant  
12 networks by applying on-line retraining consistently during damage for real-time applications in biology and machine learning.

## 13 1 Introduction

14 Living neural networks in the brain and artificial networks engineered on neuromorphic chips [1] perform an array of computational and  
15 information processing tasks [2, 3, 4]. However, both these networks are susceptible to physical perturbations that lead to a decline in  
16 functional performance [5]. Understanding how damage of neural units in a network leads to cognitive decline is of great interest to  
17 biomedical sciences as well as to AI practitioners implementing artificial networks on neuromorphic hardware. In addition, deciphering  
18 techniques to ‘search’ for neural networks that are resilient to perturbation and strategies that efficiently rescue damaged networks to  
19 compensate for performance degradation are of great interest to both the communities. So far, researchers have only focused on studying  
20 resilience of neural nets to perturbation of input signals [6] by generating adversarial examples [7] that highlight the vulnerability of neural  
21 nets. However, not much has been done towards understanding the decline in performance due to physical perturbation of neural networks  
22 [8] and unraveling repair strategies to rescue damaged networks.

23 In this paper, inspired by the powerful paradigms introduced by deep learning, we attempt to understand the computational and mathematical  
24 principles that impact the ability of neural networks to tolerate damage and be repaired. We characterize the response of two classes  
25 of neural networks, namely multilayer perceptrons (MLP’s) and convolutional neural nets (CNN’s) to node-damage and propose a *new*  
26 *framework* that identifies strategies to efficiently rescue damaged networks in a principled fashion.

27 Our key contribution is the introduction of a framework that conceptualizes damage and repair of networks as operators of a dynamical  
28 system in the high-dimensional parameter space of a neural network. The damage and repair operators are used to dynamically traverse the  
29 landscape with the goal of mapping local geometric features [9, 10] (like, fixed points, limit-cycles or point/line-attractors) of the neural  
30 networks’ loss landscape. The framework led us to discovering that the iterative application of damage and repair operators results in  
31 networks that are highly resilient to node-deletions as well as guides us to uncover the presence of geometric features that resemble a  
32 path-connected attractors set, in many respects, in the neural networks’ loss landscape. Attractor-like geometric features in the networks’  
33 loss landscape explains why the iterative damage-repair strategy always results in the rescue of damaged networks within a small number of  
34 training cycles.

## 35 2 Susceptibility of neural networks to damage

36 The first question we ask in this paper is how do neural networks respond to physical perturbations and how does it affect their functional  
37 performance. We characterize the impact of neural damage on ‘cognitive’ performance of neural networks by tracking the performance of  
38 two classes of artificial neural networks, namely MLPs and CNNs, to deletion of neural units from the network. The MLPs and CNNs were  
39 trained to perform simple cognitive tasks like image classification on MNIST and CIFAR-10 datasets respectively before the networks were  
40 perturbed.

41 To damage a node  $i$  in the hidden layer of an MLP or in the fully connected layer of a CNN, we zero all connections between node  $i$  and the  
 42 rest of the network. And, to damage a node  $j$  in the convolutional layer of a CNN, we zero the entire feature map. In this paper, we are  
 43 specifically interested in node-damage as our perturbation because of its similarity in phenomena to neuron death in biological networks  
 44 and node-failures in neuromorphic hardware.

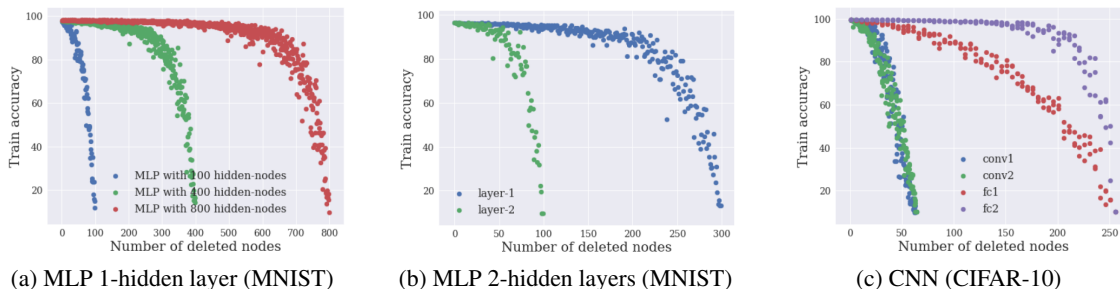


Figure 1: **Damage of neural units in artificial neural networks (Phase transition)** (a) Performance of MLP with 1 hidden layer in MNIST classification (b) Performance of MLP with 2 hidden layers in MNIST classification (c) Performance of CNN in CIFAR-10 classification

45 We observe a steep increase in the rate of decline of functional performance as we incrementally delete nodes from either an MLP with 1  
 46 hidden layer (Fig-1a), an MLP with 2 hidden layers (Fig-1b) or a CNN with 2 convolutional layers, a pooling layer and 2 fully connected  
 47 layers (Fig-1c). We refer to this discrete jump in the rate of decline of performance as a **phase transition**.

