

Going deeper with Image Transformers

Hugo Touvron^{*,†} Matthieu Cord[†] Alexandre Sablayrolles^{*}

Gabriel Synnaeve^{*} Hervé Jégou^{*}

^{*}Facebook AI [†]Sorbonne University

Abstract

Transformers have been recently adapted for large scale image classification, achieving high scores shaking up the long supremacy of convolutional neural networks. However the optimization of image transformers has been little studied so far. In this work, we build and optimize deeper transformer networks for image classification. In particular, we investigate the interplay of architecture and optimization of such dedicated transformers. We make two transformers architecture changes that significantly improve the accuracy of deep transformers. This leads us to produce models whose performance does not saturate early with more depth, for instance we obtain 86.3% top-1 accuracy on Imagenet when training with no external data. Our best model establishes the new state of the art on Imagenet with Reassessed labels and Imagenet-V2 / match frequency, in the setting with no additional training data.

1 Introduction

Residual architectures are prominent in computer vision since the advent of ResNet [27]. They are defined as a sequence of functions of the form

$$x_{l+1} = g_l(x_l) + R_l(x_l), \quad (1)$$

where the function g_l and R_l define how the network updates the input x_l at layer l . The function g_l is typically identity, while R_l is the main building block of the network: many variants in the literature essentially differ on how one defines this residual branch R_l is constructed or parametrized [50, 62, 72]. Residual architectures highlight the strong interplay between optimization and architecture design. As pointed out by He *et al.* [27], residual networks do not offer a better representational power. They achieve better performance because they are easier to train: shortly after their seminal work, He *et al.* discussed [28] the importance of having a clear path both forward and backward, and advocate setting g_l to the identity function.

The vision transformers [19] instantiate a particular form of residual architecture: after casting the input image into a set x_0 of vectors, the network alternates self-attention layers (SA) with feed-forward networks (FFN), as

$$\begin{aligned}x'_l &= x_l + \text{SA}(\eta(x_l)) \\x_{l+1} &= x'_l + \text{FFN}(\eta(x'_l))\end{aligned}\tag{2}$$

where η is the LayerNorm operator [1]. This definition follows the original architecture of Vaswani *et al.* [66], except the LayerNorm is applied before the block (*pre-norm*) in the residual branch, as advocated by He *et al.* [28]. Child *et al.* [13] adopt this choice with LayerNorm for training deeper transformers for various media, including for image generation where they train transformers with 48 layers.

How to normalize, weigh, or initialize the residual blocks of a residual architecture has received significant attention both for convolutional neural networks [7, 8, 28, 75] and for transformers applied to NLP or speech tasks [2, 34, 75]. In Section 2, we revisit this topic for transformer architectures solving image classification problems. Examples of approaches closely related to ours include Fixup [75], T-Fixup [34], ReZero [2] and SkipInit [16].

Following our analysis of the interplay between different initialization, optimization and architectural design, we propose an approach that is effective to improve the training of deeper architecture compared to current methods for image transformers. Formally, we add a learnable diagonal matrix on output of each residual block, initialized close to (but not at) 0. Adding this simple layer after each residual block improves the training dynamic, allowing us to train deeper high-capacity image transformers that benefit from depth. We refer to this approach as **LayerScale**.

Section 3 introduces our second contribution, namely **class-attention layers**, that we present in Figure 2. It is akin to an encoder/decoder architecture, in which we explicitly separate the transformer layers involving self-attention between patches, from class-attention layers that are devoted to extract the content of the processed patches into a single vector so that it can be fed to a linear classifier. This explicit separation avoids the contradictory objective of guiding the attention process while processing the class embedding. We refer to this new architecture as **CaiT** (Class-Attention in Image Transformers).

In the experimental Section 4, we empirically show the effectiveness and complementary of our approaches:

- LayerScale significantly facilitates the convergence and improves the accuracy of image transformers at larger depths. It adds a few thousands of parameters to the network at training time (negligible w.r.t. the total number of weights).
- Our architecture with specific class-attention offers a more effective processing of the class embedding.
- Our best CaiT models establish the new state of the art on Imagenet-Real [6] and Imagenet V2 matched frequency [52] with no additional

training data. On ImageNet1k-val [54], our model is close to the state of the art (86.3% vs 86.5%).

- We achieve competitive results on Transfer Learning.

We provide visualizations of the attention mechanisms in Section 5. We discuss related works along this paper and in the dedicated Section 6, before we conclude in Section 7. The appendices contain some variations we have tried during our exploration.

2 Deeper image transformers with LayerScale

Our goal is to increase the stability of the optimization when training transformers for image classification derived from the original architecture by Vaswani *et al.* [66], and especially when we increase their depth. We consider more specifically the vision transformer (ViT) architecture proposed by Dosovitskiy *et al.* [19] as the reference architecture and adopt the data-efficient image transformer (DeiT) optimization procedure of Touvron *et al.* [63]. In both works, there is no evidence that depth can bring any benefit when training on ImageNet only: the deeper ViT architectures have a low performance, while DeiT only considers transformers with 12 blocks of layers. The experimental section 4 will confirm that DeiT does not train deeper models effectively.

Figure 1 depicts the main variants that we compare for helping the optimization. They cover recent choices from the literature: as discussed in the introduction, the architecture (a) of ViT and DeiT is a pre-norm architecture [19, 63], in which the layer-normalisation η occurs at the beginning of the residual branch. Note that the original architecture of Vaswani *et al.* [66] applies the normalization after the block, but in our experiments the DeiT training does not converge with post-normalization.

Fixup [75], ReZero [2] and SkipInit [16] introduce learnable scalar weighting α_l on the output of residual blocks, while removing the pre-normalization and the warmup, see Figure 1(b). This amounts to modifying Eqn. 2 as

$$\begin{aligned} x'_l &= x_l + \alpha_l \text{SA}(x_l) \\ x_{l+1} &= x'_l + \alpha'_l \text{FFN}(x'_l). \end{aligned} \quad (3)$$

ReZero simply initializes this parameter to $\alpha = 0$. Fixup initializes this parameter to $\alpha = 1$ and makes other modifications: it adopts different policies for the initialization of the block weights, and adds several weights to the parametrization. In our experiments, these approaches do not converge even with some adjustment of the hyper-parameters.

Our empirical observation is that removing the warmup and the layer-normalization is what makes training unstable in Fixup and T-Fixup. Therefore we re-introduce these two ingredients so that Fixup and T-Fixup converge with DeiT models, see Figure 1(c). As we see in the experimental section, these

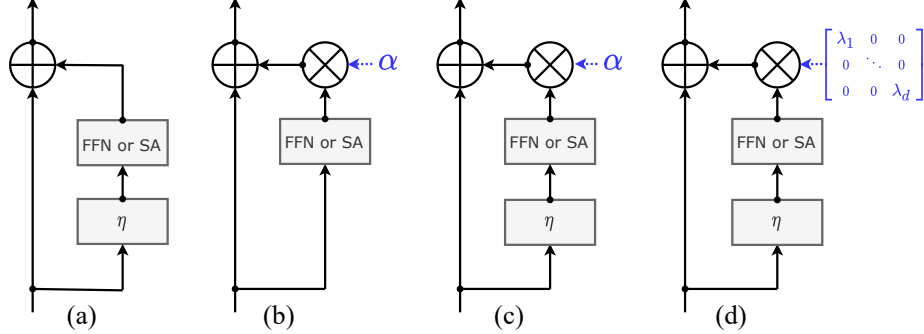


Figure 1: Normalization strategies for transformer blocks. (a) The ViT image classifier adopts pre-normalization like Child *et al.* [13]. (b) ReZero/Skipinit and Fixup remove the η normalization and the warmup (i.e., a reduced learning rate in the early training stage) and add a learnable scalar initialized to $\alpha=0$ and $\alpha=1$, respectively. Fixup additionally introduces biases and modifies the initialization of the linear layers. Since these methods do not converge with deep vision transformers, (c) we adapt them by reintroducing the pre-norm η and the warmup. Our main proposal (d) introduces a per-channel weighting (i.e, multiplication with a diagonal matrix $\text{diag}(\lambda_1, \dots, \lambda_d)$), where we initialize each weight with a small value as $\lambda_i = \varepsilon$.

amended variants of Fixup and T-Fixup are effective, mainly due to the learnable parameter α_l . When initialized at a small value, this choice does help the convergence when we increase the depth.

