

GloNets: Globally Connected Neural Networks

Antonio Di Cecco¹[0000-0002-9070-4663], Carlo Metta²[0000-0002-9325-8232],
 Marco Fantozzi³[0000-0002-0708-5495], Francesco
 Morandin³[0000-0002-2022-2300], and Maurizio Parton¹[0000-0003-4905-3544]

¹ University of Chieti-Pescara, Italy

² ISTI-CNR, Pisa, Italy

³ University of Parma, Italy

Abstract. Deep learning architectures suffer from depth-related performance degradation, limiting the effective depth of neural networks. Approaches like ResNet are able to mitigate this, but they do not completely eliminate the problem. We introduce Globally Connected Neural Networks (GloNet), a novel architecture overcoming depth-related issues, designed to be superimposed on any model, enhancing its depth without increasing complexity or reducing performance. With GloNet, the network’s head uniformly receives information from all parts of the network, regardless of their level of abstraction. This enables GloNet to self-regulate information flow during training, reducing the influence of less effective deeper layers, and allowing for stable training irrespective of network depth. This paper details GloNet’s design, its theoretical basis, and a comparison with existing similar architectures. Experiments show GloNet’s self-regulation ability and resilience to depth-related learning challenges, like performance degradation. Our findings suggest GloNet as a strong alternative to traditional architectures like ResNets.

Keywords: Neural Networks · Deep Learning · Skip Connections

1 Introduction

Deep learning’s success in AI is largely due to its hierarchical representation of data, with initial layers learning simple features and deeper ones learning more complex, nonlinear transformations of these features [3]. Increasing depth should enhance learning, but sometimes it leads to performance issues [4,9]. Techniques like normalized initialization [15,9,21,10] and normalization layers [14,1] enable up to 30-layer deep networks, but performance degradation persists at greater depths without skip connections. This issue, detailed in the original ResNet paper [11], stems from the fact that learning identity maps is not easy for a deeply nonlinear layer. ResNet idea is to focus on learning nonlinear “residual” information, with a backbone carrying the identity map. This brilliant solution has been key in training extremely deep networks that, when weight-sharing and batch normalization are used, can scale up to thousands of layers.

Deeper neural networks should not experience performance degradation. Theoretically, a deeper network could match the performance of an n -layer network

by similarly learning features $\mathcal{G}_1, \dots, \mathcal{G}_n$ in its initial layers, then minimizing the impact of additional layers. With this ability to *self-regulate*, such a network could effectively be “infinitely deep”. However, even with ResNet architectures, performance degradation persists beyond a certain depth, see Figure 5b or [7]. This issue can be due to various factors, see for instance [9,23,17], and may partly arise from the inability of modern architectures to self-regulate their depth. Our paper introduces a novel technique to enable self-regulation in neural network architectures, overcoming these depth-related performance challenges.

Novel Contributions. The main contribution of this paper is introducing and testing GloNet, an explainable-by-design layer that can be superimposed on any neural network architecture, see Section 2. GloNet’s key feature is its capacity to self-regulate information flow during training. It achieves this by reducing the influence of the deepest layers to a negligible level, thereby making the training more stable, preventing issues like vanishing gradients, and making the network trainable irrespective of its depth, see Section 5.

This self-regulation capabilities of GloNet lead to several significant benefits:

1. **Faster training:** GloNet trains in half the ResNet time while achieving comparable performance. Beyond the depth threshold where ResNet begins to degrade, GloNet trains in less than half the time and outperforms ResNet.
2. **ResNet alternative:** The inability of ResNet-based architectures to self-regulate depth makes GloNet a preferable option, particularly for applications requiring very deep architectures.
3. **No NAS needed:** GloNet networks inherently find their effective depth, eliminating the need for computationally expensive Network Architecture Search methods to determine optimal network depth.
4. **More controllable efficiency/performance trade-off:** Layers can be selectively discarded to boost efficiency, allowing a controlled trade-off between efficiency and performance, optimizing the network for specific requirements.

2 Notation and Model Definition

A feedforward neural network is described iteratively by a sequence of L blocks:

$$\mathbf{x}_{\ell+1} = \mathcal{G}_\ell(\mathbf{x}_\ell), \quad \ell = 0, \dots, L - 1, \quad (1)$$

where \mathbf{x}_0 denotes the input vector, and $\mathbf{x}_{\ell+1}$ is the output from the ℓ -th block. In this context, a “block” is a modular network unit, representing a broader concept than a traditional “layer”. Each block function \mathcal{G}_ℓ typically merges a non-linearity, such as ReLU, with an affine transformation, and may embody more complex structures, like the residual blocks in ResNet.

