

# ResMLP: Feedforward networks for image classification with data-efficient training

Hugo Touvron<sup>1,2</sup>   Piotr Bojanowski<sup>1</sup>   Mathilde Caron<sup>1,3</sup>  
 Matthieu Cord<sup>1,2</sup>   Alaaeldin El-Nouby<sup>1,3</sup>   Edouard Grave<sup>1</sup>  
 Armand Joulin<sup>1</sup>   Gabriel Synnaeve<sup>1</sup>   Jakob Verbeek<sup>1</sup>   Hervé Jégou<sup>1</sup>

<sup>1</sup>Facebook AI   <sup>2</sup>Sorbonne University   <sup>3</sup>Inria

## Abstract

We present ResMLP, an architecture built entirely upon multi-layer perceptrons for image classification. It is a simple residual network that alternates (i) a linear layer in which image patches interact, independently and identically across channels, and (ii) a two-layer feed-forward network in which channels interact independently per patch. When trained with a modern training strategy using heavy data-augmentation and optionally distillation, it attains surprisingly good accuracy/complexity trade-offs on ImageNet. We will share our code based on the Timm library and pre-trained models.

## 1 Introduction

Recently, the transformer architecture [52], adapted from its original use in natural language processing with only minor changes, has achieved performance competitive with the state of the art on ImageNet-1k [43] when pre-trained with a sufficiently large amount of data [13]. Retrospectively, this achievement is yet another step towards less priors: convolutional neural networks had removed a lot of hand-made choices compared to hand-designed pre-CNN approaches, moving the paradigm of hard-wired features to hand-designed architectural choices. Vision transformers avoid making assumptions inherent to convolutional architectures and noticeably the translation invariance.

What these recent transformer-based works suggest is that longer training schedules, more parameters, more data [13] and/or more regularization [49], are sufficient to recover the important priors for tasks as complex as ImageNet classification. See also our discussion of related work in Section 4. This concurs with recent studies [2, 12] that better disentangle the benefits from the architectures from those of the training scheme.

In this paper, we push this trend further, and propose Residual Multi-Layer Perceptrons (ResMLP): a purely multi-layer perceptron (MLP) based architecture for image classification. We outline our architecture in Figure 1 and detail it further in Section 2. It is intended to be simple: it takes flattened patches as input, projects them with a linear layer, and sequentially updates them in turn with two residual operations: (i) a simple linear layer that provides interaction between the patches, which is applied to all channels independently; and (ii) an MLP with a single hidden layer, which is independently applied to all patches. At the end of the network, the patches are average pooled, and fed to a linear classifier.

This architecture is strongly inspired by the vision transformers (ViT) [13], yet it is much simpler in several ways: we do not use any form of attention, only linear layers along with the GELU non-linearity. Since our architecture is much more stable to train than transformers, we do not need batch-specific or cross-channel normalizations such as Batch-

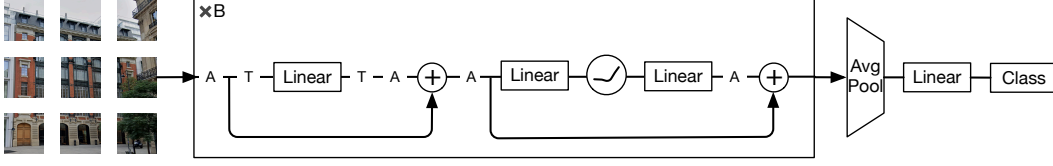


Figure 1: The ResMLP architecture: After flattening the patch into vectors, our network alternately processes them by (1) a communication layer between vectors implemented as a linear layer; (2) a two-layer residual perceptron. We denote by  $A$  the operator  $Aff$ , and by  $T$  the transposition.

Norm, GroupNorm or LayerNorm. Our training procedure mostly follows the one initially introduced for DeiT [49] and CaiT [50].

Due to its linear nature, the patch interactions in our model can be easily visualised and interpreted. While the interaction pattern learned in the first layer is very similar to a small convolutional filter, we observe more subtle interactions across patches in deeper layers. These includes some form of axial filters, and long-range interactions early in the network.

In summary, in this paper, we show that

- despite their simplicity, Residual Multi-Layer Perceptrons can reach surprisingly good accuracy/complexity trade-offs with ImageNet-1k training only<sup>1</sup>, without requiring normalization based on batch or channel statistics;
- these models benefit significantly from distillation methods [49];
- thank to its design where patch embeddings simply “communicate” through a linear layer, we can make observations on what kind of spatial interaction the network learns across layers.

## 2 Method

Our model, depicted in Figure 1, is inspired by the ViT model, of which it adopts the path flattening structure. We proceed to drastic simplifications. We refer the reader to Dosovitskiy *et al.* [13] for more details about the ViT architecture.

**The overall ResMLP architecture.** Our model, denoted by ResMLP, takes a grid of  $N \times N$  non-overlapping patches as input, where  $N$  is typically equal to 16. The patches are then independently passed through a linear layer to form a set of  $N^2$   $d$ -dimensional embeddings.

The resulting set of  $N^2$  embeddings are fed to a sequence of *Residual Multi-Layer Perceptron* layers to produce a set of  $N^2$   $d$ -dimension output embeddings. These output embeddings are then averaged as a  $d$ -dimension vector to represent the image, which is fed to a linear classifier to predict the label associated with the image. Training uses the cross-entropy loss.

**The Residual Multi-Perceptron Layer.** Our network is a sequence of layers that all have the same structure: a linear sublayer followed by a feedforward sublayer. Similar to the Transformer layer, each sublayer is paralleled with a skip-connection [19]. We do not apply Layer Normalization [1] because training is stable without it when using the following Affine transformation:

$$Aff_{\alpha, \beta}(\mathbf{x}) = \text{Diag}(\alpha)\mathbf{x} + \beta, \quad (1)$$

where  $\alpha$  and  $\beta$  are learnable vectors. This operation simply rescales and shifts the input component-wise. Moreover, it has no cost at inference time, as it can be fused in the adjacent

<sup>1</sup>Concurrent work by Tolstikhin *et al.* [48] brings complementary insights to ours: they achieve interesting performance with larger MLP models pre-trained on the larger public ImageNet-21k and even more data with the proprietary JFT-300M. In contrast, we focus on faster models trained on ImageNet-1k. Another concurrent work is the report by Melas-Kyriazi [32].

linear layer. Note, when writing  $\text{Aff}(\mathbf{X})$  the operation is applied independently to each column of  $\mathbf{X}$ . While similar to BatchNorm [24] and Layer Normalization [1], the  $\text{Aff}$  operator does not depend on any batch statistics. Therefore, it is closer to the recent LayerScale method [50], which improves the optimization of deep transformers when initializing  $\alpha$  to a small value. Note that LayerScale does not have a bias term.

