
Continual Learning with Group-wise Neuron Competition Appendix

Anonymous Author(s)

Affiliation

Address

email

1 In the following, we provide the Appendix as part of the supplementary material to the main paper.
2 **In Appendix Section A.**, we prove the alternative form for the back-propagation to make more
3 obvious at which part the gradient of the activation with respect to its arguments comes into use. **In**
4 **Appendix Section B.**, we give additional comment about applying group-wise neuron normalization
5 with nonlinear activation, providing additional results. **In Appendix Section C.**, we prove additional
6 result regarding performance comparison of existing normalization techniques and our GNC. **In**
7 **Appendix Section D.**, we give very simple pytorch code implementation of the proposed group-wise
8 neuron normalization that induces group-wise neuron competition (GNC).

9	A. Back-propagation Alternative Form	2
10	B. Additional Comment on Group-wise Neuron Competition	3
11	B.1. Relevance of the Order for the Operations Sparse Coding and GNC	3
12	B.2. Soft Sparse Coding Transform vs Hard Sparse Coding Transform	4
13	C. Additional Comparative Results with Normalization Methods	5
14	D. Group-wise Neuron Competition Implementation	6
15	D.1. Group-wise Neuron Competition (GNC) pytorch Implementation	6
16	D.2. Hard Sparse Coding Transform (sT) pytorch Implementation	6

17 A. Back-propagation Alternative Form

18 In the main text of the paper, we analyzed the influence of GNC that has during back-propagation, so
 19 we were interested in the gradient of the activation with respect to its arguments, *i.e.*, $\mathbf{g}^{L-1} = \frac{\partial \mathbf{a}^{L-1}}{\partial \mathbf{z}^{L-1}}$.
 20 In this appendix section, we give the full derivation of the equivalent form for the back-propagation,
 21 and pinpoint where $\frac{\partial \mathbf{a}^{L-1}}{\partial \mathbf{z}^{L-1}}$ comes to play.

22 Consider a neural network that maps inputs $\mathbf{x} \in \mathbb{R}^{N_0}$ to outputs $\mathbf{y} \in \mathbb{R}^{N_L}$. The forward pass reads
 23 as:

$$\begin{aligned} \mathbf{a}^0 &= \mathbf{x}, \\ \mathbf{z}^l &= \mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l, \quad l \in \{1, \dots, L\}, \\ \mathbf{a}^l &= \sigma(\mathbf{z}^l), \quad l \in \{1, \dots, L\}. \end{aligned} \quad (1)$$

24 Here, $\mathbf{W}^l \in \mathbb{R}^{N_l \times N_{l-1}}$ is the weight matrix of layer l (note that if we have convolution, we can
 25 express it similarly as a matrix product multiplication), \mathbf{b}^l the corresponding bias vector. Where N_l
 26 denotes the dimension of the layer l . The activation function is denoted by σ , it processes the layer's
 27 weighted sum \mathbf{z}^l to the layer's output \mathbf{a}^l .

28 Given data set, we define the training process of a neural network as an optimization over an objective
 29 function, usually called *loss*, that we denote as \mathcal{L} . In other word, we want to find the parameters
 30 (weights and biases) that minimize the cost function. The training is commonly done using a gradient
 31 based rule. Therefore, the update relies on the gradient of \mathcal{L} with respect to weight \mathbf{W}^l and the bias
 32 \mathbf{b}^l , that is it relies on $\nabla_{\mathbf{W}^l} \mathcal{L}$ and $\nabla_{\mathbf{b}^l} \mathcal{L}$, respectively. Back-propagation facilitates the computation of
 33 these gradients, and makes use of the chain rule to back-propagate the prediction error through the
 34 network [7]. We express the error vector at layer l as:

$$\delta^l = \nabla_{\mathbf{a}^l} \mathcal{L}, \quad (2)$$

35 and further use it to express the gradients as

$$\begin{aligned} \nabla_{\mathbf{W}^l} \mathcal{L} &= \delta^l (\mathbf{a}^{l-1})^T, \\ \nabla_{\mathbf{b}^l} \mathcal{L} &= \delta^l. \end{aligned} \quad (3)$$

36 The output layer's error is simply given by

$$\delta^L = \nabla_{\mathbf{a}^L} \mathcal{L} \odot \frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L}, \quad (4)$$

37 where \odot denotes the Hadamard element-wise product. Subsequent earlier layer's error are computed
 38 with

$$\delta^l = (\mathbf{W}^{l+1})^T \delta^{l+1} \odot \frac{\partial \mathbf{a}^l}{\partial \mathbf{z}^l}, \quad l \in \{1, \dots, L-1\}. \quad (5)$$

39 where A usual parameter update takes on the form

$$(\mathbf{W}^l)_{\text{new}} = \mathbf{W}^l - \beta \nabla_{\mathbf{W}^l} \mathcal{L}, \quad (6)$$

40 where β is a positive learning rate.