Our proposal LayerScale is a per-channel multiplication of the vector produced by each residual block, as opposed to a single scalar, see Figure 1(d). Our objective is to group the updates of the weights associated with the same output channel. Formally, LayerScale is a multiplication by a diagonal matrix on output of each residual block. In other terms, we modify Eqn. 2 as

$$\begin{aligned} x'_l &= x_l + \text{diag}(\lambda_{l,1}, \dots, \lambda_{l,d}) \times \text{SA}(\eta(x_l)) \\ x'_{l+1} &= x'_l + \text{diag}(\lambda'_{l,1}, \dots, \lambda'_{l,d}) \times \text{FFN}(\eta(x'_l)), \end{aligned} \quad (4)$$

where the parameters $\lambda_{l,i}$ and $\lambda'_{l,i}$ are learnable weights. The diagonal values are all initialized to a fixed small value ε : we set it to $\varepsilon = 0.1$ until depth 18, $\varepsilon = 10^{-5}$ for depth 24 and $\varepsilon = 10^{-6}$ for deeper networks. This formula is akin to other normalization strategies ActNorm [37] or LayerNorm but executed on output of the residual block. Yet we seek a different effect: ActNorm is a data-dependent initialization that calibrates activations so that they have zero-mean and unit variance, like batchnorm [35]. In contrast, we initialize the diagonal with small values so that the initial contribution of the residual branches to the function implemented by the transformer is small. In that respect our motivation is therefore closer to that of ReZero [2], SkipInit [16], Fixup [75] and T-Fixup [34]: to train closer to the identity function and let the network integrate the additional parameters progressively during the training. LayerScale

offers more diversity in the optimization than just adjusting the whole layer by a single learnable scalar as in ReZero/SkipInit, Fixup and T-Fixup. As we will show empirically, offering the degrees of freedom to do so per channel is a decisive advantage of LayerScale over existing approaches. In Appendix A, we present other variants or intermediate choices that support our proposal, and a control experiment that aims at disentangling the specific weighting of the branches of LayerScale from its impact on optimization.

Formally, adding these weights does not change the expressive power of the architecture since, as they can be integrated into the previous matrix of the SA and FFN layers without changing the function implemented by the network.

3 Specializing layers for class attention

In this section, we introduce the CaiT architecture, depicted in Figure 2 (right). This design aims at circumventing one of the problems of the ViT architecture: the learned weights are asked to optimize two contradictory objectives: (1) guiding the self-attention between patches while (2) summarizing the information useful to the linear classifier. Our proposal is to explicitly separate the two stages, in the spirit of an encoder-decoder architecture, see Section 6.

Later class token. As an intermediate step towards our proposal, we insert the so-called class token, denoted by CLS, later in the transformer. This choice eliminates the discrepancy on the first layers of the transformer, which are therefore fully employed for performing self-attention between patches only. As a baseline that does not suffer from the contradictory objectives, we also consider **average pooling** of all the patches on output of the transformers, as typically employed in convolutional architectures.

Architecture. Our CaiT network consists of two distinct processing stages visible in Figure 2:

1. The *self-attention* stage is identical to the ViT transformer, but with no class embedding (CLS).
2. The *class-attention* stage is a set of layers that compiles the set of patch embeddings into a class embedding CLS that is subsequently fed to a linear classifier.

This class-attention alternates in turn a layer that we refer to as a multi-head class-attention (CA), and a FFN layer. In this stage, only the class embedding is updated. Similar to the one fed in ViT and DeiT on input of the transformer, it is a learnable vector. The main difference is that, in our architecture, we do not copy information from the class embedding to the patch embeddings during the forward pass. Only the class embedding is updated by residual in the CA and FFN processing of the class-attention stage.

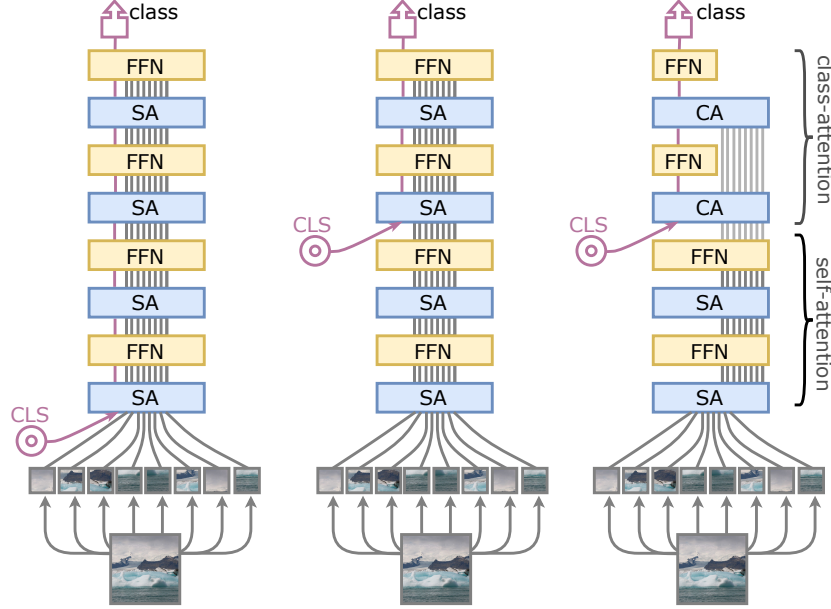


Figure 2: In the ViT transformer (*left*), the class embedding (CLS) is inserted along with the patch embeddings. This choice is detrimental, as the same weights are used for two different purposes: helping the attention process, and preparing the vector to be fed to the classifier. We put this problem in evidence by showing that inserting CLS later improves performance (*middle*). In the CaiT architecture (*right*), we further propose to freeze the patch embeddings when inserting CLS to save compute, so that the last part of the network (typically 2 layers) is fully devoted to summarizing the information to be fed to the linear classifier.

Multi-heads class attention. The role of the CA layer is to extract the information from the set of processed patches. It is identical to a SA layer, except that it relies on the attention between (i) the class embedding x_{class} (initialized at CLS in the first CA) and (ii) itself plus the set of frozen patch embeddings x_{patches} . We discuss why we include x_{class} in the keys in Appendix B.

Considering a network with h heads and p patches, and denoting by d the embedding size, we parametrize the multi-head class-attention with several projection matrices, $W_q, W_k, W_v, W_o \in \mathbf{R}^{d \times d}$, and the corresponding biases $b_q, b_k, b_v, b_o \in \mathbf{R}^d$. With this notation, the computation of the CA residual block proceeds as follows. We first augment the patch embeddings (in matrix form) as $z = [x_{\text{class}}, x_{\text{patches}}]$ (see Appendix B for results when $z = x_{\text{patches}}$). We then perform the projections:

$$Q = W_q x_{\text{class}} + b_q, \quad (5)$$

$$K = W_k z + b_k, \quad (6)$$

$$V = W_v z + b_v. \quad (7)$$

The class-attention weights are given by

$$A = \text{Softmax}(Q.K^T / \sqrt{d/h}) \quad (8)$$

where $Q.K^T \in \mathbf{R}^{h \times 1 \times p}$. This attention is involved in the weighted sum $A \times V$ to produce the residual output vector

$$\text{out}_{\text{CA}} = W_o A V + b_o, \quad (9)$$

which is in turn added to x_{class} for subsequent processing.

The CA layers extract the useful information from the patches embedding to the class embedding. In preliminary experiments, we empirically observed that the first CA and FFN give the main boost, and a set of 2 blocks of layers (2 CA and 2 FFN) is sufficient to cap the performance. In the experimental section, we denote by 12+2 a transformer when it consists of 12 blocks of SA+FFN layers and 2 blocks of CA+FFN layers.

Complexity. The layers contain the same number of parameters in the class-attention and self-attention stages: CA is identical to SA in that respect, and we use the same parametrization for the FFNs. However the processing of these layers is much faster: the FFN only processes matrix-vector multiplications.

The CA function is also less expensive than SA in term of memory and computation because it computes the attention between the class vector and the set of patch embeddings: $Q \in \mathbf{R}^d$ means that $Q.K^T \in \mathbf{R}^{h \times 1 \times p}$. In contrast, in the “regular self-attention” layers SA, we have $Q \in \mathbf{R}^{p \times d}$ and therefore $Q.K^T \in \mathbf{R}^{h \times p \times p}$. In other words, the initially quadratic complexity in the number of patches becomes linear in our extra CaiT layers.

4 Experiments

In this section, we report our experimental results related to LayerScale and CaiT. We first study strategies to train at deeper scale in Section 4.1, including our LayerScale method. Section 4.2 shows the interest of our class-attention design. We present our models in Subsection 4.3. Section 4.4 details our results on Imagenet and Transfer learning. We provide an ablation of hyper-parameter and ingredients in Section 4.5. Note, in Appendix A and B, we provide variations on our methods and corresponding results.

Experimental setting. Our implementation is based on the Timm library [68]. Unless specified otherwise, for this analysis we make minimal changes to hyper-parameters compared to the DeiT training scheme [63], although we have tried minor adjustments of the learning rate in case we did not observe convergence. In order to speed up training and optimize memory consumption we have used a sharded training provided by the Fairscale library¹ with fp16 precision.

¹<https://pypi.org/project/fairscale/>

4.1 Preliminary analysis with deeper architectures

In our early experiments, we observe that Vision Transformers become increasingly more difficult to train when we scale architectures. Depth is one of the main source of instability. For instance the DeiT procedure [63] fails to properly converge above 18 layers without adjusting hyper-parameters. Large ViT [19] models with 24 and 32 layers were trained with large training datasets, but when trained on Imagenet only the larger models are not competitive.

In the following, we analyse various ways to stabilize the training with different architectures. At this stage we consider a DeiT-Small model² during 300 epochs to allow a direct comparison with the results reports by Touvron *et al.* [63]. We measure the performance on the Imagenet1k [17, 54] classification dataset as a function of the depth.