We apply this transformation twice for each residual block. As as a pre-normalization  $\text{Aff}$  replaces the LayerNormalization, and avoids using channel-wise statistics. Here, we initialize  $\alpha = 1$ , and  $\beta = 0$ . As a post-processing of the residual block,  $\text{Aff}$  implements LayerScale and therefore we follow the same small value initialization for  $\alpha$  as in [50] for the post-normalization. Both transformations are integrated to the linear layers at inference.

Finally, we follow the same structure for the feedforward sublayer as in the Transformer, and we only replace the ReLU non-linearity by a GELU function [21].

Overall, our Multi-perceptron layer takes a set of  $N^2$   $d$ -dimensional input features stacked in a  $d \times N^2$  matrix  $\mathbf{X}$ , and outputs a set of  $N^2$   $d$ -dimension output features, stacked in a matrix  $\mathbf{Y}$  with the following set of transformations:

$$\mathbf{Z} = \mathbf{X} + \text{Aff} \left( (\mathbf{A} \text{Aff}(\mathbf{X})^\top)^\top \right), \quad (2)$$

$$\mathbf{Y} = \mathbf{Z} + \text{Aff}(\mathbf{C} \text{GELU}(\mathbf{B} \text{Aff}(\mathbf{Z}))), \quad (3)$$

where  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are the main learnable parameters of the layer. The dimensions of the parameter matrix  $\mathbf{A}$  are  $N^2 \times N^2$ , i.e, this sublayer mixes the information from all the locations, while the feedforward sublayer works per location. As a consequence, the intermediate activation matrix  $\mathbf{Z}$  has the same dimensions as the matrices  $\mathbf{X}$  and  $\mathbf{Y}$ . Finally, the parameter matrix  $\mathbf{B}$  and  $\mathbf{C}$  have the same dimensions as in a Transformer layer, that are  $4d \times d$  and  $d \times 4d$  respectively.

The main difference compared to a Transformer layer is that we replace the self-attention by the linear interaction defined in Eq. (2). While self-attention computes a convex combination of other features with coefficients that are data dependent, the linear interaction layer in Eq. (2) computes a general linear combination using learned coefficients that are not data dependent. As compared to a convolutional layers which have local support and share weights across space, our linear patch interaction layer offers a global support and does not share weights, moreover it is applied independently across channels.

**Relationship to the Vision Transformer.** Our model can be regarded as a drastic simplification of the ViT model by Dosovitskiy *et al.* [13]. We depart from this model as follows:

- We do not include any self-attention block. Instead we have a linear patch interaction layer without non-linearity.
- We do not have the extra “class” token that is typically used in these models to aggregate information via attention. Instead, we simply use average pooling. We do, however, also consider a specific aggregation layer as a variant, which we describe in the next paragraph.
- Similarly, we do not include any form of positional embedding: it is not required as the linear communication module between patches implicitly takes into account the patch position.
- Instead of pre-LayerNormalization, we use a simple learnable affine transform, thus avoiding any form of batch and channel-wise statistics.

**Class-MLP: MLP with class embedding.** As an alternative to average pooling, we also experimented with an adaptation of the class-attention introduced in CaiT [50]. It consists of two layers that have the same structure as the transformer, but in which only the class token is updated based on the frozen patch embeddings. We translate this method to our network, by replacing the attention-based interaction between the class and patch embeddings by simple linear layers. This increases the performance, at the expense of adding some parameters and computational cost. We refer to this pooling variant as “class-MLP”.

### 3 Experiments

In this section, we present experimental results for our ResMLP architecture for image classification. We also study the impact of the different components in the ResMLP architecture in a series of ablations.

#### 3.1 Experimental setting

**Datasets.** We train our models on the ImageNet-1k dataset [43], that contains 1.2M images evenly spread over 1,000 object categories. In the absence of an available test set for this benchmark, we follow the standard practice in the community by reporting performance on the validation set. This is not ideal since the validation set was originally designed to select hyper-parameters. Comparing methods on this set may not be conclusive enough because an improvement in performance may not be caused by better modeling, but by a better selection of hyper-parameters. To mitigate this risk, we report additional results on two alternative versions of ImageNet that have been built to have distinct validation and test sets, namely the ImageNet-real [3] and ImageNet-v2 [42] datasets. Our hyper-parameters are mostly adopted from Touvron et al. [49, 50].

**Training paradigms.** We consider two training paradigms in our experiments:

- *Supervised learning:* We train ResMLP from labeled images with a softmax classifier and cross-entropy loss. This paradigm is the main focus of our work.
- *Knowledge distillation:* We employ the knowledge distillation procedure proposed by Touvron et al. [49] to guide the training of ResMLP with a convnet. Note that in this paradigm, the ResMLP architecture does not require access to the labels during training, only to an existing pre-trained model.

**Hyper-parameter setting.** In the case of supervised learning, we train our network with Lamb optimizer [55] with a learning rate of  $5 \times 10^{-3}$  and weight decay 0.2. We initialize the LayerScale parameters as a function of the depth by following off-the-shelf those proposed by Touvron et al. [50] for CaiT. The rest of the hyper-parameters follow the default setting used in DeiT [49]. For the knowledge distillation paradigm, we use the same RegNety-16GF [41] as in DeiT with the same training schedule.

#### 3.2 Main Results

In this section, we compare our architecture with standard neural networks of comparable size and throughput on ImageNet.

**Comparison with Transformers and convnets in a supervised setting.** In Table 1, we compare ResMLP with different convolutional and Transformer architectures. For completeness, we also report the best-published numbers obtained with a model trained on ImageNet alone. As expected, in terms of the trade-off between accuracy, FLOPs, and throughput, ResMLP is not as good as convolutional networks or Transformers. However, their accuracy is very encouraging. Indeed, we compare them with architectures that have benefited from years of research and careful optimization towards these trade-offs. Overall, our results suggest that the structural constraints imposed by the layer design do not have a drastic influence on performance, especially when training models with enough data and modern advances in training and regularization.