41 B. Additional Comment on Group-wise Neuron Competition

42 B.1. Relevance of the Order for the Operations Sparse Coding and GNC

Table 1: Results using the CIFAR100 data set in the single pass continual learning regime. At the penultimate layer we used two different operating orders (i) GNC followed by sT and (ii) sT followed by GNC, showing that the difference between the two is negligible.

ST ACTIVATION FOLLOWED BY GNC		GNC FOLLOWED BY ST ACTIVATION	
RA	BTI	RA	BTI
61.47 ± 1.07	-04.34 ± 1.34	61.30 ± 0.56	-04.07 ± 0.94

43 **Existing Methods (Normalization Followed by Nonlinear Activation).** In the existing normaliza-
 44 tion techniques like batch norm [4], group norm [9] and layer norm [2], it is common that we first
 45 apply the normalization and then the nonlinear activation, *i.e.*, $\mathbf{a}^l = \text{ReLU}(\text{bn}(\mathbf{z}^l))$, where $\text{bn}(\mathbf{z}^l)$ is
 46 the batch norm in this example. It turn out that when we use the existing normalization techniques
 47 [4, 9, 2] it makes a big difference whether (i) we first apply the nonlinear activation and then the apply
 48 the normalization, or (ii) we first apply the normalization and then apply the nonlinear activation. In
 49 the later case, using normalization with learnable parameters [4, 9, 2], we might have degradation of
 50 performance. We mention this in the paper as well as give additional results in Section C.

51 **GNC.** In our approach we found out that it dose not play a role if we first apply the nonlinear activation
 52 (at least in the case of sT) and then the GNC, or the other way around. In the following we will revisit
 53 the definition about GNC and give explanation.

54 In the main text, we construct \mathbf{a}^l as $\mathbf{a}^l = [(\mathbf{a}_1^l)^T, (\mathbf{a}_2^l)^T, \dots, (\mathbf{a}_J^l)^T]^T$ and define the group-wise
 55 normalization as:

$$\mathbf{a}_j^l = v_{GNC}(\mathbf{z}_j^l) = \frac{\mathbf{z}_j^l}{\|\mathbf{z}_j^l\|}, \forall j \in \{1, \dots, J\}, \quad (7)$$

56 where \mathbf{z}_j^l are linear activation's that belong to group j and $\|\mathbf{z}_j^l\|$ is the ℓ_2 -norm of \mathbf{z}_j^l . We also
 57 mentioned that before the normalization, we can apply nonlinear activation, such as ReLU [1], so we
 58 have $\mathbf{a}_j^l = v_{GNC}(\text{ReLU}(\mathbf{z}_j^l))$, or with sparse coding transform sT [6] we have:

$$\mathbf{a}_j^l = v_{GNC}(\text{sT}(\mathbf{z}_j^l)). \quad (8)$$

59 Note that the group-wise neuron competition is about the max magnitude neuron response that should
 60 lay in the group with the highest sparsity. Therefore, when we apply the sT, we pronounce the sparsity
 61 (ether before or after) within the group of neuron responses, and this should help the group-wise
 62 neuron competition, since we remove (threshold out) some not-relevant "information".

63 In Tab. 1, we give results that evaluated the above line of tough. We experimented using the
 64 CIFAR100 data set in the task-aware single pass regime. We set the reservoir size of 200 and use
 65 a reservoir batch size of 20. As we can see in Tab. 1, the results suggest that the order of these
 66 operation dose not play a role. When we use the sT after the GNC, very slightly we have higher
 67 retained accuracy, while when we use the sT before the GNC, very slightly we have better backward
 68 transfer. So, the difference in results between the two are negligible.

69 B.2. Soft Sparse Coding Transform vs Hard Sparse Coding Transform

Table 2: Results using the CIFAR100 data set in the single pass continual learning regime. At the penultimate layer we used two different sparse coding transforms (sT) (i) "hard" sT followed by GNC and (ii) "soft" sT followed by GNC.