48 The existence of a phase transition shows that neural nets (MLP's and CNN's) damaged above their respective critical thresholds are not  
 49 resilient to any further perturbation. We are interested in deciphering strategies that enable the quick rescue of damaged neural nets and also  
 50 want to identify networks that are more resilient to perturbation.

### 51 3 Can we rescue these damaged networks?

52 We ask whether it is fundamentally possible to rescue damaged networks in order to compensate for their performance degradation. To do  
 53 so, we re-train damaged networks via two strategies mentioned below: **(Strategy-1: Functional sub-network retraining)** Purely re-train  
 54 the 'functioning' sub-network, ie the weights connecting damaged nodes are kept at zero, while enabling plasticity for weights connecting  
 55 the remaining undamaged nodes. **(Strategy-2: Node replacement)** Replace the damaged units with 'embryonic' nodes and retrain the  
 56 network such that 'embryonic' nodes are more plastic than the nodes in the functioning sub-network.

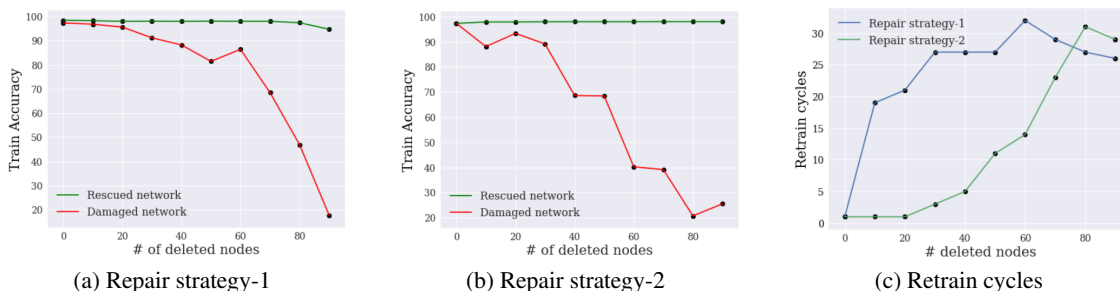


Figure 2: **Repairing damaged networks** (a,b) The red line is the accuracy of the network after  $m$  nodes have been damaged. The green line is the accuracy of network after it has been retrained using one of two strategies. (a) [Strategy-1] Purely retraining the functioning sub-network. (b) [Strategy-2] Replacing damaged nodes with 'embryonic' nodes and selectively retraining only the newly replaced nodes. (c) The number of training cycles required for repairing a damaged network with  $m$  nodes via both strategies.

57 The plots in figure-2 show that damaged neural networks can be rescued to regain their original functional performance when re-trained  
 58 via both strategies 1 and 2. However, they require a large number of training cycles (epochs) to be effectively rescued (figure-2c). The  
 59 requirement of a large number of training cycles for the effective rescue of a neural network *reduces the feasibility* of either strategy as it  
 60 isn't ideal for both, living neural networks in the brain or artificial networks implemented on neuromorphic hardware to be re-trained for  
 61 extended periods of time to recover from small damages to its network.

### 62 4 Iterative repair of networks v/s batch repair of networks

63 Inspired by the dynamic recovery paradigm adopted by most biological systems, where networks are constantly being perturbed and  
 64 repaired, we propose an iterative damage-repair strategy and test whether this produces networks that are more resilient to perturbation as  
 65 well as if it allows us to rescue damaged networks with much lesser training cycles.

66 Figure-3 demonstrates that the iterative damage-repair paradigm can rescue neural networks to their  
 67 original functional performance within 15 training cycles! This is in stark contrast to the batch recovery  
 68 of networks, either via strategy 1 or 2, as they need up to 35 training cycles to repair networks with small  
 69 damages. It is important to note that although iterative-rescue constantly damages and repairs networks,  
 70 the repair operation doesn't revive any of the damaged nodes, ie the damaged nodes and its weights  
 71 remain 0. We stress that the constant perturbation and repair strategy allows us to reach a favorable  
 72 'space' in the networks' loss manifold that contains high performing, more resilient, sparser networks.

73 As the iterative process of damage and repair always enabled the fast recovery of a damaged network  
 74 (irrespective of the number of damaged units), this was surprising to us and we were interested in  
 75 determining if the loss landscape manifold had 'special' geometric features that enabled this rescue.

76 To map geometric features of a neural networks' loss landscape, we formally conceptualize the iterative  
 77 damage-repair paradigm as a dynamical system that involves the application of a damage and repair  
 78 operator ( $r$ ) on a neural network ( $w$ ).