**Improving model convergence with knowledge distillation.** We also study our model when training following the knowledge distillation paradigm from Touvron et al. [49]. In their work, the authors show the impact of training a ViT model by distilling it from an EfficientNet. In this experiment, we explore if ResMLP also benefits from this procedure and summarize our results in Table 2. We observe that similar to DeiT models, ResMLP

	Arch.	#params ( $\times 10^6$ )	throughput (im/s)	FLOPS ( $\times 10^9$ )	Peak Mem (MB)	Top-1 Acc.
<i>State of the art</i>	CaiT-M48 $\uparrow$ 448T [50]	356	5.4	329.6	5477.8	86.5
	NfNet-F6 SAM [5]	438	16.0	377.3	5519.3	86.5
<i>Convolutional networks</i>	EfficientNet-B3 [46]	12	661.8	1.8	1174.0	81.1
	EfficientNet-B4 [46]	19	349.4	4.2	1898.9	82.6
	EfficientNet-B5 [46]	30	169.1	9.9	2734.9	83.3
	RegNetY-4GF [40]	21	861.0	4.0	568.4	80.0
	RegNetY-8GF [40]	39	534.4	8.0	841.6	81.7
	RegNetY-16GF [40]	84	334.7	16.0	1329.6	82.9
<i>Transformer networks</i>	DeiT-S [49]	22	940.4	4.6	217.2	79.8
	DeiT-B [49]	86	292.3	17.5	573.7	81.8
	CaiT-XS24 [50]	27	447.6	5.4	245.5	81.8
<i>Feedforward networks</i>	ResMLP-12	15	1415.1	3.0	179.5	76.6
	ResMLP-24	30	715.4	6.0	235.3	79.4
	ResMLP-36	45	478.7	8.9	291.3	79.7

Table 1: **Comparison between architectures on ImageNet classification.** We compare different architectures for images based on convolutional networks, Transformers and feedforward networks with comparable FLOPS and number of parameters. We report Top-1 accuracy on the validation set of ImageNet with different measure of complexity: throughput, FLOPS, number of parameters and peak memory usage. All the models use  $224 \times 224$  images as input. All the Transformers and feedforward networks uses  $14 \times 14$  patches of size  $16 \times 16$ . Feedforward networks use a working dimension of  $d=384$ . The throughput is measured on a single V100 GPU with batch size fixed to 32. For reference, we include the state of the art with ImageNet training only (86.5% Top-1 with very large models).

Arch.	Teacher	w/o Dist.	w/ Dist.
DeiT-S	RegNety-16GF	79.9	81.2
ViT-B	RegNety-16GF	81.8	83.4
ResMLP-12	RegNety-16GF	76.6	77.8

Table 2: **Impact of distillation on ResMLP.** We report Top-1 accuracy on the validation set of ImageNet when training models with and without knowledge distillation using a pre-trained model. Accuracies for Transformer-based models are taken from Touvron *et al.* [49].

greatly benefits from distilling from a convnet. This result concurs with the observations made by d’Ascoli *et al.* [11], who used convnets to initialize feedforward networks. Even though our setting differs from theirs in scale, the problem of overfitting for feedforward networks is still very much present on ImageNet. The additional regularization obtained from the distillation is a possible explanation for this improvement.

### 3.3 Transfer learning

We evaluate the quality of features obtained from a ResMLP architecture when transferring them to other domains. The goal is to assess if the features generated from a feedforward network are more prone to overfitting on the training data distribution.

We adopt the typical setting where we pre-train a model on ImageNet-1k and fine-tune it on the training set associated with a specific domain. We report the performance with different architectures on different image benchmarks in Table 3, namely CIFAR-10 and CIFAR-100 [28], Flowers-1022 [35], Stanford Cars [27] and iNaturalist [22]. We refer the reader to the corresponding references for a more detailed description of the different datasets. We observe that the performance of our ResMLP are competitive with the existing architectures, showing that pretraining feedforward models with enough data and regularization via data augmentation greatly reduces their tendency to overfit on the original distribution. Interestingly, this regularization also prevents them from overfitting on the training set of smaller dataset during the fine-tuning stage.

Architecture	FLOPs	res.	CIFAR <sub>10</sub>	CIFAR <sub>100</sub>	Flowers102	Cars	iNat <sub>18</sub>	iNat <sub>19</sub>
EfficientNet-B7	37.0B	600	98.9	91.7	98.8	94.7	–	–
ViT-B/16	55.5B	384	98.1	87.1	89.5	–	–	–
ViT-L/16	190.7B	384	97.9	86.4	89.7	–	–	–
DeiT-B/16	17.5B	224	99.1	90.8	98.4	92.1	73.2	77.7
ResNet50	4.1B	224	–	–	96.2	90.0	68.4	73.7
Grafit/ResNet50	4.1B	224	–	–	97.6	92.7	68.5	74.6
ResMLP-12	3.0B	224	98.1	87.0	97.4	84.6	–	–
ResMLP-24	6.0B	224	98.7	89.5	97.9	89.5	64.3	72.5

Table 3: **Evaluation on transfer learning.** We compare models trained on ImageNet on for transfer to datasets covering different domains. The ResMLP architecture takes  $224 \times 224$  images during training and transfer, while the ViTs and EfficientNet-B7 work with higher resolutions, see “res.” column.