HARD ST ACTIVATION		SOFT ST ACTIVATION	
RA	BTI	RA	BTI
61.47 ± 1.07	-04.34 ± 1.34	60.86 ± 0.37	-04.29 ± 0.18

70 **Hard Sparse Coding Transform.** We define the "hard" sparse coding transform (sT) [6], as,

$$\text{sT}(\mathbf{z}_j^l) = \mathbf{z}_j^l \odot \mathbf{t}_j^l, \quad (9)$$

71 where $t_{j,i}^l = 1$ if $|z_{j,i}^l| \geq \lambda_j$, and $t_{j,i}^l = 0$ otherwise, while $\lambda_j = \frac{J}{N_l} \sum_{i=1}^{N_l/J} |z_{j,i}^l|$ is a sparsifying
 72 threshold. We have \mathbf{z}_j^l as the linear activation's that belong to group $j \in \{1, 2, \dots, J\}$, $\mathbf{z}^l =$
 73 $[(\mathbf{z}_1^l)^T, (\mathbf{z}_2^l)^T, \dots, (\mathbf{z}_J^l)^T]^T \in \mathbb{R}^{N_l}$.

74 **Soft Sparse Coding Transform.** The "soft" sparse coding transform (sT) is similar to the "hard"
 75 sparse coding transform, *i.e.*:

$$\text{sT}_{\text{soft}}(\mathbf{z}_j^l) = \text{sign}(\mathbf{z}_j^l) \odot \max(|\mathbf{z}_j^l| - \lambda_j, 0), \quad (10)$$

76 where we define the sparsifying threshold as above. The difference between (10) and (9) is that in
 77 (10) in addition to the tresholding to zero the magnitude of the nonzero values are also reduced by
 78 the sparsifying threshold λ_j .

79 Although (10) might have better bias/variance properties in classical signal processing applications,
 80 we empirically found out that (9) works better in our continual learning regimes. In Tab. 2, we
 81 verify this and give comparative results between "hard" and "soft" sT. We experimented using the
 82 CIFAR100 data set in the task-aware single pass regime. We set the reservoir size of 200 and use a
 83 reservoir batch size of 20. We use sT ("hard" or "soft") followed by GNC. As we can see in Tab. 2,
 84 the "hard" sT has an edge over "soft" sT, with higher retained accuracy of .6%, while the backward
 85 trasferability score is very similar between the "hard" sT and the "soft" sT.

86 C. Additional Comparative Results with Normalization Methods

Table 3: Results using the CIFAR100 data set in the single pass continual learning regime. We show the performance when we use BN [4], GN [9], LN [2] and CN [5] at the penultimate layer. We experimented with ReLu before and after the normalization.

NORM	LINEAR		ReLU		NORMALIZATION	
	ACTIVATION FOLLOWED		ACTIVATION FOLLOWED		FOLLOWED BY	
	RA	BTI	RA	BTI	ReLU	BTI
BN	61.01 ± 0.984	-04.30 ± 0.319	56.90 ± 0.05	-06.80 ± 0.52	60.65 ± 0.92	-06.93 ± 0.54
GN	59.47 ± 0.134	-02.20 ± 0.026	56.49 ± 0.50	-04.50 ± 0.020	58.73 ± 0.67	-04.66 ± 2.89
LN	60.46 ± 3.76	-02.46 ± 0.035	58.55 ± 0.58	-05.53 ± 0.013	58.14 ± 1.72	-04.93 ± 2.08
CN	60.18 ± 0.089	-04.50 ± 2.65	58.26 ± 0.68	-04.50 ± 0.20	58.88 ± 1.65	-07.75 ± 3.08

Table 4: Results using the CIFAR100 data set in the single pass continual learning regime. We show the performance when we different activation functions.

f_n	ACTIVATION FOLLOWED		ACTIVATION WITHOUT	
	RA	BTI	RA	BTI
LIN	60.14 ± 1.80	-06.22 ± 0.37	55.89 ± 0.38	-08.09 ± 1.06
ReLU	54.61 ± 0.98	-10.04 ± 0.00	53.50 ± 0.05	-09.81 ± 0.37
ST _{soft}	60.86 ± 0.37	-04.29 ± 0.18	53.63 ± 0.75	-09.25 ± 0.61
ST	61.76 ± 2.03	-03.99 ± 1.09	54.74 ± 0.80	-09.18 ± 0.03

87 In Tab. 3, we show the results in the single-pass regime using the CIFAR100 data set. We compare
88 against existing normalization techniques, including including BN [4], GN [9], LN [2] and CN [5].
89 We used equivalent convolution neural network as in [3], and we apply the normalization at the
90 penultimate layer. We use reservoir of size 200 and reservoir batch size of 20. In the GN, CN and
91 GNC, we use 32 groups with 10 neuron activation's each for fair comparison.