At the end of the sequence (1), a classification or regression head \mathcal{H} is applied to \mathbf{x}_L . For instance, a convolutional architecture could use a head with average pooling and a fully connected classifier. The fundamental principle in deep learning is that $\mathcal{G}_0, \dots, \mathcal{G}_{\ell-1}$ hierarchically extract meaningful features from the input \mathbf{x}_0 , that can then be leveraged by computing the output of the network:

$$\text{output} = \mathcal{H}(\mathbf{x}_L).$$

When the blocks in (1) are simple layers like an affine map followed by a non-linearity (this description comprises, for instance, fully connected and convolutional neural networks), all features extracted at different depths are exposed to the head by a single feature vector \mathbf{x}_L that has gone through several non-linearities. This fact leads to several well-known drawbacks, like vanishing gradients or difficulty in learning when the task requires more direct access to low-level features. When using ReLU and shared biases, some low-level information could actually be destroyed. Several excellent solutions have been proposed to these drawbacks, like for instance residual networks [11,12], DenseNets [13], and preactivated units with non-shared biases [17].

We propose an alternative solution: a modification to (1), consisting of a simple layer between the feature-extraction sequence and the head, computing the sum of every feature vector. The architecture is designed to receive information uniformly from all parts of the network, regardless of their level of abstraction:

$$\begin{cases} \mathbf{x}_{\ell+1} = \mathcal{G}_\ell(\mathbf{x}_\ell), & \ell = 0, \dots, L-1 \\ \mathbf{x}_{L+1} = \sum_{\ell=1}^L \mathbf{x}_\ell = \sum_{\ell=0}^{L-1} \mathcal{G}_\ell(\mathbf{x}_\ell) \\ \text{output} = \mathcal{H}(\mathbf{x}_{L+1}) \end{cases} \quad (2)$$

When feature vectors have different dimensions, adaptation to a common dimension is required before the sum, as happens in ResNet. If $\mathbf{x}_\ell \in \mathbb{R}^{n_\ell}$, one can use embeddings in $\mathbb{R}^{\max\{n_\ell\}}$ to maximally preserve information, and embeddings or projections to \mathbb{R}^{n_L} to maintain the same parameters for the head. We refer to the additional layer in (2) as a *GloNet layer*, because all the features \mathbf{x}_ℓ that without GloNet would be preserved only “locally” up to the next $\mathbf{x}_{\ell+1} = \mathcal{G}_\ell(\mathbf{x}_\ell)$, appear now in the “global” feature vector $\mathbf{x}_{L+1} = \sum_{\ell=1}^L \mathbf{x}_\ell$ as summands.

Remark 1 (Theoretical support for GloNet self-regulation). GloNet provides skip connections solely to the head, and intermediate blocks are not required to learn a residual map, see (2) and Figure 1. This ensures direct and simultaneous back-propagation pathways from each block, enabling uniform information distribution across the network to the head. Due to SGD-like training’s preference for shorter paths [25,26,16], *GloNet is expected to accumulate information mostly in the initial blocks rather than the latter ones*, by reducing the influence of the deepest layers to a negligible level. Consequently, GloNet self-regulates its depth during training, rendering it akin to an “infinitely deep” architecture. Empirical evidence supporting this claim is presented in Section 5.

Remark 2 (Explainability-by-design). Given the linearity of the GloNet layer, the global feature vector \mathbf{x}_{L+1} provides the contribution that each layer makes to the neural network’s prediction. In a feedforward neural network, GloNet enables the analysis of an ensemble of networks represented by the outputs of each network block. This ensemble, comprising blocks $(\mathcal{G}_\ell)_{\ell=0, \dots, k}$, functions as individual neural networks, with GloNet integrating their outputs through a linear (currently unweighted) combination. This architecture allows each sub-network to specialize in learning features at varying levels of granularity, from