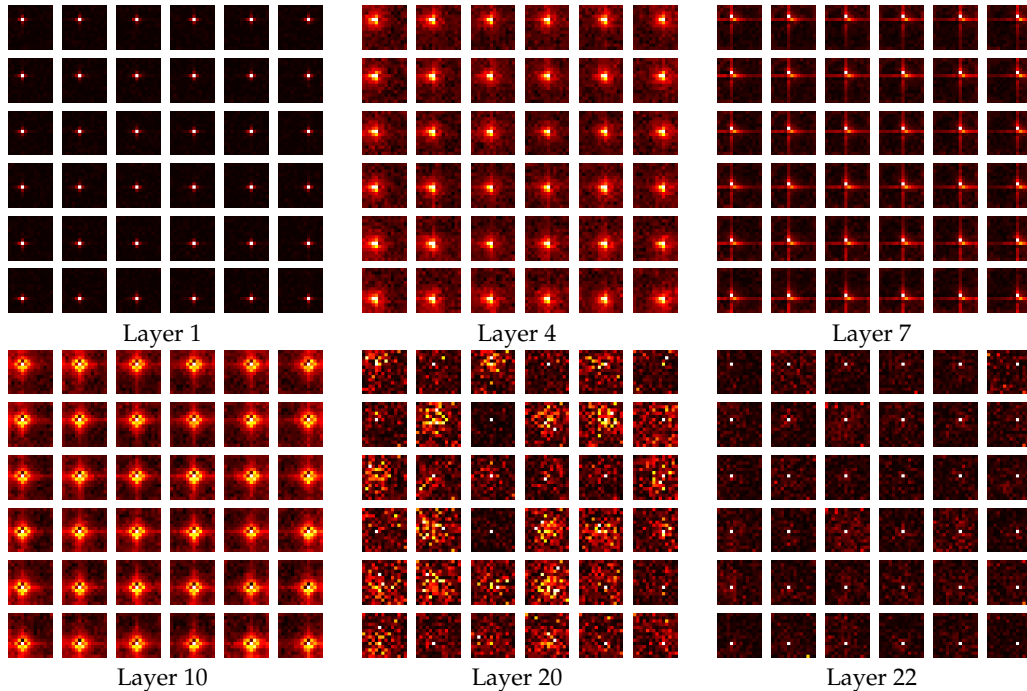


Figure 2: **Visualisation of the linear layers in ResMLP-24.** For each layer we visualise the rows of the matrix  $\mathbf{A}$  as a set of  $14 \times 14$  pixel images, for sake of space we only show the rows corresponding to the  $6 \times 6$  central patches. We observe patterns in the linear layers that share similarities with convolutions. In appendix B we provide comparable visualizations for all layers of a ResMLP-12 model.

### 3.4 Vizualization and analysis

**Visualisation.** Because they are linear, our patch interaction layers from Eq. (2) are easily interpretable. In Fig 2 we visualise the rows of the interaction matrices  $\mathbf{A}$  as  $N \times N$  images, for our ResMLP-24 model. The early layers show convolution-like patterns: the learned weights resemble shifted versions of each other and have local support. Interestingly, in many layers, the support also extends along both axes, most prominently seen in layer seven. The last seven layers of the network are different: they consist of a spike for the patch itself and a diffuse response across other patches with larger or smaller magnitude; see layers 20 and 22.

**Measuring sparsity of the weights.** The visualizations described above suggest that the linear communication layers are sparse. We analyze this quantitatively in more detail in Fig. 3. We measure the sparsity of the matrix  $\mathbf{A}$ , and compare it to the sparsity of  $\mathbf{B}$  and  $\mathbf{C}$



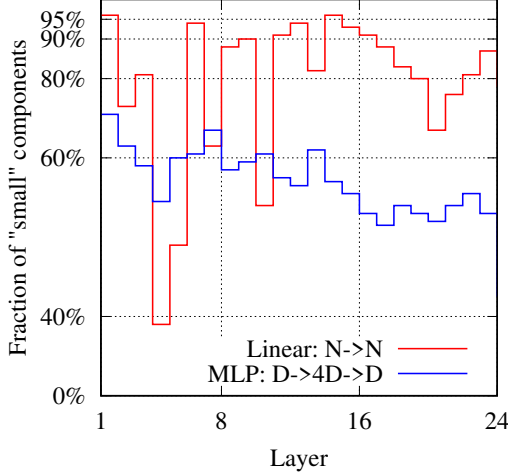


Figure 3: How sparse are the linear layers? For each layer (linear and MLP), we show the rate of components whose values is lower than 1% of the maximum value. Observe how most linear interaction layers are significantly sparser than the matrices involved in the patch update MLP.

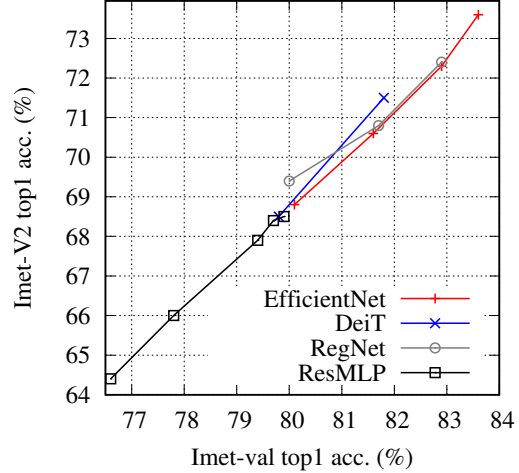


Figure 4: Accuracy on ImageNet-V2 vs measured on ImageNet-val for different models. The relative offsets between curves reflect to which extent models are overfitted to ImageNet-val w.r.t. hyper-parameter selection. Our ResMLPs are comparable to other networks in that respect.

from the per-patch MLP. Since there are no exact zeros, we measure the rate of components whose absolute value is lower than 5% of the maximum value. Note, discarding the small values is analogous to the case where we normalize the matrix by its maximum and use a finite-precision representation of weights. For instance, with a 4-bits representation of weight, one would typically round to zero all weights whose absolute value is below 6.25% of the maximum value.

The measurements in Fig. 3 show that all three matrices are sparse, with the layers implementing the patch communication being significantly more so. This suggests that they may be compatible with parameter pruning, or better, with modern quantization techniques that induce sparsity at training time, such as Quant-Noise [17] and DiffQ [16]. The sparsity structure, in particular in earlier layers (see Fig. 2), hints that we could implement the patches mixing linear layer with a convolution. This line of research on network compression is beyond the scope of our paper, yet we believe it worth investigating in the future.

**Control of overfitting.** Since MLPs are subject to overfitting, we show in Fig. 4 a control experiment to probe for problems with generalization. We explicitly analyze the differential of performance between the ImageNet-val and the distinct ImageNet-V2 test set. The degree of overfitting of our MLP-based model is overall comparable to that of other transformer-based architectures or convnets.

### 3.5 Ablation studies

Table 4 reports the ablation study of our base network and a summary of our preliminary exploratory studies. We discuss the ablation below and give more detail about the early experiments in Appendix A.

**Communication with low-resolution convolutions.** As discussed when presenting the visualization, the linear layers (that implicitly exploit the patch position) look like convolutions for most of the layers. In this experiment, we have replaced the linear layer with a  $3 \times 3$  convolution operating on  $14 \times 14$  patches of dimension  $D = 384$ . Our ablation shows that this choice improves the performance, showing that low-resolution convolutions at all layers is an interesting alternative to the most common design of convnets, where early layers operate at high resolution and small feature dimension.

#layers	Supervision	Norm.	Pooling layer	Patch commun.	top-1 acc. on ImageNet		
					Inet-val	Inet-real	Inet-V2
12	regular	$\times$	average	linear	76.6	83.3	64.4
24	regular	$\times$	average	linear	79.4	85.3	67.9
36	regular	$\times$	average	linear	79.7	85.6	68.4
12	distillation	$\times$	average	linear	77.8	84.6	66.0
12	regular	$\times$	class-MLP	linear	77.5	84.3	65.9
24	regular	$\times$	class-MLP	linear	79.9	85.8	68.5
12	regular	LayerNorm	average	linear	77.7	84.1	65.7
12	regular	$\times$	average	conv 3x3	78.6	85.3	67.1

Table 4: Main experiments in classification and ablation.