4.1.1 Adjusting the drop-rate of stochastic depth.

The first step to improve convergence is to adapt the hyper-parameters that interact the most with depth, in particular Stochastic depth [33]. This method is already popular in NLP [21, 22] to train deeper architectures. For ViT, it was first proposed by Wightman *et al.* [68] in the Timm implementation, and subsequently adopted in DeiT [63]. The per-layer drop-rate depends linearly on the layer depth, but in our experiments this choice does not provide an advantage compared to the simpler choice of a uniform drop-rate d_r . In Table 1 we show that the default stochastic depth of DeiT allows us to train up to 18 blocks of SA+FFN. After that the training becomes unstable. By increasing the drop-rate hyper-parameter d_r , the performance increases until 24 layers. It saturates at 36 layers (we measured that it drops to 80.7% at 48 layers).

4.1.2 Comparison of normalization strategies

We carry out an empirical study of the normalization methods discussed in Section 2. As previously indicated, Rezero, Fixup and T-Fixup do not converge when training DeiT off-the-shelf. However, if we re-introduce Layer-Norm³ and warmup, Fixup and T-Fixup achieve convergence and even improve training compared to the baseline DeiT. We report the results for these “adaptations” of Fixup and T-Fixup in Table 1.

The modified methods are able to converge with more layers without saturating too early. ReZero converges, we show (column $\alpha = \epsilon$) that it is better to initialize α to a small value instead of 0, as in LayerScale. All the methods have a beneficial effect on convergence and they tend to reduce the need for stochastic depth, therefore we adjust these drop rate accordingly per method. Figure 3 provides the performance as the function of the drop rate d_r for LayerScale. We empirically use the following formula to set up the drop-rate for the

²<https://github.com/facebookresearch/deit>

³Bachlechner *et al.* report that batchnorm is complementary to ReZero, while removing Layer-Norm in the case of transformers.

Table 1: **Improving convergence at depth** on ImageNet-1k. The baseline is DeiT-S with uniform drop rate of $d = 0.05$ (same expected drop rate and performance as progressive stochastic depth of 0.1). Several methods include a fix scalar learnable weight α per layer as in Figure 1(c). We have adapted Rezero, Fixup, T-Fixup, since the original methods do not converge: we have re-introduced the Layer-normalization η and warmup. We have adapted the drop rate d_r for all the methods, including the baseline. The column $\alpha = \varepsilon$ reports the performance when initializing the scalar with the same value as for LayerScale. †: *failed before the end of the training*.

depth	baseline		scalar α weighting				LayerScale
	$d_r = 0.05$	adjust [d_r]	Rezero	T-Fixup	Fixup	$\alpha = \varepsilon$	
12	79.9	79.9 [0.05]	78.3	79.4	80.7	80.4	80.5
18	80.1	80.7 [0.10]	80.1	81.7	82.0	81.6	81.7
24	78.9†	81.0 [0.20]	80.8	81.5	82.3	81.1	82.4
36	78.9†	81.9 [0.25]	81.6	82.1	82.4	81.6	82.9

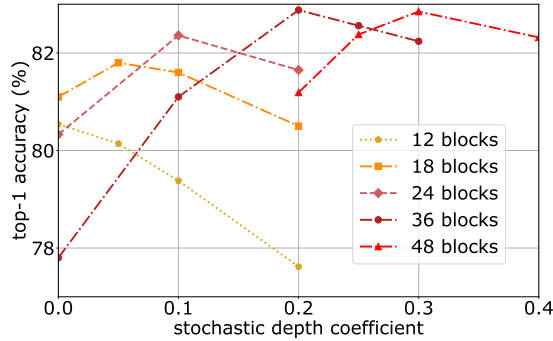


Figure 3: We measure the impact of stochastic depth on ImageNet with a DeiT-S with LayerScale for different depths. The drop rate of stochastic depth needs to be adapted to the network depth.

CaiT-S models derived from on DeiT-S: $d_r = \min(0.1 \times \frac{\text{depth}}{12} - 1, 0)$. This formulaic choice avoids cross-validating this parameter and overfitting, yet it does not generalize to models with different d : We further increase (resp. decrease) it by a constant for larger (resp. smaller) working dimensionality d .

Fixup and T-Fixup are competitive with LayerScale in the regime of a relatively low number of blocks (12–18). However, they are more complex than LayerScale: they employ different initialization rules depending of the type of layers, and they require more changes to the transformer architecture. Therefore we only use LayerScale in subsequent experiments. It is much simpler and parametrized by a single hyper-parameter ε , and it offers a better performance for the deepest models that we consider, which are also the more accurate.

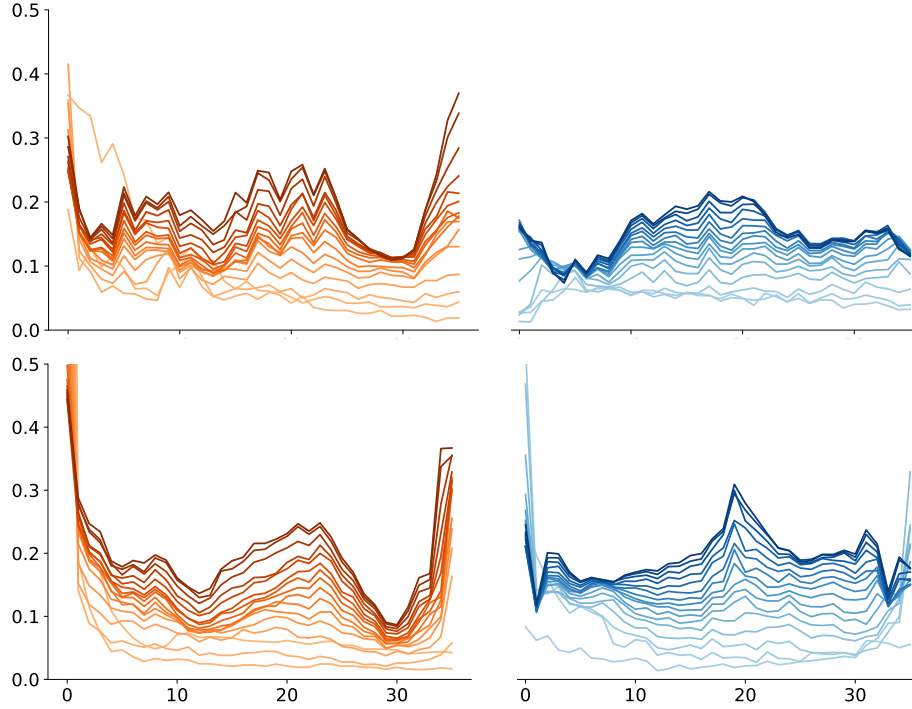


Figure 4: Analysis of the contribution of the residual branches (*Top*: Self-attention ; *Bottom*: FFN) for a network comprising 36 layers, **without** (red) or **with** (blue) Layerscale. The ratio between the norm of the residual and the norm of the main branch is shown for each layer of the transformer and for various epochs (darker shades correspond to the last epochs). For the model trained with layerscale, the norm of the residual branch is on average 20% of the norm of the main branch. We observe that the contribution of the residual blocks fluctuates more for the model trained without layerscale and in particular is lower for some of the deeper layers.

4.1.3 Analysis of Layerscale

Statistics of branch weighting. We evaluate the impact of Layerscale for a 36-blocks transformer by measuring the ratio between the norm of the residual activations and the norm of the activations of the main branch $\|g_l(x)\|_2/\|x\|_2$. The results are shown in Figure 4. We can see that training a model with Layerscale makes this ratio more uniform across layers, and seems to prevent some layers from having a disproportionate impact on the activations. Similar to prior works [2, 75] we hypothesize that the benefit is mostly the impact on optimization. This hypothesis is supported by the control experiment that we detail in Appendix A.

Table 2: Variations on CLS with DeiT-Small (no LayerScale): we change the layer at which the class embedding is inserted. In ViT and DeiT, it is inserted at layer 0 jointly with the projected patches. We evaluate a late insertion of the class embedding, as well as our design choice to introduce specific class-attention layers.

depth: SA+CA	insertion layer	top-1 acc.	#params	FLOPs
Baselines: DeiT-S and average pooling				
12: 12 + 0	0	79.9	22M	4.6B
12: 12 + 0	n/a	80.3	22M	4.6B
Late insertion of class embedding				
12: 12 + 0	2	80.0	22M	4.6B
12: 12 + 0	4	80.0	22M	4.6B
12: 12 + 0	8	80.0	22M	4.6B
12: 12 + 0	10	80.5	22M	4.6B
12: 12 + 0	11	80.3	22M	4.6B
DeiT-S with class-attention stage (SA+FFN)				
12: 9 + 3	9	79.6	22M	3.6B
12: 10 + 2	10	80.3	22M	4.0B
12: 11 + 1	11	80.6	22M	4.3B
13: 12 + 1	12	80.8	24M	4.7B
14: 12 + 2	12	80.8	26M	4.7B
15: 12 + 3	12	80.6	27M	4.8B

4.2 Class-attention layers

In Table 2 we study the impact on performance of the design choices related to class embedding. We depict some of them in Figure 2. As a baseline, average pooling of patches embeddings with a vanilla DeiT-Small achieves a better performance than using a class token. This choice, which does not employ any class embedding, is typical in convolutional networks, but possibly weaker with transformers when transferring to other tasks [20].