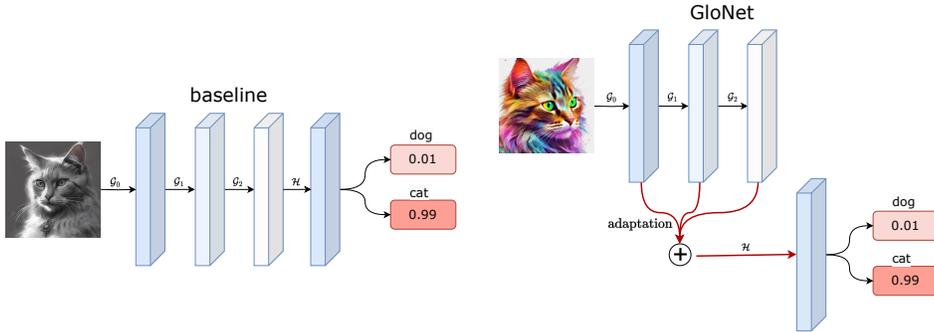


Fig. 1: On the left, a 3-blocks neural network followed by a classification head. On the right, the same architecture with GloNet modification in red.

low-level in early blocks to more complex, large-scale features in later blocks. The linear nature of the GloNet layer facilitates the attribution of importance scores to these features, effectively creating an ‘explainable-by-design’ tool [8].

3 Related Work

In a ResNet with an activation-free backbone, also known as ResNetv2 [12], blocks \mathcal{G}_ℓ are defined as $\text{Id} + \mathcal{F}_\ell$, where Id is the skip connection and \mathcal{F}_ℓ is the “residual block”, computing two times $\text{Conv} \circ \text{ReLU} \circ \text{BN}$, where BN is batch normalization. Unrolling the ResNetv2 equation $\mathbf{x}_{\ell+1} = \mathbf{x}_\ell + \mathcal{F}_\ell(\mathbf{x}_\ell)$ from any block output \mathbf{x}_ℓ (see [12, Equation 4]) gives:

$$\mathbf{x}_L = \mathbf{x}_\ell + \sum_{i=\ell}^{L-1} \mathcal{F}_i(\mathbf{x}_i) \quad (3)$$

This equation shows that ResNet, much like GloNet, passes the output \mathbf{x}_ℓ of each block directly to the head. However, a key distinction lies in how the head accesses these outputs: ResNet requires distinct pathways for simultaneous access to different outputs, whereas GloNet’s head achieves simultaneous access to each output via the GloNet layer. This unique capability of GloNet may contribute to its additional features compared to ResNetv2, as explored in Section 5.

Moreover, GloNet is faster than ResNet (not requiring batch normalization), and can be seen as an ensemble computing the sum of models of increasing complexity, giving an explainable-by-design model (differently from ResNet).

Unlike DenseNet [13], aggregating each block with all subsequent blocks through concatenation, GloNet connects only to the last block, with summation for aggregation. This approach avoids the parameter explosion given by concatenation in DenseNet, and maintains the original complexity of the model.

Finally, GloNet can be viewed as a network with early exits at every block, adapted and aggregated before the head. See [24,18,19] for the early exit idea.

4 Implementing GloNet

Implementing GloNet within a certain architecture may not always be as straightforward as described in Section 2. In this Section we describe how GloNet can be implemented in common scenarios.

GloNet and Skip Connections. If the original architecture includes skip connections (such as ResNet or DenseNet), these should be removed and replaced with the GloNet connection. Otherwise, putting GloNet on top of skip connections, training would not converge as we would be adding the identity to the output multiple times. Note that GloNet provides only skip connections to the GloNet layer, and does not ask the blocks to learn a residual map.

GloNet and Batch Normalization. Batch normalization plays a major role in enhancing and stabilizing neural network training by normalizing the inputs of each block. Its positive impact is widely acknowledged, though the specific mechanisms of its benefits are still debated [14,20]. Despite these advantages, batch normalization poses challenges, particularly in its interaction with GloNet. GloNet is designed to dynamically regulate the outputs of different blocks, based on their contribution to the task. It makes negligible the outputs of deeper blocks, a strategy that conflicts with the objectives of batch normalization, which strives to maintain a consistent mean and variance for block inputs. Consequently, batch normalization, and similarly layer normalization, should be removed prior to the GloNet layer’s aggregation. GloNet introduces an alternative form of regularization, which, as we demonstrate in Section 5, is capable of achieving comparable performances without the need for batch normalization.

GloNet into Residual Networks. Architectures using residual blocks feature both skip connections and normalization. Once skip connections and normalizations are removed from a ResNetv2 block computing $\text{Id} + \text{affine map} \circ \text{ReLU} \circ \text{BN} \circ \text{affine map} \circ \text{ReLU} \circ \text{BN}$, one is left with two simpler blocks $\text{affine map} \circ \text{ReLU}$, and each of those blocks can potentially be aggregated into the GloNet layer. In this case, one ResNetv2 block corresponds to two simpler blocks. This is what we do in this paper, and for this reason when GloNet has n blocks, its equivalent ResNetv2 architecture has $n/2$ blocks.