**Normalization.** Our primary network configuration does not contain any batch normalizations. Instead, we use affine per-channel transforms such as Layer Normalization [1], typically used in transformers. In preliminary experiments with pre-norm and post-norm [20], we observed that both choices lead to convergence. Pre-normalization in conjunction with Batch Normalization could provide an accuracy gain in some cases (see Appendix A). However, for the sake of simplicity, we preferred not to introduce any dependency on batch statistics, which is why we resorted on the `Aff` operator only.

**Preliminary experiments.** In our early exploration, we evaluated several alternative design choices. We summarize our main findings below:

- *Block design.* We have tried several variants for the patch interaction layer. Amongst them, using the same MLP structure as for patch processing. In our experiments, the simpler choice of a single linear layer led to better performance while being more efficient. Moreover, it requires fewer parameters than a residual MLP block.
- *Activation.* We choose to use a GELU [21] function. ReLU [18] also gives a good performance, but we observed that it was a bit more unstable in some settings. We did not manage to get good results with SiLU [21] and HardSwish [23].
- *Positional encoding and class token.* As in transformers, we could use positional embeddings mixed with the input patches. In our experiments, we did not see any benefit from using these features. This observation suggests that our linear patch interaction layer provides sufficient spatial communication. Referencing absolute positions obviates the need for any form of positional encoding.
- *Class-MLP.* In contrast, specialized layers extracting information from image patches with a class embedding increase performance by +0.5% top-1 acc. This improvement is comparable to adding the same number of layers in the main network. However, using these specialized layers is more efficient.

## 4 Related work

We review the research on applying Fully Connected Network (FCN) for computer vision problems as well as other architectures that shares common modules with our model.

**Fully-connected network for images.** Many studies have shown that FCNs are competitive with convnets for the tasks of digit recognition [44, 10], keyword spotting [6] and handwriting recognition [4]. Several works [51, 33, 31] have questioned if FCNs are also competitive on natural image datasets, such as CIFAR-10 [28]. More recently, d’Ascoli *et al.* [11] have shown that a FCN initialized with the weights of a pretrained convnet achieves performance that are superior than the original convnet. Neyshabur [34] further extend this line of work by achieving competitive performance by training an FCN from scratch but with a regularizer that constraint the models to be close to a convnet. These studies have



been conducted on small scale datasets with the purpose of studying the impact of architectures on generalization in terms of sample complexity [15] and energy landscape [25]. In our work, we show that, in the larger scale setting of ImageNet, FCNs can attain surprising accuracy without any constraint or initialization inspired by convnets.

Finally, the application of FCN networks in computer vision have also emerged in the study of the properties of networks with infinite width [36], or for inverse scattering problems [26]. More interestingly, the Tensorizing Network [37] is an approximation of very large FCN that shares similarity with our model, in that they intend to remove prior by approximating even more general tensor operations, *i.e.*, not arbitrarily marginalized along some pre-defined sharing dimensions. However, their method is designed to compress the MLP layers of a standard convnets.

**Other architectures with similar components.** Our FCN architecture shares several components with other architectures, such as convnets [30, 29] or Transformers [52]. A fully connected layer is equivalent to a convolution layer with a  $1 \times 1$  receptive field, and several work have explored convnet architectures with small receptive fields. For instance, the VGG model [45] uses  $3 \times 3$  convolutions, and later, other architectures such as the ResNext [54] or the Xception [9] mix  $1 \times 1$  and  $3 \times 3$  convolutions. In contrast to convnets, in our model interaction between patches is obtained via a linear layer that is shared across channels, and that relies on absolute rather than relative positions.

More recently, Transformers have emerged as a promising architecture for computer vision [8, 14, 39, 49, 56]. In particular, our architecture takes inspiration from the structure used in the Vision Transformer (ViT) [14], and as consequence, shares many components. Our model takes a set of non-overlapping patches as input and passes them through a series of MLP layers that share the same structure as the Transformer, replacing the self-attention layer with a linear patch interaction layer. Both layers have a global field-of-view, unlike convolutional layers. Whereas in self-attention the weights to aggregate information from other patches are data dependent through queries and keys, in ResMLP the weights are not data dependent and only based on absolute positions of patches. In our implementation we follow the improvements of DeiT [49] to train vision transformers, use the skip-connections from the ResNet [19] with pre-normalization of the layers [7, 20].

Finally, our work questions the importance of self-attention in the performance of vision transformers, or at least whether the performance increase they provide justify the training challenges that they raise. Similar observations have been made in natural language processing. Notably, the Synthesizer [47] shows that dot-product self-attention can be replaced by a feedforward network, with competitive performance on sentence representation benchmarks. As opposed to our work, the Synthesizer does use data dependent weights, but in contrast to transformers, the weights are determined from the query point only.

## 5 Conclusion

In this paper we have shown that a simple residual architecture, whose residual blocks consist of a one-hidden layer feed-forward network and a linear patch interaction layer, achieves an unexpectedly high performance on ImageNet classification benchmarks, provided that we adopt a modern training strategy such as those recently introduced for transformer-based architecture. Thanks to their simple structure, with linear layers as the main mean of communication between patches, we can visualize the filters inherently learned by this simple MLP. While some of the layers are similar to convolutional filters, we also observe sparse long-range interactions as early as the second layer of the network. We hope that our spatial prior-free model will contribute to further understanding of what networks with less priors learn, and potentially guide the design choices of future networks without the pyramidal design prior adopted by most convolutional neural networks.

## 6 Acknowledgments

We would like to thank Mark Tygert for relevant references. This work builds upon the Timm library [53] by Ross Wightman.