79 We define  $w$  to be a feed-forward neural network with  $n$  nodes and  $N$  total connections.

$$w = [\vec{w}_1, \vec{w}_2, \dots, \vec{w}_i, \dots, \vec{w}_n]$$

80 Here,  $\vec{w}_i$  is the set of connections made by node  $i$  with the previous layer in the network. By definition,  $\vec{w}_i = \phi$ , if node  $i$  is in the first layer.  
 81 We also have:

$$\sum_{i=1}^n \text{Dim}(\vec{w}_i) = N \text{ and } w \in \mathbb{R}^N$$

82 To damage a neural network, we define a damage operator  $D_i$ , that damages node  $i$  in the network.

$$D_i : \mathbb{R}^N \longrightarrow \mathbb{R}^N$$

$$w' = D_i(w) \begin{cases} \vec{w}'_i = \mathbf{0}, \\ \vec{w}'_j = \vec{w}_j \end{cases}$$

83 To repair a neural network, we define a rescue operator  $r_{\{i,j\}}$ . Here  $\{i,j\}$  refers to the set of damaged nodes. The rescue operator forces the  
 84 network to descend the loss manifold, while fixing nodes within the set and their connections to zero. Rescue of the network is achieved by  
 85 performing a constrained gradient descent on the networks' loss manifold.

$$r_{\{i,j\}} : \mathbb{R}^N \longrightarrow \mathbb{R}^N$$

$$w' = r_{\{i,j\}}(w) \begin{cases} \vec{w}'_i = \mathbf{0}, \vec{w}'_j = \mathbf{0}, \\ \vec{w}'_k = \vec{w}_k - \eta \frac{\partial L}{\partial \vec{w}_k} \end{cases}$$

86 where,  $\eta$  is the gradient step-size and  $\frac{\partial L}{\partial \vec{w}_k}$  is the gradient of the loss function of the neural network along  $\vec{w}_k$

87 A damage-repair sequence involves the application of a damage operator followed by a repair operator.

$$w' = r_{\{i\}}(D_i(w))$$

88 A stochastic damage-repair sequence involves the random sampling of a damage operator from  $\mathbf{D}$ , followed by the application of an  
 89 appropriate repair operator (ensuring that gradient descent is performed on remaining undamaged nodes).

$$w' = r_{\{i\}}(D_i(w)) \text{ where, } i \sim P(i) = \frac{1}{n}$$

90 We define a random variable  $\mathbf{D}$  to sample an operator  $D_i$  from the set of all possible damage operators =  $\{D_i : i \in \{1, \dots, n\}\}$ . An iterative  
 91 damage-repair sequence is the repeated application of a random damage operator  $\mathbf{D}$  coupled with a deterministic repair operator  $r_{\{i,j,k,\dots\}}$ ,  
 92 that ensures all damaged nodes maintain a zero edge-weight, while other weights are plastic. Here, we show the long-hand and short-hand  
 93 notation for the iterative application of damage-repair operators.

$$w' = r_{\{i,j,k\}}(D_k(r_{\{i,j\}}(D_j(r_{\{i\}}(D_i(w)))))) \text{ where, } i, j, k \sim P(i, j, k) = \frac{1}{n^3}$$

$$w' = (r \circ \mathbf{D})^3(w)$$

94 We hypothesize that  $\exists$  an open set of networks  $U$ , that constitutes an **invariant set**, where:

$$\text{if } w \in U, \text{ then}$$

$$(r \circ \mathbf{D})^m(w) \in U \quad \forall m$$

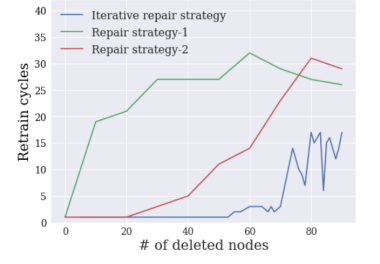


Figure 3: Iterative damage-repair strategy enables swift recovery of the network, when compared to batch-recovery of the network achieved by either strategy 1 or 2.

95 We also claim that the invariant set  $U$  is **path-connected**, ie given any two points from this topological space ( $U$ ), there exists a path ( $\gamma$ )  
 96 that connects the two points, starting at one point and ending at the other.