GloNet into Vision Transformers. When using more complex architectures like transformers, several different choices can be made, each one potentially affecting the final performance of the GloNet-enhanced model. In this initial exploration of GloNet, we propose a straightforward integration with a Vision Transformer (ViT) [6] adapted to CIFAR-10. The image is segmented into 4x4 patches, its class encoded, and then concatenated with the patch encoding and a positional embedding. This series is then fed into a cascade of n encoders with 4 attention heads each, which are accumulated into a GloNet layer, and passed to a classification head. In our experiment, we compared $n = 4, 5$ and 6.

5 Experiments

In this section we provide experiments supporting the core claims of our paper, as stated in the Introduction. In particular, we focus on showing that GloNet trains

much faster than ResNet, that GloNet performances are on par with ResNet’s ones, that GloNet can self-regulate its depth, that GloNet does not need batch normalization, and that GloNet is virtually immune to depth-related problems. All the experiments can be reproduced using the source code provided at [5].

SGEMM Fully Connected Regression. We experimented with a regression task from the UCI repository [2,22], focused on predicting the execution time of matrix multiplication on an SGEMM GPU kernel. See [2,22] for details on this task and SGEMM dataset.

Since GloNet has skip connections, to obtain a fair comparison we used a ResNetv2-like baseline. For comparison, we used also a vanilla baseline, identical to the ResNetv2 baseline but without skip connections. Moreover, since GloNet does not use batch normalization (in fact, GloNet self-regulation capabilities can be tampered by normalization, see Section 4), we also experimented with a vanilla and a ResNetv2 baseline with the BN layer removed. GloNet and the corresponding baselines (denoted by vanilla, ResNetv2, vanilla-no-BN, and ResNetv2-no-BN in figures) have a similar amount of parameters, the only difference being given by the trainable BN parameters. All blocks have 16 units. All models starts with a linear layer mapping the 14-dimensional input to \mathbb{R}^{16} , and ends with the head, a linear layer with 1 unit. GloNet models have an additional GloNet layer before the head *with no additional parameters*. The number of blocks ranges in [10, 24, 50, 100, 200], respectively (halved for the ResNets because every block is twice the layers of the corresponding non-ResNet model).

For training, we used MSE loss, L^2 -regularization with a coefficient of 10^{-5} (we also tried 10^{-4} without improvements), Adam optimizer with a batch size of 1024, learning rate set to 0.01, He normal initializer for weights, and zero initializer for biases. We trained all models for 200 epochs, the point at which baselines plateaued, potentially favoring them over GloNet. The first thing to notice is that with GloNet *training takes almost half or less than half the time of ResNet*, see Table 1. This is because GloNet, differently from ResNet, does not need batch normalization.

Table 1: Average epoch’s training time in seconds at different depths, for GloNet and its equivalent ResNetv2 baseline, on SGEMM regression task.

architecture	depth						
	10	24	50	100	200	600	1000
ResNetv2	42	65	96	147	231	410	675
Glonet	27	39	55	86	121	175	289

After 200 training epochs, we compared the best test errors and learning curves across different block configurations, see Figure 2 for learning curves and Table 2 for best test errors. At 200 blocks, GloNet surpassed ResNet in both best test error and learning curve shape. See the caption of Figure 2 for details on the results of this experiment.