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 2, 3, 8
- [2] Irwan Bello, William Fedus, Xianzhi Du, Ekin D Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, and Barret Zoph. Revisiting resnets: Improved training and scaling strategies. *arXiv preprint arXiv:2103.07579*, 2021. 1
- [3] Lucas Beyer, Olivier J. Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aaron van den Oord. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020. 4
- [4] Théodore Bluche. *Deep neural networks for large vocabulary handwritten text recognition*. PhD thesis, Université Paris-Sud, 2015. 8
- [5] A. Brock, Soham De, S. L. Smith, and K. Simonyan. High-performance large-scale image recognition without normalization. *arXiv preprint arXiv:2102.06171*, 2021. 5
- [6] Clément Chatelain. *Extraction de séquences numériques dans des documents manuscrits quelconques*. PhD thesis, Université de Rouen, 2006. 8
- [7] Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Niki Parmar, Mike Schuster, Zhifeng Chen, et al. The best of both worlds: Combining recent advances in neural machine translation. In *Annual Meeting of the Association for Computational Linguistics*, 2018. 9
- [8] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019. 9
- [9] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 9
- [10] Dan Claudiu Cireşan, Ueli Meier, Luca Maria Gambardella, and Jürgen Schmidhuber. Deep big multilayer perceptrons for digit recognition. In *Neural networks: tricks of the trade*, pages 581–598. Springer, 2012. 8
- [11] Stéphane d’Ascoli, Levent Sagun, Joan Bruna, and Giulio Biroli. Finding the needle in the haystack with convolutions: on the benefits of architectural bias. In *NeurIPS*, 2019. 5, 8
- [12] Piotr Dollár, Mannat Singh, and Ross Girshick. Fast and accurate model scaling. *arXiv preprint arXiv:2103.06877*, 2021. 1
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*, 2021. 1, 2, 3
- [14] Alexey Dosovitskiy, Jost Tobias Springenberg, Martin A. Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. In *NeurIPS*, 2014. 9
- [15] Simon S Du, Yining Wang, Xiyu Zhai, Sivaraman Balakrishnan, Ruslan Salakhutdinov, and Aarti Singh. How many samples are needed to estimate a convolutional neural network? In *NeurIPS*, pages 371–381, 2018. 9
- [16] Alexandre Défossez, Yossi Adi, and Gabriel Synnaeve. Differentiable model compression via pseudo quantization noise, 2021. 7
- [17] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Rémi Gribonval, Hervé Jégou, and Armand Joulin. Training with quantization noise for extreme model compression. In *International Conference on Learning Representations*, 2021. 7
- [18] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *International Conference on Machine Learning*, 2011. 8
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 9
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027*, 2016. 8, 9
- [21] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 3, 8
- [22] Grant Van Horn, Oisín Mac Aodha, Yang Song, Alexander Shepard, Hartwig Adam, Pietro Perona, and Serge J. Belongie. The inaturalist species classification and detection dataset. *arXiv preprint arXiv:1707.06642*, 2017. 5

- [23] A. Howard, Mark Sandler, G. Chu, Liang-Chieh Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan, Quoc V. Le, and H. Adam. Searching for mobilenetv3. *International Conference on Computer Vision*, 2019. 8
- [24] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 2015. 3
- [25] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017. 9
- [26] Yuehaw Khoo and Lexing Ying. Switchnet: a neural network model for forward and inverse scattering problems. *SIAM Journal on Scientific Computing*, 41(5):A3182–A3201, 2019. 9
- [27] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, 2013. 5
- [28] A. Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, University of Toronto, 2009. 5, 8
- [29] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 9
- [30] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 9
- [31] Zhouhan Lin, Roland Memisevic, and Kishore Konda. How far can we go without convolution: Improving fully-connected networks. *arXiv preprint arXiv:1511.02580*, 2015. 8
- [32] Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly well on imagenet. *arXiv preprint arXiv:2105.02723*, 2021. 2
- [33] Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):1–12, 2018. 8
- [34] Behnam Neyshabur. Towards learning convolutions from scratch. In *NeurIPS*, 2020. 8
- [35] M-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing*, 2008. 5
- [36] Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In *International Conference on Learning Representations*, 2019. 9
- [37] Alexander Novikov, Dmitry Podoprikin, Anton Osokin, and Dmitry Vetrov. Tensorizing neural networks. *arXiv preprint arXiv:1509.06569*, 2015. 9
- [38] Alexander Novikov, Dmitrii Podoprikin, Anton Osokin, and Dmitry P Vetrov. Tensorizing neural networks. In *Neurips*, 2015.
- [39] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, 2018. 9
- [40] Filip Radenović, Giorgos Tolias, and Ondrej Chum. Fine-tuning CNN image retrieval with no human annotation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. 5
- [41] Ilija Radosavovic, Raj Prateek Kosaraju, Ross B. Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. *Conference on Computer Vision and Pattern Recognition*, 2020. 4
- [42] B. Recht, Rebecca Roelofs, L. Schmidt, and V. Shankar. Do ImageNet classifiers generalize to ImageNet? In *International Conference on Machine Learning*, 2019. 4
- [43] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International journal of Computer Vision*, 2015. 1, 4
- [44] Patrice Y Simard, David Steinkraus, John C Platt, et al. Best practices for convolutional neural networks applied to visual document analysis. In *ICDAR*, 2003. 8
- [45] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 9
- [46] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019. 5
- [47] Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer: Rethinking self-attention in transformer models. *arXiv preprint arXiv:2005.00743*, 2020. 9
- [48] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. MLP-Mixer: An all-MLP architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021. 2
- [49] Hugo Touvron, M. Cord, M. Douze, F. Massa, Alexandre Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. 1, 2, 4, 5, 9, I

- [50] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. *arXiv preprint arXiv:2103.17239*, 2021. 2, 3, 4, 5, 1
- [51] Gregor Urban, Krzysztof J Geras, Samira Ebrahimi Kahou, Ozlem Aslan, Shengjie Wang, Rich Caruana, Abdelrahman Mohamed, Matthai Philipose, and Matt Richardson. Do deep convolutional nets really need to be deep and convolutional? In *International Conference on Learning Representations*, 2017. 8
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 9
- [53] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. 10
- [54] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 9
- [55] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training BERT in 76 minutes. In *International Conference on Learning Representations*, 2020. 4
- [56] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2020. 9

# ResMLP: Feedforward networks for image classification with data-efficient training

## Appendix

### A Report on our exploration phase

As discussed in the main paper, our work on designing a residual multi-layer perceptron was inspired by the Vision Transformer. For our exploration, we have adopted the recent CaiT variant [50] as a starting point. This transformer-based architecture achieves state-of performance with Imagenet-training only (achieving 86.5% top-1 accuracy on Imagenet-val for the best model). Most importantly, the training is relatively stable with increasing depth.

In our exploration phase, our objective was to radically simplify this model. For this purpose, we have considered the CaiT-S24 model for faster iterations. This network consists of 24-layer with a working dimension of 384. All our experiments below were carried out with images in resolution  $224 \times 224$  and  $N = 16 \times 16$  patches. Trained with regular supervision, CaiT-S24 attains 82.7% top-1 acc. on Imagenet.