Late insertion. The performance increases when we insert the class embedding later in the transformer. It is maximized two layers before the output. Our interpretation is that the attention process is less perturbed in the 10 first layers, yet it is best to keep 2 layers for compiling the patches embedding into the class embedding via class-attention, otherwise the processing gets closer to a weighted average.

Our class-attention layers are designed on the assumption that there is no benefit in copying information from the class embedding back to the patch embeddings in the forward pass. Table 2 supports that hypothesis: if we compare the performance for a total number of layers fixed to 12, the performance of CaiT with 10 SA and 2 CA layers is identical to average pooling and better

Table 3: CaiT models: The design parameters are depth and d . The mem columns correspond to the memory usage. All models are initially trained at resolution 224 during 400 epochs. We also fine-tune these models at resolution 384 (identified by $\uparrow 384$) or train them with distillation (Υ). The FLOPs are reported for each resolution.

CAiT model	depth (SA+CA)	d	#params ($\times 10^6$)	FLOPs ($\times 10^9$)		Top-1 acc. (%): Imagenet1k-val			
				@224	@384	@224	$\uparrow 384$	@224 Υ	$\uparrow 384\Upsilon$
XXS-24	24+2	192	12.0	2.5	9.5	77.6	80.4	78.4	80.9
XXS-36	36+2	192	17.3	3.8	14.2	79.1	81.8	79.7	82.2
XS-24	24+2	288	26.6	5.4	19.3	81.8	83.8	82.0	84.1
XS-36	36+2	288	38.6	8.1	28.8	82.6	84.3	82.9	84.8
S-24	24+2	384	46.9	9.4	32.2	82.7	84.3	83.5	85.1
S-36	36+2	384	68.2	13.9	48.0	83.3	85.0	84.0	85.4
S-48	48+2	384	89.5	18.6	63.8	83.5	85.1	83.9	85.3
M-24	24+2	768	185.9	36.0	116.1	83.4	84.5	84.7	85.8
M-36	36+2	768	270.9	53.7	173.3	83.8	84.9	85.1	86.1

than the DeiT-Small baseline with a lower number of FLOPs. If we set 12 layers in the self-attention stage, which dominates the complexity, we increase the performance significantly by adding two blocks of CA+FFN.

4.3 Our CaiT models

Our CaiT models are built upon ViT: the only difference is that we incorporate LayerScale in each residual block (see Section 2) and the two-stages architecture with class-attention layers described in Section 3. Table 3 describes our different models. The design parameters governing the capacity are the depth and the working dimensionality d . In our case d is related to the number of heads h as $d = 48 \times h$, since we fix the **number of components per head** to 48. This choice is a bit smaller than the value used in DeiT. We also adopt the **crop-ratio** of 1.0 optimized for DeiT by Wightman [68]. Table 9 and 10 in the ablation section 4.5 support these choices.

We incorporate talking-heads attention [55] into our model. It increases the performance on Imagenet of DeiT-Small from 79.9% to 80.3%.

The hyper-parameters are identical to those provided in DeiT [63], except mentioned otherwise. We use a batch size of 1024 samples and train during 400 epochs with repeated augmentation [5, 29]. The learning rate of the AdamW optimizer [44] is set to 0.001 and associated with a cosine training schedule, 5 epochs of warmup and a weight decay of 0.05. We report in Table 4 the two hyper-parameters that we modify depending on the model complexity, namely the drop rate d_r associated with uniform stochastic depth, and the initialization value ε of LayerScale.

Table 4: Hyper-parameters for training CaiT models: The only parameters that we adjust per model are the drop rate d_r of stochastic depth and the LayerScale initialization ε .

CAiT model	hparams	
	d_r	ε
XXS-24	0.05	10^{-5}
XXS-36	0.1	10^{-6}
XS-24	0.05	10^{-5}
XS-36	0.1	10^{-6}
S-24	0.1	10^{-5}
S-36	0.2	10^{-6}
S-48	0.3	10^{-6}
M-24	0.2	10^{-5}
M-36	0.3	10^{-6}

Fine-tuning at higher resolution (\uparrow) and distillation (Υ). We train all our models at resolution 224, and optionally fine-tune them at a higher resolution to trade performance against accuracy [19, 63, 64]: we denote the model by $\uparrow 384$ models fine-tuned at resolution 384×384 . We also train models with distillation (Υ) as suggested by Touvron *et al.* [63]. We use a RegNet-16GF [50] as teacher and adopt the “hard distillation” [63] for its simplicity.

4.4 Results

4.4.1 Performance/complexity of CaiT models

Table 3 provides different complexity measures for our models. As a general observation, we observe a subtle interplay between the width and the depth, both contribute to the performance as reported by Dosovitskiy *et al.* [19] with longer training schedules. But if one parameter is too small the gain brought by increasing the other is not worth the additional complexity.

Fine-tuning to size 384 (\uparrow) systematically offers a large boost in performance without changing the number of parameters. It also comes with a higher computational cost. In contrast, leveraging a pre-trained convnet teacher with hard distillation as suggested by Touvron *et al.* [63] provides a boost in accuracy without affecting the number of parameters nor the speed.

4.4.2 Comparison with the state of the art on Imagenet

Our main classification experiments are carried out on ImageNet [54], and also evaluated on two variations of this dataset: ImageNet-Real [6] that corrects and give a more detailed annotation, and ImageNet-V2 [52] (matched frequency) that provides a separate test set. In Table 5 we compare some of our models with the state of the art on Imagenet classification when training without

Table 5: **Throughput (images/s) vs accuracy** on Imagenet [54], Imagenet Real [6] and Imagenet V2 matched frequency [52] for models trained without external data. We compare CaiT with DeiT [63], ViT-B [19], TNT [26], T2T [74] and to several state-of-the-art convnets: Regnet [50] improved by Touvron et al. [63], EfficientNet [14, 62, 71], Fix-EfficientNet [65] and NFNet [8]. Most reported results are from corresponding papers, and therefore the training procedure differs for the different models. For Imagenet V2 matched frequency and Imagenet Real we report the results provided by the authors. When not available (like NFNet), we report the results measured by Wighthman [68] with converted models. The RegNetY-16GF is the teacher model that we trained for distillation. We report the best result in **bold** and the second best result(s) underlined.

Network	nb of param.	nb of FLOPs	image size		ImNet top-1	Real top-1	V2 top-1
RegNetY-16GF	84M	16.0B	224	224	82.9	88.1	72.4
EfficientNet-B5	30M	9.9B	456	456	83.6	88.3	73.6
EfficientNet-B7	66M	37.0B	600	600	84.3	-	-
EfficientNet-B5 RA	30M	9.9B	456	456	83.7	-	-
EfficientNet-B7 RA	66M	37.0B	600	600	84.7	-	-
EfficientNet-B7 AdvProp	66M	37.0B	600	600	85.2	89.4	76.0
Fix-EfficientNet-B8	87M	89.5B	672	800	85.7	<u>90.0</u>	75.9
NFNet-F0	72M	12.4B	192	256	83.6	88.1	72.6
NFNet-F1	133M	35.5B	224	320	84.7	88.9	74.4
NFNet-F2	194M	62.6B	256	352	85.1	88.9	74.3
NFNet-F3	255M	114.8B	320	416	85.7	89.4	75.2
NFNet-F4	316M	215.3B	384	512	85.9	89.4	75.2
NFNet-F5	377M	289.8B	416	544	86.0	89.2	74.6
NFNet-F6+SAM	438M	377.3B	448	576	86.5	89.9	75.8
Transformers							
ViT-B/16	86M	55.4B	24	384	77.9	83.6	-
ViT-L/16	307M	190.7B	224	384	76.5	82.2	-
T2T-ViT t-14	21M	5.2B	224	224	80.7	-	-
TNT-S	24M	5.2B	224	224	81.3	-	-
TNT-S + SE	25M	5.2B	224	224	81.6	-	-
TNT-B	66M	14.1B	224	224	82.8	-	-
DeiT-S	22M	4.6B	224	224	79.8	85.7	68.5
DeiT-B	86M	17.5B	224	224	81.8	86.7	71.5
DeiT-B \uparrow 384	86M	55.4B	224	384	83.1	87.7	72.4
DeiT-B \uparrow 384 Υ 1000 epochs	87M	55.5B	224	384	85.2	89.3	75.2
Our deep transformers							
CaiT-S36	68M	13.9B	224	224	83.3	88.0	72.5
CaiT-S36 \uparrow 384	68M	48.0B	224	384	85.0	89.2	75.0
CaiT-S48 \uparrow 384	89M	63.8B	224	384	85.1	89.5	75.5
CaiT-S36 Υ	68M	13.9B	224	224	84.0	88.9	74.1
CaiT-S36 \uparrow 384 Υ	68M	48.0B	224	384	85.4	89.8	76.2
CaiT-M36 Υ	271M	53.7B	224	224	85.1	89.3	74.9
CaiT-M36 \uparrow 384 Υ	271M	173.3B	224	384	86.1	<u>90.0</u>	<u>76.3</u>
CaiT-M36 \uparrow 448 Υ	271M	247.8B	224	448	<u>86.3</u>	90.2	76.7