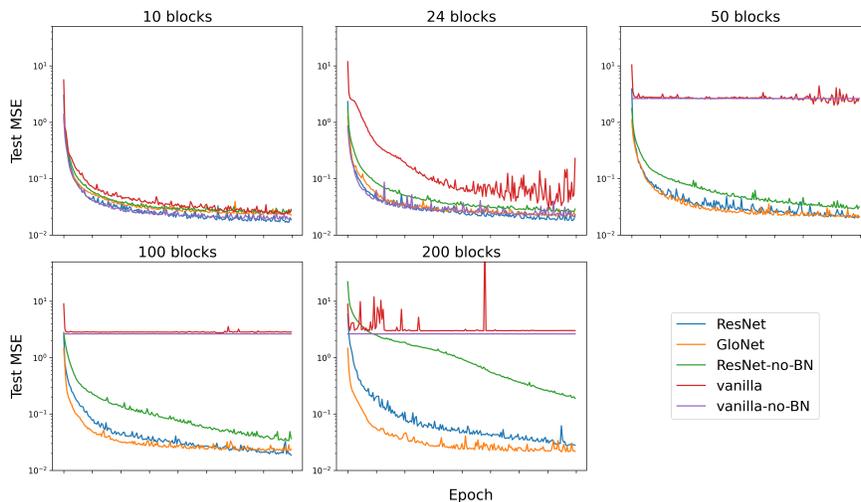


Fig. 2: Test errors while training for 200 epochs GloNet and four corresponding baselines at different depths on SGEMM. Starting at 24 blocks, GloNet and ResNetv2 outperform vanilla, as expected due to skip connections. At 50 and 100 blocks, both ResNetv2 and GloNet show depth resilience, with ResNetv2 slightly leading in best test error. ResNetv2-no-BN exhibits significantly higher error, highlighting ResNetv2’s reliance on batch normalization. At 100 blocks GloNet curve shows a markedly higher curvature than ResNetv2. At 200 blocks, GloNet outperforms ResNetv2 in best test error and curve shape.

The training shapes in Figure 2 suggested a unique aspect of GloNet not present in the ResNetv2 baseline. GloNet’s training was unaffected by the increasing depth, as shown by the shape of the learning curve that remained consistent whether the network had 10, 24, 50, 100, or 200 blocks. On the contrary, ResNet learning curve became flatter when depth is increasing.

To further explore this feature, GloNet was tested with even deeper models (600 and 1000 blocks), and compared against the corresponding ResNetv2 baseline. Even at these substantial depths, GloNet’s learning curve maintained its shape, as shown in Figure 3. Moreover, GloNet’s performance remained stable across these varying depths, maintaining a best test error of around 0.02 regardless of the number of blocks (10, 24, 50, 100, 200, 600, or 1000), see Table 2.

In contrast, ResNet’s showed a clear decline as the network depth increased. While its best test error remained around 0.02 up to 200 blocks, this error increased to 0.04 and almost 0.05 at 600 and 1000 blocks, respectively, see again Table 2. As happened with 100 and 200 blocks, the ResNet learning curve was flatter, diverging from GloNet’s more consistent curve shape as depth increased.

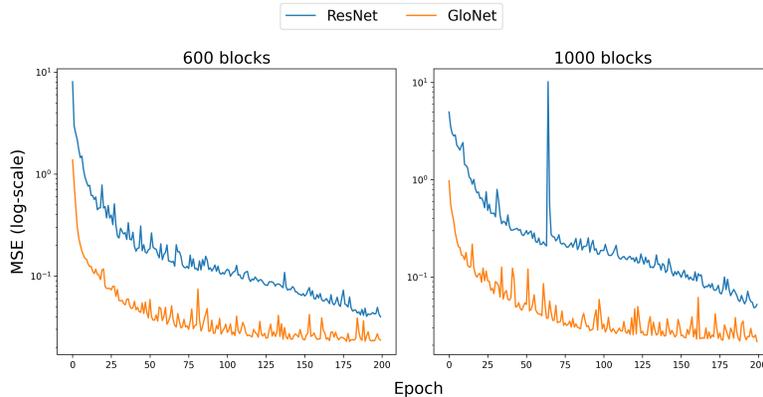


Fig. 3: Test MSE learning curves of GloNet (orange) with 600 blocks (left) and 1000 blocks (right), against the equivalent ResNetv2 baseline (blue). Comparing with Figure 2 shows these substantial depths severely degrade ResNet performance, but do not affect at all GloNet performance.

GloNet accumulates information in the first few blocks and uses only the required capacity for a specific task and architecture, leading to minimal output from subsequent blocks, see Remark 1. This is not observed in baseline models with or without batch normalization, see Figure 4, and likely contributes to GloNet’s stable performance as network depth increases, in contrast to the degradation observed in the baseline models under similar depth conditions.

MNIST Fully Connected Classification. To confirm that GloNet automatic choice of optimal depth and GloNet training resilience to depth were not associated to the particular SGEMM regression task, we performed a series of experiments with identical architecture on a completely different task: image classification with MNIST. Although using fully connected architectures for image classification is generally not the best approach, with this task we have been able to significantly increase the number of input features, which theoretically could pose a greater challenge to models that are not very deep.