92 As we can see in Tab. 3 on the left side on the top, when we do not apply non-linearity, BN, GN, LN
93 and CN have good performance in both RA and BTI. While when we used ReLu activation followed
94 by BN, GN, LN or CN we have lower retained accuracy (RA) and backward transferability (BTI).
95 That is the retained accuracy over BN, GN, LN and CN drops by an average of around 3.5%. In the
96 case when we preform the operations in the reversed order, that is BN, GN, LN or CN followed by
97 ReLu the results have a bit higher retained accuracy compared to the previous ones. However, we
98 still have lower retained accuracy compared to the linear activation by an average of around 2%. We
99 also observed that when we used a sT activation followed by BN, GN, LN or CN, we had a strong
100 degradation of performance (during training, the objective value got stuck at a certain value and was
101 not able to improve/learn for the experimented normalization methods). We also tried out instance
102 normalization (IN) [8] in a equivalent setup as before, but similarly as above, we were not able to
103 make the model to learn.

104 In Tab. 4 on the bottom, we show the influence of the used nonlinear activation function before the
105 GNC. As activation function we experimented with linear, ReLu, and sT. As reported in the main
106 text, we can see that not all combinations of activation function followed by GNC are useful.

107 In Tab. 3 and Tab. 4, we can also see that our approach with "hard" sT and GNC has the best
108 performance.

109 **D. Group-wise Neuron Competition Implementation**

110 **D.1. Group-wise Neuron Competition (GNC) pytorch Implementation**

111 In the following we prove a **pytorch** implementation of the group-wise normalization that induces
112 group-wise neuron competition (GNC).

```
def gnc(a, group_size, num_of_groups):  
    # GNC implementation  
    # assuming shape_1 = group_size*num_of_groups  
    shape_0, shape_1 = a.shape  
    a = a.reshape(shape_0, group_size, num_of_groups)  
    a = torch.nn.functional.normalize(a, p=2, dim=1).reshape(shape_0, shape_1)  
  
    return a
```

113 We note that before applying the GNC, we also apply "hard" sparse coding transform (sT) (for more
114 details please see Appendix Section B.2.).

115 **D.2. Hard Sparse Coding Transform (sT) pytorch Implementation**

116 In the following, we give the **pytorch** implementation for the "hard" sparse coding transform.

```
def st(z, group_size, num_of_groups):  
    # sT implementation  
    shape_0, shape_1 = z.shape  
    # compute the threshold per group of neuron activation  
    r = z.reshape(shape_0, group_size, num_of_groups).detach().clone()  
    tr = torch.abs( r ) >= torch.mean( torch.abs(r), 1).unsqueeze(1)  
    # set to zero the components below the threshold  
    a = z * tr.float().reshape(shape_0, shape_1)  
  
    return a
```

References

- [1] Abien Fred Agarap. Deep learning using rectified linear units (relu). *CoRR*, abs/1803.08375, 2018.
- [2] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450, 2016.
- [3] G. Gupta, Karmesh Yadav, and L. Paull. Look-ahead meta learning for continual learning. In *NeurIPS*, 2020.
- [4] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 448–456. PMLR, 07–09 Jul 2015.
- [5] Quang Pham, Chenghao Liu, and Steven HOI. Continual normalization: Rethinking batch normalization for online continual learning. In *International Conference on Learning Representations*, 2022.
- [6] Saiprasad Ravishankar and Yoram Bresler. Sparsifying transform learning with efficient optimal updates and convergence guarantees. *IEEE Transactions on Signal Processing*, 63(9):2389–2404, may 2015.
- [7] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986. Publisher: Nature Publishing Group.
- [8] Dmitry Ulyanov, Andrea Vedaldi, and Victor S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *CoRR*, abs/1607.08022, 2016.
- [9] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.