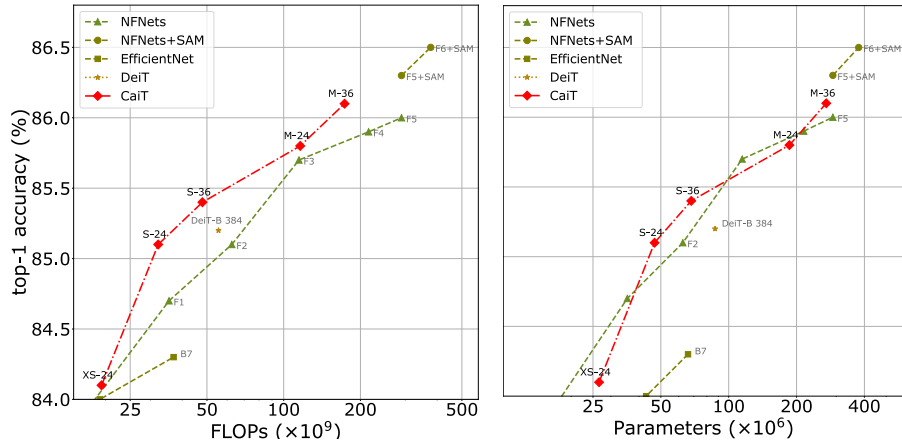


Figure 5: We represent FLOPs and parameters for our best CaiT \uparrow 384 Υ models trained with distillation. They are competitive on Imagenet1k-val with the state of the art in the high accuracy regime, from XS-24 to M-36. Convolution-based neural networks like NFNets and EfficientNet are better in low-FLOPs and low-parameters regimes.

external data. We focus on the models CaiT-S36 and CaiT-M36, at different resolutions and with or without distillation.

On Imagenet1k-val, we achieve 86.3% of top-1 accuracy, which is a significant improvement over DeiT (85.2%). It is the state of the art except for a very recent concurrent work [8] that reports a top-1 accuracy of 86.5%. Furthermore, our CaiT-M36 \uparrow 384 Υ model obtains 86.1% top-1 accuracy on Imagenet-1k-val: this network is more efficient than the best performing NFNet convnets w.r.t. FLOPs and even more in terms of throughput (images processed per second on a V100 GPU), thanks to the lower memory usage. Our approach outperforms the state of the art on Imagenet with reassessed labels, and on Imagenet-V2, which has a distinct validation set which makes it harder to overfit.

4.4.3 Transfer learning

We evaluated our method on transfer learning tasks by fine-tuning on the datasets in Table 6.

Fine-tuning procedure. For fine-tuning we use the same hyperparameters as for training. We only decrease the learning rates by a factor 10 (for CARS, Flowers, iNaturalist), 100 (for CIFAR-100, CIFAR-10) and adapt the number of epochs (1000 for CIFAR-100, CIFAR-10, Flowers-102 and Cars-196, 360 for iNaturalist 2018 and 2019). We have not used distillation for this finetuning.

Results. Table 7 compares CaiT transfer learning results to those of EfficientNet [62], ViT [19] and DeiT [63]. These results show the excellent generalization

Table 6: Datasets used for our different tasks.

Dataset	Train size	Test size	#classes
ImageNet [54]	1,281,167	50,000	1000
iNaturalist 2018 [30]	437,513	24,426	8,142
iNaturalist 2019 [31]	265,240	3,003	1,010
Flowers-102 [46]	2,040	6,149	102
Stanford Cars [38]	8,144	8,041	196
CIFAR-100 [39]	50,000	10,000	100
CIFAR-10 [39]	50,000	10,000	10

Table 7: Results in transfer learning. All our CaiT models are trained and evaluated at resolution 224.

Model	ImageNet	CIFAR-10	CIFAR-100	Flowers	Cars	iNat-18	iNat-19	FLOPs
EfficientNet-B7	84.3	98.9	91.7	98.8	94.7	-	-	37.0B
ViT-B/16	77.9	98.1	87.1	89.5	-	-	-	55.5B
ViT-L/16	76.5	97.9	86.4	89.7	-	-	-	190.7B
DeiT-B 224	81.8	99.1	90.8	98.4	92.1	73.2	77.7	17.5B
CaiT-S-36 224	83.4	99.2	92.2	98.8	93.5	77.1	80.6	13.9B
CaiT-M-36 224	83.7	99.3	93.3	99.0	93.5	76.9	81.7	53.7B
CaiT-S-36 Υ 224	83.7	99.2	92.2	99.0	94.1	77.0	81.4	13.9B
CaiT-M-36 Υ 224	84.8	99.4	93.1	99.1	94.2	78.0	81.8	53.7B

of the transformers-based models in general. Our CaiT models achieve excellent results, as shown by the overall better performance than EfficientNet-B7 across datasets.

4.5 Ablation

In this section we provide different sets of ablation, in the form of a transition from DeiT to CaiT. Then we provide experiments that have guided our hyperparameter optimization. As mentioned in the main paper, we use the same hyperparameters as in DeiT [63] everywhere except stated otherwise. We have only changed the number of attention for a given working dimension (see Section 4.5.2), and changed the crop-ratio (see Section 4.5.3).

4.5.1 Step by step from DeiT-Small to CaiT-S36

In Table 8 we present how to gradually transform the DeiT-S [63] architecture into CaiT-36, and measure at each step the performance/complexity changes.

Table 8: Ablation: we present the ablation path from DeiT-S to our CaiT models. We highlight the complementarity of our approaches and optimized hyper-parameters. Note, Fine-tuning at higher resolution supersedes the inference at higher resolution. †: training failed.

Improvement	top-1 acc.	#params	FLOPs
DeiT-S [$d=384, 300$ epochs]	79.9	22M	4.6B
+ More heads [8]	80.0	22M	4.6B
+ Talking-heads	80.5	22M	4.6B
+ Depth [36 blocks]	69.9†	64M	13.8B
+ Layer-scale [$init\ \varepsilon = 10^{-6}$]	80.5	64M	13.8B
+ Stch depth. adaptation [$d_r=0.2$]	83.0	64M	13.8B
+ CaiT architecture [<i>specialized class-attention layers</i>]	83.2	68M	13.9B
+ Longer training [400 epochs]	83.4	68M	13.9B
+ Inference at higher resolution [256]	83.8	68M	18.6B
+ Fine-tuning at higher resolution [384]	84.8	68M	48.0B
+ Hard distillation [<i>teacher: RegNetY-16GF</i>]	85.2	68M	48.0B
+ Adjust crop ratio [$0.875 \rightarrow 1.0$]	85.4	68M	48.0B

One can see that CaiT is complementary with LayerScale and offers an improvement without significantly increasing the FLOPs. As already reported in the literature, the resolution is another important step for improving the performance and fine-tuning instead of training the model from scratch saves a lot of computation at training time. Last but not least, our models benefit from longer training schedules.

4.5.2 Optimization of the number of heads

In Table 9 we study the impact of the number of heads for a fixed working dimensionality. This architectural parameter has an impact on both the accuracy, and the efficiency: while the number of FLOPs remain roughly the same, the compute is more fragmented when increasing this number of heads and on typical hardware this leads to a lower effective throughput. Choosing 8 heads in the self-attention offers a good compromise between accuracy and speed. In DeiT-Small, this parameter was set to 6.

4.5.3 Adaptation of the crop-ratio

In the typical (“center-crop”) evaluation setting, most convolutional neural networks crop a subimage with a given ratio, typically extracting a 224×224 center crop from a 256×256 resized image, leading to the typical ratio of 0.875. Wightman *et al.* [68] notice that setting this crop ratio to 1.0 for transformer models has a positive impact: the distilled DeiT-B \uparrow 384 reaches a top1-accuracy on Imagenet1k-val of 85.42% in this setting, which is a gain of +0.2% compared to the accuracy of 85.2% reported by Touvron *et al.* [63].

Table 9: DeiT-Small: for a fixed 384 working dimensionality and number of parameters, impact of the number of heads on the accuracy and throughput (images processed per second at inference time on a single V100 GPU).

# heads	dim/head	throughput (im/s)	GFLOPs	top-1 acc.
1	384	1079	4.6	76.80
2	192	1056	4.6	78.06
3	128	1043	4.6	79.35
6	64	989	4.6	79.90
8	48	971	4.6	80.02
12	32	927	4.6	80.08
16	24	860	4.6	80.04
24	16	763	4.6	79.60

Our measurements concur with this observation: We observe a gain for almost all our models and most of the evaluation benchmarks. For instance our model M36 \uparrow 384 Υ increases to 86.1% top-1 accuracy on Imagenet-val1k.

4.5.4 Longer training schedules

As shown in Table 8, increasing the number of training epochs from 300 to 400 improves the performance of CaiT-S-36. However, increasing the number of training epochs from 400 to 500 does not change performance significantly (83.44 with 400 epochs 83.42 with 500 epochs). This is consistent with the observation of the DeiT [63] paper, which notes a saturation of performance from 400 epochs for the models trained without distillation.