Table 2: Best test errors across different network depths.

architecture	depth						
	10	24	50	100	200	600	1000
vanilla	0.023	0.030	1.997	2.750	2.985	-	-
vanilla-no-BN	0.018	0.021	2.624	2.624	2.624	-	-
ResNetv2	0.018	0.019	0.020	0.019	0.027	0.040	0.048
ResNetv2-no-BN	0.024	0.026	0.029	0.033	0.189	-	-
Glonet	0.021	0.021	0.020	0.022	0.021	0.022	0.022

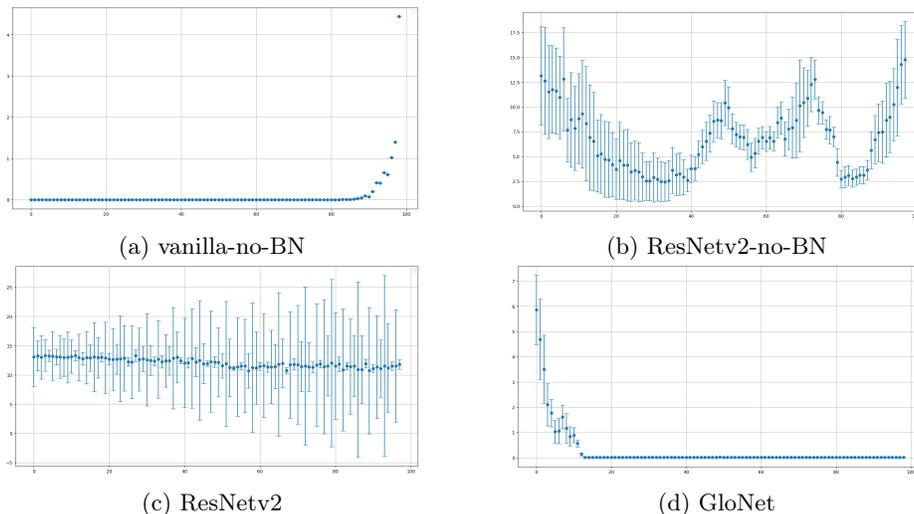


Fig. 4: Block’s outputs L1 norm, 100 blocks, 200 epochs, SGEMM. Mean and standard deviation from 5000 samples. GloNet (d) uses only the first 12 blocks, and achieves the least variance, indicating a more consistent learning of features. In ResNetv2 (c) all blocks give a similar contribution, as expected because of BN. ResNetv2-no-BN (b) uses all outputs in different ways, because of the residual maps mixed with the skip connections. The baseline vanilla-no-BN (a) shows the opposite behavior to GloNet, emphasizing late outputs and diminishing earlier ones, indicating shorter gradient paths to the last blocks.

The only difference from the architecture used in SGEMM is the head, which in this case is a fully connected layer followed by a SoftMax layer on 10 classes. We tested architectures with 6, 10, 24, 50, 60, 80, 100, and 200 blocks for GloNet and a convolutional baseline, halved for the corresponding ResNet baseline. Figure 5a shows that GloNet automatically chooses the optimal number of blocks. However, notice that in this case, differently from Figure 4c, ResNetv2 outputs show a decreasing shape. This is probably due to the trainable parameters of the batch normalization, that in this case are able to force a small mean and variance on the last blocks. This indicates that also with batch normalization the network struggles to self-regulate its depth. Figure 5b confirms GloNet resilience to an increasing depth, and shows a severe performance degradation of ResNetv2 when depth goes above 50 blocks.

CIFAR10 Convolutional Classification. We further experimented with a ResNet20 on CIFAR10. ResNet20 is a ResNetv2 with 3 stages of 3 residual blocks each, described in [11]. We compared a ResNet20 architecture with its GloNet version, obtained by removing the backbone and adding a GloNet layer before the classification head, as detailed in Section 4. We trained for 200 epochs. Learning curves are completely overlapping, with best test errors 91.12% and 91.08% for ResNet and GloNet respectively. This experiment shows that also

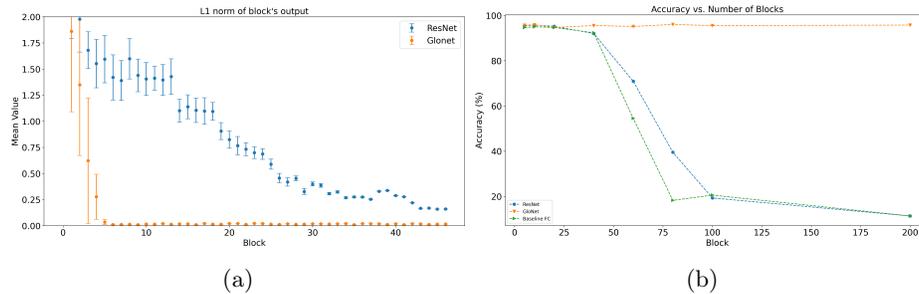


Fig. 5: (a) Block’s outputs L1 norm, GloNet (orange), and ResNetv2 (blue), 50 blocks, 200 epochs, MNIST. Mean and standard deviation from 5000 samples. GloNet uses only the first 4 blocks. In ResNetv2, all blocks give a significant contribution, as expected because of BN. (b) Accuracies of GloNet (orange), ResNetv2 (blue), and vanilla (green), 200 epochs, MNIST, against 6, 10, 24, 50, 60, 80, 100, and 200 blocks. GloNet is consistent across depths, while vanilla and ResNetv2 show a rapid decline with increased depth.