**SA  $\rightarrow$  MLP.** The self-attention can be seen a weight generator for a linear transformation on the values. Therefore, our first design modification was to get rid of the self-attention by replacing it by a residual feed-forward network, which takes as input the *transposed* set of patches instead of the patches. In other terms, in this case we alternate residual blocks operating along the channel dimension with some operating along the patch dimension. In that case, the MLP replacing the self-attention consists of the sequence of operations

$$(\cdot)^T \text{ — linear } N \times 4N \text{ — GELU — linear } 4N \times N \text{ — } (\cdot)^T$$

Hence this network is symmetrical in  $N$  and  $d$ . By keeping the other elements identical to CaiT, the accuracy drops to 80.2% (-2.5%) when replacing self-attention layers.

**Class-attention  $\rightarrow$  class-MLP.** If we further replace the class-attention layer of CaiT by a MLP as described in our paper, then we obtain an attention-free network whose top-1 accuracy on Imagenet-val is 79.2%, which is comparable to a ResNet-50 trained with a modern training strategy. This network has served as our baseline for subsequent ablations. Note that, at this stage, we still include LayerScale, a class embedding (in the class-MLP stage) and positional encodings.

**Distillation.** The same model trained with distillation inspired by Touvron et al. [49] achieves 81.5%. The distillation variant we choose corresponds to the “hard-distillation”, whose main advantage is that it does not require any parameter-tuning compared to vanilla cross-entropy. Note that, in all our experiments, this distillation method seems to bring a gain that is complementary and seemingly almost orthogonal to other modifications.

**Activation: LayerNorm  $\rightarrow$  X.** We have tried different activations on top of the aforementioned MLP-based baseline, and kept GeLU for its accuracy and to be consistent with the transformer choice.

Activation	top-1 acc.
GeLU (baseline)	79.2%
SILU	78.7%
Hard Swish	78.8%
ReLU	79.1%

**Ablation on the size of the communication MLP.** For the MLP that replaced the class-attention, we have explored different sizes of the latent layer, by adjusting the expansion factor  $e$  in the sequence: linear  $N \times e \times N$  — GELU — linear  $e \times N \times N$ . For this experiment we used average pooling to aggregating the patches before the classification layer.

expansion factor $\times e$	$\times 0.25$	$\times 0.5$	$\times 1$	$\times 2$	$\times 3$	$\times 4$
Imnet-val top-1 acc.	78.6	79.2	79.2	79.3	78.8	78.8

We observe that a large expansion factor is detrimental in the patch communication, possibly because we should not introduce too much capacity in this residual block. This has motivated the choice of adopting a simple linear layer of size  $N \times N$ : This subsequently improved performance to 79.5% in a setting comparable to the table above. Additionally, as shown earlier this choice allows visualizations of the interaction between patches.

**Normalization.** On top of our MLP baseline, we have tested different variations for normalization layers. We report the variation in performance below.

Pre-normalization	top-1 acc.
Layernorm (baseline)	79.2%
Batch-Norm	+0.8%
$\ell_2$ -norm	+0.4%
no norm (Aff)	+0.4%

For the sake of simplicity, we therefore adopted only the Aff transformation so as to not depend on any batch or channel statistics.

**Position encoding.** In our experiments, removing the position encoding does not change the results when using a MLP or a simple linear layer as a communication mean across patch embeddings. This is not surprising considering that the linear layer implicitly encodes each patch identity as one of the dimension, and that additionally the linear includes a bias that makes it possible to differentiate the patch positions before the shared linear layer.

## B Analysis of interaction layers in 12-layer networks

In this section we further analyze the linear interaction layers in 12-layer models.

In Figure B.1 we consider a ResMLP-12 model trained on the ImageNet-1k dataset, as explained in Section 3.1, and show all the 12 linear patch interaction layers. The linear interaction layers in the supervised 12-layer model are similar to those observed in the 24-layer model in Figure 2.

We also provide the corresponding sparsity measurements for this model in Figure B.2, analogous to the measurements in Figure 3 for the supervised 24-layer model. The sparsity levels in the supervised 12-layer model (left panel) are similar to those observed in the supervised 24-layer model, cf. Figure 3. In the right panel of Figure B.2 we consider the sparsity levels of the Distilled 12-layer model, which are overall similar to those observed for supervised the 12-layer and 24-layer models.