5 Visualizations

5.1 Attention map

In Figure 6 we show the attention maps associated with the individual 4 heads of a XXS CaiT model, and for the two layers of class-attention. In CaiT and in contrast to ViT, the class-attention stage is the only one where there is some interaction between the class token and the patches, therefore it conveniently concentrates all the spatial-class relationship. We make two observations:

- The first class-attention layer clearly focuses on the object of interest, corresponding to the main part of the image on which the classification decision is performed (either correct or incorrect). In this layer, the different heads focus either on the same or on complementary parts of the objects. This is especially visible for the waterfall image;
- The second class-attention layer seems to focus more on the context, or at least the image more globally.

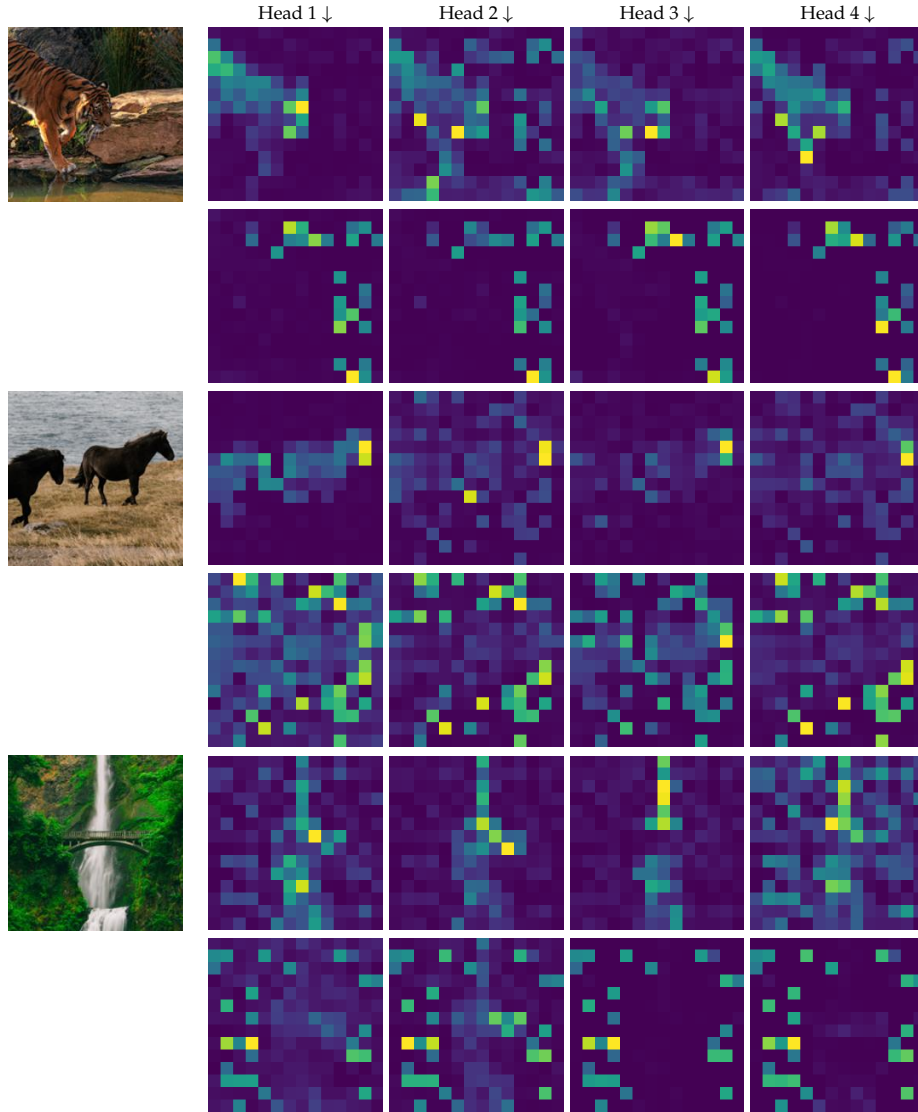


Figure 6: Visualization of the attention maps in the class-attention stage, obtained with a XXS model. For each image we present two rows: the top row correspond to the four heads of the attention maps associated with the first CA layer. The bottom row correspond to the four heads of the second CA layer.

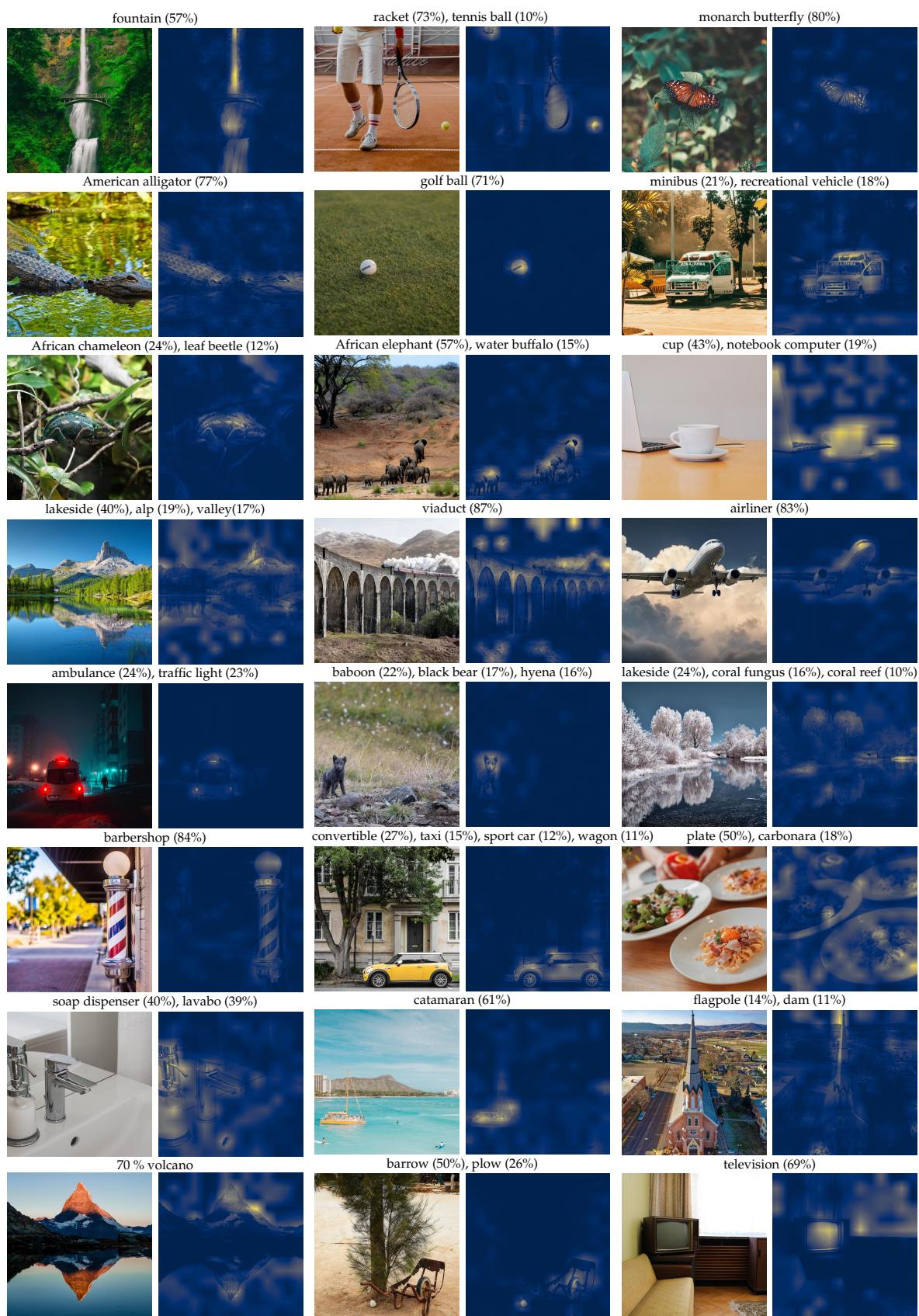


Figure 7: Illustration of the regions of focus of a CaiT-XXS model, according to the response of the first class-attention layer.

Table 10: We compare performance with the default crop-ratio of 0.875 usually used with convnets, and the crop-ratio of 1.0 [68] that we adopt for CaiT.

Network	Crop Ratio		ImNet	Real	V2
	0.875	1.0	top-1	top-1	top-1
S36	✓	-	83.4	88.1	73.0
	-	✓	83.3	88.0	72.5
S36↑384	✓	-	84.8	88.9	74.7
	-	✓	85.0	89.2	75.0
S36Υ	✓	-	83.7	88.9	74.1
	-	✓	84.0	88.9	74.1
M36Υ	✓	-	84.8	89.2	74.9
	-	✓	84.9	89.2	75.0
S36↑384Υ	✓	-	85.2	89.7	75.7
	-	✓	85.4	89.8	76.2
M36↑384Υ	✓	-	85.9	89.9	76.1
	-	✓	86.1	90.0	76.3

5.2 Illustration of saliency in class-attention

In figure 7 we provide more visualisations for a XXS model. They are just illustration of the saliency that one may extract from the first class-attention layer. As discussed previously this layer is the one that, empirically, is the most related to the object of interest. To produce these visual representations we simply average the attention maps from the different heads (depicted in Figure 6), and upsample the resulting map to the image size. We then modulate the gray-level image with the strength of the attention after normalizing it with a simple rule of the form $(x - x_{\min}) / (x_{\max} - x_{\min})$. We display the resulting image with `cividis` colormap.