with convolutional architectures, GloNet performs on par with the traditional ResNet architecture, despite taking half the time for training.

GloNet for Vision Transformer Classification. A Visual Transformer (ViT) is a transformer applied to sequences of feature vectors extracted from image patches [6]. The n encoders outputs can be accumulated into a GloNet layer before going to the classification head, see Section 4 for details. In this experiment we compare ViT, with and without GloNet, on CIFAR-10, with 4, 5 and 6 encoders. Training plateaus at around 500 epochs, and final accuracies align with those from literature. With 4 and 6 encoders, accuracies overlap for ViT and GloNet-ViT. With 5 encoders, GloNet-ViT appears to improve over ViT, see Figure 6b. ViT best accuracies are 0.707, 0.709, and 0.725, and GloNet-ViT best accuracies are 0.709, 0.727, and 0.729, for 4, 5, and 6 encoders, respectively. This is a proof-of-concept experiment showcasing the robustness and versatility of GloNet for complex architectures like transformers.

Controllable Efficiency/Performance Trade-Off. In an experiment to demonstrate how GloNet can be used to choose an optimal efficiency/performance trade-off, we trained a 50-block GloNet fully connected architecture on MNIST for 200 epochs. After training, we progressively removed the last block, adjusted accordingly the GloNet layer to sum fewer blocks, and evaluated the shallower model without retraining. As shown in Figure 6a, removing up to 42 blocks did not significantly impact accuracy, illustrating GloNet’s ability to balance efficiency and performance.

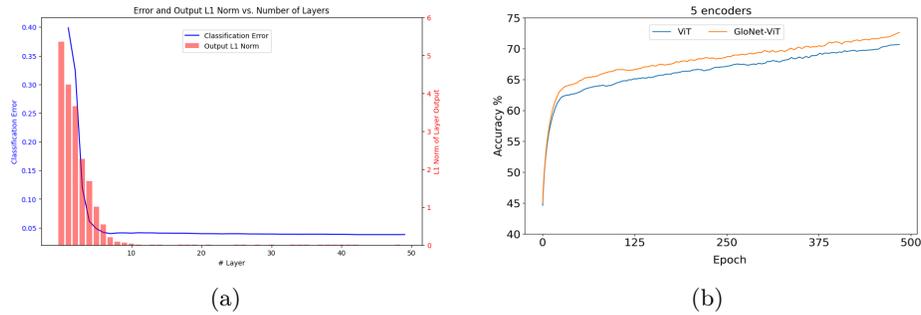


Fig. 6: (a) Block’s outputs L1 norm of GloNet (red), smaller networks performance (blue), 50 blocks, 200 epochs, MNIST. Performance plateaus after block eight. (b) Accuracy of ViT (blue), and GloNet-ViT (orange) with 5 encoders. Best accuracies 0.709, 0.727 for ViT, GloNet-ViT, respectively.

6 Conclusions

We introduce GloNet, a method designed to augment existing architectures without adding complexity or reducing performance. It effectively renders the architecture resilient to depth-related learning issues. As an alternative to ResNet, GloNet offers advantages, without any disadvantage: it achieves similar training outcomes in nearly half the time at depths where ResNet remains stable, and maintains consistent performance at greater depths where ResNet falters.