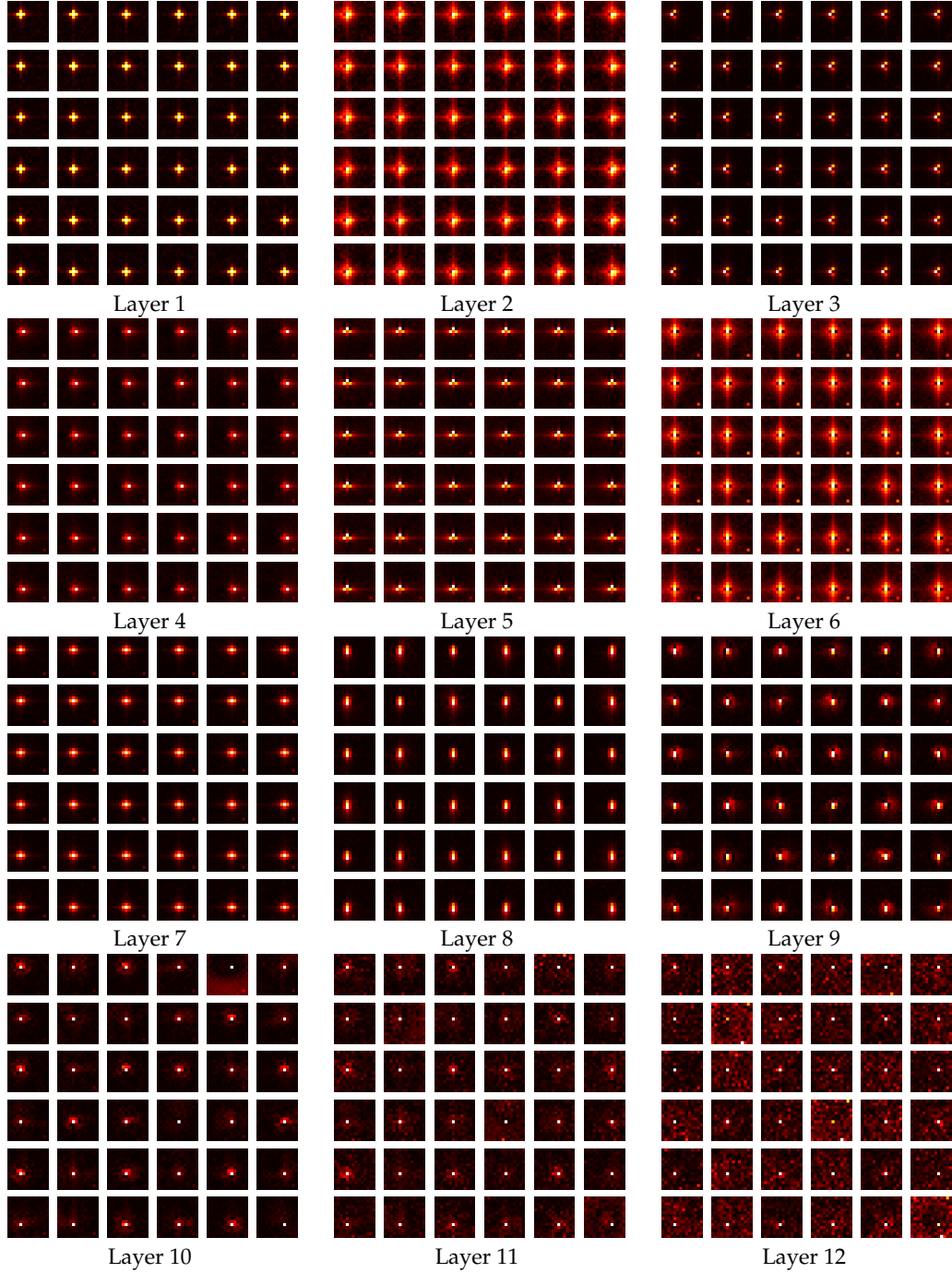


Figure B.1: **Visualisation of the linear interaction layers in the supervised ResMLP-12 model.** For each layer we visualise the rows of the matrix  $\mathbf{A}$  as a set of  $14 \times 14$  pixel images, for sake of space we only show the rows corresponding to the  $6 \times 6$  central patches.

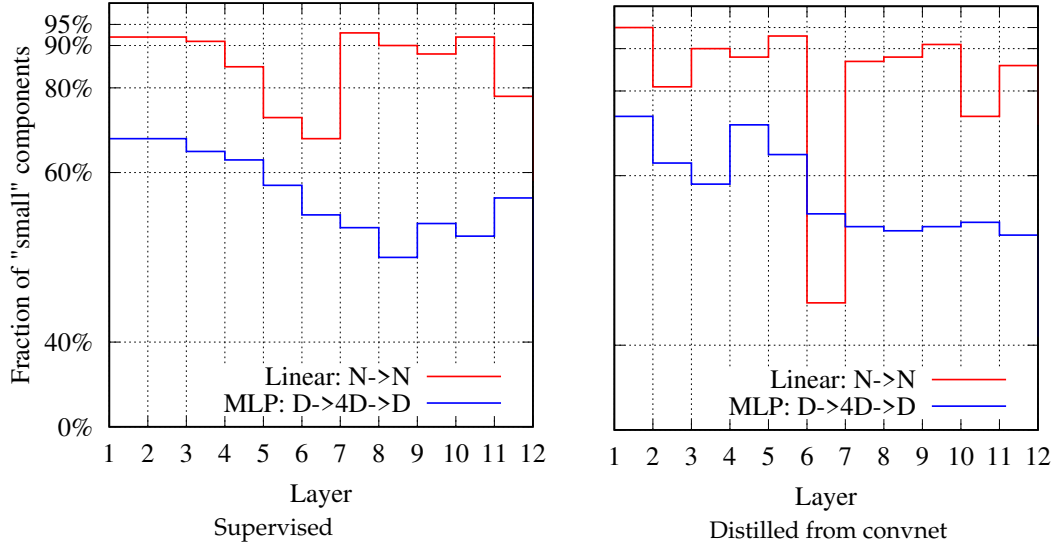


Figure B.2: Degree of sparsity (fraction of values small than 5% of the maximum) for **Linear** and **MLP** layers, for ResMLP-12 networks. The network trained in supervised mode and the one learned with distillation overall have a comparable degree of sparsity. The self-supervised model, trained during 300 epochs *vs.* 400 for the other ones, is less sparse on the patch communication linear layer.

## C Model definition in Pytorch

In Algorithm 1 we provide the pseudo-pytorch-code associated with our model.

---

**Algorithm 1** Pseudocode of ResMLP in PyTorch-like style

---

```
# No norm layer
class Affine(nn.Module):
    def __init__(self, dim):
        super().__init__()
        self.alpha = nn.Parameter(torch.ones(dim))
        self.beta = nn.Parameter(torch.zeros(dim))
    def forward(self, x):
        return self.alpha * x + self.beta

# MLP on channels
class Mlp(nn.Module):
    def __init__(self, dim):
        super().__init__()
        self.fc1 = nn.Linear(dim, 4 * dim)
        self.act = nn.GELU()
        self.fc2 = nn.Linear(4 * dim, dim)
    def forward(self, x):
        x = self.fc1(x)
        x = self.act(x)
        x = self.fc2(x)
        return x

# ResMLP blocks: a linear between patches + a MLP to process them independently
class ResMLP_Blocks(nn.Module):
    def __init__(self, nb_patches, dim, layerscale_init):
        super().__init__()
        self.affine_1 = Affine(dim)
        self.affine_2 = Affine(dim)
        self.linear_patches = nn.Linear(nb_patches, nb_patches) #Linear layer on patches
        self.mlp_channels = Mlp(dim) #MLP on channels
        self.layerscale_1 = nn.Parameter(layerscale_init * torch.ones((dim))) #LayerScale
        self.layerscale_2 = nn.Parameter(layerscale_init * torch.ones((dim))) # parameters

    def forward(self, x):
        res_1 = self.linear_patches(self.affine_1(x).transpose(1,2)).transpose(1,2)
        x = x + self.layerscale_1 * res_1
        res_2 = self.mlp_channels(self.affine_2(x))
        x = x + self.layerscale_2 * res_2
        return x

# ResMLP model: Stacking the full network
class ResMLP_models(nn.Module):
    def __init__(self, dim, depth, nb_patches, layerscale_init, num_classes):
        super().__init__()
        self.patch_projector = Patch_projector()
        self.blocks = nn.ModuleList([
            ResMLP_Blocks(nb_patches, dim, layerscale_init)
            for i in range(depth)])
        self.affine = Affine(dim)
        self.linear_classifier = nn.Linear(dim, num_classes)

    def forward(self, x):
        B, C, H, W = x.shape
        x = self.patch_projector(x)
        for blk in self.blocks:
            x = blk(x)
        x = self.affine(x)
        x = x.mean(dim=1).reshape(B,-1) #average pooling
        return self.linear_classifier(x)
```

---