For each image we show this saliency map and provides all the class for which the model assigns a probability higher than 10%. These visualizations illustrate how the model can focus on two distinct regions (like racket and tennis ball on the top row/center). We can also observe some failure cases, like the top of the church classified as a flagpole.

6 Related work

Since AlexNet [40], convolutional neural networks (CNN) are the standard in image classification [27, 62, 64], and more generally in computer vision. While a deep CNN can theoretically model long range interaction between pixels across many layers, there has been research in increasing the range of interactions within a single layer. Some approaches adapt the receptive field of convolutions dynamically [15, 42]. At another end of the spectrum, attention

can be viewed as a general form of non-local means, which was used in filtering (e.g. denoising [10]), and more recently in conjunction with convolutions [67]. Various other attention mechanism have been used successfully to give a global view in conjunction with (local) convolutions [4, 51, 12, 76, 78], most mimic *squeeze-and-excite* [32] for leveraging global features. Lastly, LambdaNetworks [3] decomposes attention into an approximated content attention and a batch-amortized positional attention component.

Hybrid architectures combining CNNs and transformers blocks have also been used on ImageNet [57, 69] and on COCO [11]. Originally, transformers without convolutions were applied on pixels directly [47], even scaling to hundred of layers [13], but did not perform at CNNs levels. More recently, a transformer architecture working directly on small patches has obtained state of the art results on ImageNet [19]. Nevertheless, the state of the art has since returned to CNNs [8, 48]. While some small improvements have been applied on the transformer architecture with encouraging results [74], their performance is below the one of DeiT [63], which uses a vanilla ViT architecture.

Encoder/decoder architectures. Transformers were originally introduced for machine translation [66] with encoder-decoder models, and gained popularity as masked language model encoders (BERT) [18, 43]. They yielded impressive results as scaled up language models, e.g. GPT-2 and 3 [49, 9]. They became a staple in speech recognition too [45, 36], being it in encoder and sequence criterion or encoder-decoder seq2seq [60] conformations, and hold the state of the art to this day [73, 77] with models 36 blocks deep. Note, transforming only the class token with frozen trunk embeddings in CaiT is reminiscent of non-autoregressive encoder-decoders [25, 41], where a whole sequence (we have only one prediction) is produced at once by iterative refinements.

Deeper architectures usually lead to better performance [27, 56, 61], however this complicates their training process [58, 59]. One must adapt the architecture and the optimization procedure to train them correctly. Some approaches focus on the initialization schemes [24, 27, 70], others on multiple stages training [53, 56], multiple loss at different depth [61], adding components in the architecture [2, 75] or regularization [33]. As pointed in our paper, in that respect our LayerScale approach is more related to Rezero [2] and Skip-init [16], Fixup [75], and T-Fixup [34].

7 Conclusion

In this paper, we have shown how train deeper transformer-based image classification neural networks when training on Imagenet only. We have also introduced the simple yet effective CaiT architecture designed in the spirit of encoder/decoder architectures. Our work further demonstrates that transformer models offer a competitive alternative to the best convolutional neural networks when considering trade-offs between accuracy and complexity.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 2
- [2] Thomas C. Bachlechner, Bodhisattwa Prasad Majumder, H. H. Mao, G. Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth. *arXiv preprint arXiv:2003.04887*, 2020. 2, 3, 4, 10, 22
- [3] Irwan Bello. Lambdanetworks: Modeling long-range interactions without attention. In *International Conference on Learning Representations*, 2021. 22
- [4] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *Conference on Computer Vision and Pattern Recognition*, 2019. 22
- [5] Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain: a unified image embedding for classes and instances. *arXiv preprint arXiv:1902.05509*, 2019. 12
- [6] 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. 2, 13, 14
- [7] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. *arXiv preprint arXiv:2101.08692*, 2021. 2
- [8] A. Brock, Soham De, S. L. Smith, and K. Simonyan. High-performance large-scale image recognition without normalization. *arXiv preprint arXiv:2102.06171*, 2021. 2, 14, 15, 22
- [9] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 22
- [10] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *Conference on Computer Vision and Pattern Recognition*, 2005. 22
- [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, 2020. 22
- [12] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Conference on Computer Vision and Pattern Recognition*, 2020. 22
- [13] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019. 2, 4, 22
- [14] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical automated data augmentation with a reduced search space. *arXiv preprint arXiv:1909.13719*, 2019. 14
- [15] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Conference on Computer Vision and Pattern Recognition*, 2017. 21
- [16] Soham De and Samuel L Smith. Batch normalization biases residual blocks towards the identity function in deep networks. *arXiv e-prints*, pages arXiv–2002, 2020. 2, 3, 4, 22
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 8

- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 22
- [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua 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. 2, 3, 8, 13, 14, 15, 22
- [20] Alaaeldin El-Nouby, Natalia Neverova, Ivan Laptev, and Hervé Jégou. Training vision transformers for image retrieval. *arXiv preprint arXiv:2102.05644*, 2021. 11
- [21] Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. *arXiv preprint arXiv:1909.11556*, 2019. ICLR 2020. 8
- [22] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Rémi Gribonval, Hervé Jégou, and Armand Joulin. Training with quantization noise for extreme model compression. *arXiv preprint arXiv:2004.07320*, 2020. 8
- [23] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. 1
- [24] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 22
- [25] Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK Li, and Richard Socher. Non-autoregressive neural machine translation. *arXiv preprint arXiv:1711.02281*, 2017. 22
- [26] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *arXiv preprint arXiv:2103.00112*, 2021. 14
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, June 2016. 1, 21, 22
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027*, 2016. 1, 2
- [29] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoeffler, and Daniel Soudry. Augment your batch: Improving generalization through instance repetition. In *Conference on Computer Vision and Pattern Recognition*, 2020. 12
- [30] Grant Van Horn, Oisin Mac Aodha, Yang Song, Alexander Shepard, Hartwig Adam, Pietro Perona, and Serge J. Belongie. The inaturalist challenge 2018 dataset. *arXiv preprint arXiv:1707.06642*, 2018. 16
- [31] Grant Van Horn, Oisin Mac Aodha, Yang Song, Alexander Shepard, Hartwig Adam, Pietro Perona, and Serge J. Belongie. The inaturalist challenge 2019 dataset. *arXiv preprint arXiv:1707.06642*, 2019. 16
- [32] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. *arXiv preprint arXiv:1709.01507*, 2017. 22
- [33] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. Deep networks with stochastic depth. In *European Conference on Computer Vision*, 2016. 8, 22
- [34] Xiao Shi Huang, Felipe Perez, Jimmy Ba, and Maksims Volkovs. Improving transformer optimization through better initialization. In *International Conference on Machine Learning*, pages 4475–4483. PMLR, 2020. 2, 4, 22
- [35] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 2015. 4

- [36] Shigeki Karita, Nanxin Chen, Tomoki Hayashi, et al. A comparative study on transformer vs rnn in speech applications. *arXiv preprint arXiv:1909.06317*, 2019. 22
- [37] Diederik P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *arXiv preprint arXiv:1807.03039*, 2018. 4
- [38] 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. 16
- [39] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, CIFAR, 2009. 16
- [40] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 21
- [41] Jason Lee, Elman Mansimov, and Kyunghyun Cho. Deterministic non-autoregressive neural sequence modeling by iterative refinement. *arXiv preprint arXiv:1802.06901*, 2018. 22
- [42] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. Selective kernel networks. *Conference on Computer Vision and Pattern Recognition*, 2019. 21
- [43] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 22
- [44] I. Loshchilov and F. Hutter. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 2017. 12
- [45] Christoph Lüscher, Eugen Beck, Kazuki Irie, et al. Rwth asr systems for librispeech: Hybrid vs attention. *Interspeech 2019*, Sep 2019. 22
- [46] 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. 16
- [47] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, pages 4055–4064. PMLR, 2018. 22
- [48] H. Pham, Qizhe Xie, Zihang Dai, and Quoc V. Le. Meta pseudo labels. *arXiv preprint arXiv:2003.10580*, 2020. 22
- [49] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 22
- [50] 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. 1, 13, 14
- [51] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, I. Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. In *Neurips*, 2019. 22
- [52] B. Recht, Rebecca Roelofs, L. Schmidt, and V. Shankar. Do imagenet classifiers generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019. 2, 13, 14
- [53] A. Romero, Nicolas Ballas, S. Kahou, Antoine Chassang, C. Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2015. 22
- [54] 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. 3, 8, 13, 14, 16