References

1. Ba, L.J., Kiros, J.R., Hinton, G.E.: Layer normalization. CoRR **abs/1607.06450** (2016), <http://arxiv.org/abs/1607.06450> 1
2. Ballester-Ripoll, R., Paredes, E.G., Pajarola, R.: Sobol tensor trains for global sensitivity analysis. Reliability Engineering & System Safety **183**, 311–322 (2019), <https://www.sciencedirect.com/science/article/pii/S0951832018303132> 6
3. Bengio, Y., Courville, A., Vincent, P.: Representation learning: A review and new perspectives. IEEE Trans. Pattern Anal. Mach. Intell. **35**(8), 1798–1828 (2013). <https://doi.org/10.1109/TPAMI.2013.50> 1
4. Bengio, Y., Simard, P., Frasconi, P.: Learning long-term dependencies with gradient descent is difficult. IEEE transactions on neural networks **5**(2), 157–166 (1994) 1
5. Di Cecco, A.: GloNet, <https://github.com/AntonioDiCecco/GloNet> 6
6. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021), <https://openreview.net/forum?id=YicbFdNTTy> 5, 10
7. Ebrahimi, M.S., Abadi, H.K.: Study of residual networks for image recognition. Lecture Notes in Networks and Systems, vol. 284, pp. 754–763. Springer (2021). https://doi.org/10.1007/978-3-030-80126-7_53 2
8. Gianfagna, L., Di Cecco, A.: Explainable AI with python. Springer (2021) 4

9. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: AISTATS. JMLR, vol. 9, pp. 249–256 (2010) [1](#), [2](#)
10. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: ICCV. pp. 1026–1034 (2015). <https://doi.org/10.1109/ICCV.2015.123> [1](#)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE CVPR. pp. 770–778 (2016) [1](#), [3](#), [9](#)
12. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: Computer Vision - ECCV 2016 - 14th European Conference Proceedings, Part IV. Lecture Notes in Computer Science, vol. 9908, pp. 630–645. Springer (2016). https://doi.org/10.1007/978-3-319-46493-0_38 [3](#), [4](#)
13. Huang, G., Liu, Z., Maaten, L.V.D., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR. pp. 2261–2269 (2017). <https://doi.org/10.1109/CVPR.2017.243> [3](#), [4](#)
14. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML. vol. 37, p. 448–456 (2015) [1](#), [5](#)
15. LeCun, Y., Bottou, L., Orr, G.B., Müller, K.: Efficient backprop. In: Neural Networks: Tricks of the Trade - Second Edition, Lecture Notes in Computer Science, vol. 7700, pp. 9–48 (2012). https://doi.org/10.1007/978-3-642-35289-8_3 [1](#)
16. Martin, C.H., Mahoney, M.W.: Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. J. Mach. Learn. Res. **22**, 165:1–165:73 (2021), <http://jmlr.org/papers/v22/20-410.html> [3](#)
17. Metta, C., Fantozzi, M., Papini, A., Amato, G., Bergamaschi, M., Galfrè, S.G., Marchetti, A., Vegliò, M., Parton, M., Morandin, F.: Increasing biases can be more efficient than increasing weights - extended version. In: IEEE/CVF, WACV (2024), <https://arxiv.org/abs/2301.00924> [2](#), [3](#)
18. Panda, P., Sengupta, A., Roy, K.: Conditional deep learning for energy-efficient and enhanced pattern recognition. In: Proceedings of DATE. p. 475–480 (2016) [4](#)
19. Rajesh, G.: A benchmark repository of Early Exit Neural Networks in .onnx format. https://github.com/gorakraj/earlyexit_onnx (2021) [4](#)
20. Santurkar, S., Tsipras, D., Ilyas, A., Madry, A.: How does batch normalization help optimization? In: NeurIPS 31. pp. 2488–2498 (2018), <https://proceedings.neurips.cc/paper/2018/hash/905056c1ac1dad141560467e0a99e1cf-Abstract.html> [5](#)
21. Saxe, A.M., McClelland, J.L., Ganguli, S.: Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In: ICLR (2014), <http://arxiv.org/abs/1312.6120> [1](#)
22. SGEMM: SGEMM GPU Kernel Performance. <https://archive.ics.uci.edu/ml/datasets/SGEMM+GPU+kernel+performance> (2018), accessed: 2023-05-03 [6](#)
23. Srivastava, R.K., Greff, K., Schmidhuber, J.: Training very deep networks. In: NeurIPS. pp. 2377–2385 (2015), <https://proceedings.neurips.cc/paper/2015/hash/215a71a12769b056c3c32e7299f1c5ed-Abstract.html> [2](#)
24. Teerapittayanon, S., McDanel, B., Kung, H.T.: Branchynet: Fast inference via early exiting from deep neural networks. In: ICPR. pp. 2464–2469 (2016). <https://doi.org/10.1109/ICPR.2016.7900006> [4](#)
25. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. In: ICLR (2017), <https://openreview.net/forum?id=Sy8gdB9xx> [3](#)
26. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning (still) requires rethinking generalization. Commun. ACM **64**(3), 107–115 (2021). <https://doi.org/10.1145/3446776> [3](#)