For any two points,  $w_1$  and  $w_2 \exists \gamma : [0, 1] \rightarrow U$ , such that:  
 $\gamma(0) = w_1, \gamma(1) = w_2$  and  $\gamma(t) \in U \forall t \in [0, 1]$

97 Our numerical results strongly suggests the presence of an invariant, path-connected topological space  $U$  in the neural networks' loss  
 98 manifold. In our experiments, the invariant, path-connected set is a collection of trained networks, whose image corresponding to the  
 99 application of a damage and repair operator lies in the same set, visualized by the thick black arc (as shown in figure-4) obtained by tSNE  
 100 embeddings of the high-dimensional network ( $w$ ). We observe that iterative application of the damage-repair operator on a network sampled  
 101 from  $U$  results in a series of networks that belong to the same set  $U$ . This is observed in fig-4b & fig-4d. The red lines indicate damage of  
 102 network, while the green lines correspond to repair of damaged networks. This hints at the possibility that  $U$  is an invariant set. We also  
 103 interpolated between all pairs of networks sampled from  $U$  and observed that all the interpolated networks were present in  $U$  as well.

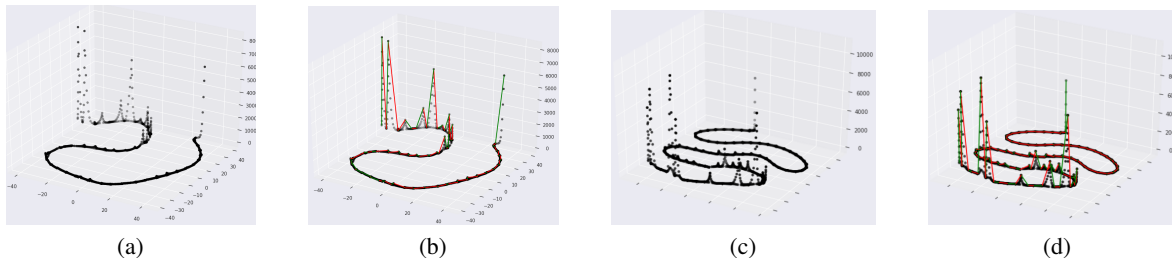


Figure 4: **Geometric features of the neural networks' loss landscape:** The x,y axes are tSNE embedding and z axes is the loss of the network (a,b) [tSNE] MLP with 1-hidden layer network (c,d) [tSNE] of MLP 2-hidden layers network (b,d) The green and red lines refer to repair and damage of networks respectively.

## 104 5 Discussion

105 In this paper, we address a pertinent question of how neural networks in the brain, or in engineered systems respond to damage of their units  
 106 and whether there exists efficient strategies to repair damaged networks. We observe a phase transition behavior as we incrementally delete  
 107 nodes from the neural network as the rate of decline of performance steeply increases after crossing a critical number of node deletions.  
 108 We discover that damaged networks can be rescued and the iterative damage-rescue strategy produces networks that are highly resilient  
 109 to perturbations, and can be rescued within a small number of training cycles. This is enabled by the putative presence of an invariant,  
 110 path-connected set in the networks' loss manifold. Although we have shown numerical results that strongly suggest the presence of invariant  
 111 sets in the loss manifold, our future work will focus on analytically proving the presence of these topological spaces in the loss manifold,  
 112 through the formalization presented in the paper, and the use of the Koopman operator machinery, amongst others.

## 113 References

- 114 [1] Carver Mead. "Neuromorphic electronic systems". In: *Proceedings of the IEEE* 78.10 (1990), pp. 1629–1636.  
 115 [2] Lindsey L Glickfeld and Shawn R Olsen. "Higher-order areas of the mouse visual cortex". In: *Annual review of vision science* 3  
 116 (2017), pp. 251–273.  
 117 [3] Mike Davies et al. "Loihi: A neuromorphic manycore processor with on-chip learning". In: *IEEE Micro* 38.1 (2018), pp. 82–99.  
 118 [4] Timothy D Hanks and Christopher Summerfield. "Perceptual decision making in rodents, monkeys, and humans". In: *Neuron* 93.1  
 119 (2017), pp. 15–31.  
 120 [5] Alison Duffy et al. "Variation in sequence dynamics improves maintenance of stereotyped behavior in an example from bird song".  
 121 In: *Proceedings of the National Academy of Sciences* 116.19 (2019), pp. 9592–9597.  
 122 [6] Chih-Hong Cheng, Georg Nührenberg, and Harald Ruess. "Maximum resilience of artificial neural networks". In: *International*  
 123 *Symposium on Automated Technology for Verification and Analysis*. Springer, 2017, pp. 251–268.  
 124 [7] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples". In: *arXiv preprint*  
 125 *arXiv:1412.6572* (2014).  
 126 [8] Ari S Morcos et al. "On the importance of single directions for generalization". In: *arXiv preprint arXiv:1803.06959* (2018).  
 127 [9] John Milnor. "On the concept of attractor". In: *The theory of chaotic attractors*. Springer, 1985, pp. 243–264.  
 128 [10] Morris W Hirsch, Stephen Smale, and Robert L Devaney. *Differential equations, dynamical systems, and an introduction to chaos*.  
 129 Academic press, 2012.