- [55] Noam Shazeer, Zhenzhong Lan, Youlong Cheng, N. Ding, and L. Hou. Talking-heads attention. *arXiv preprint arXiv:2003.02436*, 2020. 12
- [56] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 22
- [57] A. Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, P. Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. *arXiv preprint arXiv:2101.11605*, 2021. 22
- [58] R. Srivastava, Klaus Greff, and J. Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387*, 2015. 22
- [59] R. Srivastava, Klaus Greff, and J. Schmidhuber. Training very deep networks. In *NIPS*, 2015. 22
- [60] Gabriel Synnaeve, Qiantong Xu, Jacob Kahn, Tatiana Likhomanenko, Edouard Grave, Vineel Pratap, Anuroop Sriram, Vitaliy Liptchinsky, and Ronan Collobert. End-to-end asr: from supervised to semi-supervised learning with modern architectures. *arXiv preprint arXiv:1911.08460*, 2019. 22
- [61] C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Conference on Computer Vision and Pattern Recognition*, 2015. 22
- [62] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019. 1, 14, 15, 21
- [63] 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. 3, 7, 8, 12, 13, 14, 15, 16, 17, 18, 22, III
- [64] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Herve Jegou. Fixing the train-test resolution discrepancy. *Neurips*, 2019. 13, 21
- [65] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy: Fixefficientnet. *arXiv preprint arXiv:2003.08237*, 2020. 14
- [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 2, 3, 22
- [67] X. Wang, Ross B. Girshick, A. Gupta, and Kaiming He. Non-local neural networks. *Conference on Computer Vision and Pattern Recognition*, 2018. 22
- [68] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. 7, 8, 12, 14, 17, 21
- [69] Bichen Wu, Chenfeng Xu, Xiaoliang Dai, Alvin Wan, Peizhao Zhang, Masayoshi Tomizuka, Kurt Keutzer, and Peter Vajda. Visual transformers: Token-based image representation and processing for computer vision. *arXiv preprint arXiv:2006.03677*, 2020. 22
- [70] L. Xiao, Y. Bahri, Jascha Sohl-Dickstein, S. Schoenholz, and Jeffrey Pennington. Dynamical isometry and a mean field theory of cnns: How to train 10, 000-layer vanilla convolutional neural networks. *arXiv preprint arXiv:1806.05393*, 2018. 22
- [71] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, A. Yuille, and Quoc V. Le. Adversarial examples improve image recognition. *Conference on Computer Vision and Pattern Recognition*, 2020. 14
- [72] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. *Conference on Computer Vision and Pattern Recognition*, 2017. 1

- [73] Qiantong Xu, Alexei Baevski, Tatiana Likhomanenko, Paden Tomasello, Alexis Conneau, Ronan Collobert, Gabriel Synnaeve, and Michael Auli. Self-training and pre-training are complementary for speech recognition. *arXiv preprint arXiv:2010.11430*, 2020. [22](#)
- [74] L. Yuan, Y. Chen, Tao Wang, Weihao Yu, Yujun Shi, F. Tay, Jiashi Feng, and S. Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [14](#), [22](#)
- [75] Hongyi Zhang, Yann Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv preprint arXiv:1901.09321*, 2019. [2](#), [3](#), [4](#), [10](#), [22](#)
- [76] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Zhi-Li Zhang, Haibin Lin, Yu e Sun, Tong He, Jonas Mueller, R. Manmatha, M. Li, and Alex Smola. Resnest: Split-attention networks. *arXiv preprint arXiv:2004.08955*, 2020. [22](#)
- [77] Yu Zhang, James Qin, Daniel S Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V Le, and Yonghui Wu. Pushing the limits of semi-supervised learning for automatic speech recognition. *arXiv preprint arXiv:2010.10504*, 2020. [22](#)
- [78] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2020. [22](#)

Going deeper with Image Transformers

Appendix

In this supplemental material, we provide variations on the architecture presented in our main paper. These experiments have guided some of our choices for the design of class-attention layers and LayerScale.

A Variations on LayerScale init

For the sake of simplicity and to avoid overfitting per model, we have chosen to do a constant initialization with small values depending on the model depth. In order to give additional insight on the importance of this initialization we compare in Table A.1 other possible choices.

LayerScale with 0 init. We initialize all coefficients of LayerScale to 0. This resembles Rezero, but in this case we have distinct learnable parameters for each channel. We make two observations. First, this choice, which also starts with residual branches that output 0 the beginning of the training, gives a clear boost compared to the block-wise scaling done by our adapted ReZero. This confirms the advantage of introducing a learnable parameter per channel and not only per residual layer. Second, LayerScale is better: it is best to initialize to a small ε different from zero.

Random init. We have tested a version in which we try a different initial weight per channel, but with the same average contribution of each residual block as in LayerScale. For this purpose we initialize the channel-scaling values with the Uniform law ($\mathcal{U}[0, 2\varepsilon]$). This simple choice ensures that the expectation of the scaling factor is equal to the value of the classical initialization of LayerScale. This choice is overall comparable to the initialization to 0 of the diagonal, and inferior to LayerScale.

Re-training. LayerScale makes it possible to get increased performance by training deeper models. At the end of training we obtain a specific set of scaling factors for each layer. Inspired by the lottery ticket hypothesis [23], one question that arises is whether what matters is to have the right scaling factors, or to include these learnable weights in the optimization procedure. In other terms, what happens if we re-train the network with the scaling factors obtained by a previous training?

In this experiment below, we try to empirically answer that question. We compare the performance (top-1 validation accuracy, %) on ImageNet-1k with DeiT-S architectures of different depths. Everything being identical otherwise,

Table A.1: Performance when increasing the depth. We compare different strategies and report the top-1 accuracy (%) on ImageNet-1k for the DeiT training (Baseline) with and without adapting the stochastic depth rate d_r (uniform drop-rate), and a modified version of Rezero with LayerNorm and warmup. We compare different initialisation of the diagonal matrix for LayerScale. We also report results with 0 initialization, Uniform initialisation and small constant initialisation. Except for the baseline $d_r = 0.1$, we have adapted the stochastic depth rate d_r .

depth	baseline $d_r = 0.1$	baseline $[d_r]$	ReZero $\alpha = 0$	LayerScale $[\varepsilon]$		
				$\lambda_i = 0$	$\lambda_i = \mathcal{U}[0, 2\varepsilon]$	$\lambda_i = \varepsilon$
12	79.9	79.9 [0.05]	78.3	79.7	80.2 [0.1]	80.5 [0.1]
18	80.1	80.7 [0.10]	80.1	81.5	80.8 [0.1]	81.7 [0.1]
24	78.9†	81.0 [0.20]	80.8	82.1	82.1 [10^{-5}]	82.4 [10^{-5}]
36	78.9†	81.9 [0.25]	81.6	82.7	82.6 [10^{-6}]	82.9 [10^{-6}]

in the first experiment we use LayerScale, i.e. we have learnable weights initialized at a small value ε . In the control experiment we use fixed scaling factors initialised at values obtained by the LayerScale training.

Depth →	12	18	24	36
LayerScale	80.5	81.7	82.4	82.9
Re-trained with fixed weights	80.6	81.5	81.2	81.6

We can see that the control training with fixed weights also converges without suffering from instabilities with the deepest architectures. Nevertheless, the results are lower than those obtained with the learnable weighting factors. This suggests that the evolution of the parameters during training has a beneficial effect on the deepest models.

B Design of the class-attention stage

In this subsection we report some results obtained when considering alternative choices for the class-attention stage.

Not including class embedding in keys of class-attention. In our approach we chose to insert the class embedding in the class-attention: By defining

$$z = [x_{\text{class}}, x_{\text{patches}}], \quad (\text{B.1})$$

we include x_{class} in the keys and therefore the class-attention includes attention on the class embedding itself in Eqn. 6 and Eqn. 7. This is not a requirement as we could simply use a pure cross-attention between the class embedding and the set of frozen patches.

If we do not include the class token in the keys of the class-attention layers, i.e., if we define $z = x_{\text{patches}}$, we reach 83.31% (top-1 acc. on ImageNet1k-val)

Table B.1: CaiT models with and without distillation token. All these models are trained with the same setting during 400 epochs.

Model	Distillation token	
	✗	✓
XXS-24T	78.4	78.5
M-24T	84.8	84.7

with CaiT-S-36, versus 83.44% for the choice adopted in our main paper. This difference of +0.13% is likely not significant, therefore either choice is reasonable. In order to be more consistent with the self-attention layer SA, in the sense that each query has its key counterpart, we have kept the class embedding in the keys of the CA layers as stated in our paper.

Remove LayerScale in Class-Attention. If we remove LayerScale in the Class-Attention blocks in the CaiT-S-36 model, we obtain a top-1 accuracy of 83.36% on ImageNet1k-val, versus 83.44% with LayerScale. The difference of +0.08% is not significant enough to conclude on a clear advantage. For the sake of consistency we have used LayerScale after all residual blocks of the network.

Distillation with class-attention In the main paper we report results with the hard distillation proposed by Touvron *et al.* [63], which in essence replaces the label by the average of the label and the prediction of the teacher output. This is the choice we adopted in our main paper, since it provides better performance than traditional distillation.

The DeiT authors also show the advantage of considering an additional “distillation token”. In their case, employed with the ViT/DeiT architecture, this choice improves the performance compared to hard distillation. Noticeably it accelerates convergence.

In Table B.1 we report the results obtained when inserting a distillation token at the same layer as the class token, i.e., on input of the class-attention stage. In our case we do not observe an advantage of this choice over hard distillation when using class-attention layers. Therefore in our paper we have only considered the hard distillation also employed by Touvron *et al.* [63].