Transformed CNNs: recasting pre-trained convolutional layers with self-attention

Stéphane d'Ascoli^{*1,2}, Levent Sagun², Giulio Biroli¹ and Ari Morcos²

¹Laboratoire de Physique de l'Ecole Normale Supérieure, Université PSL, CNRS, Sorbonne Université, Université Paris-Diderot, Sorbonne Paris Cité, Paris, France ²Facebook AI Research, Paris, France

Abstract

Vision Transformers (ViT) have recently emerged as a powerful alternative to convolutional networks (CNNs). Although hybrid models attempt to bridge the gap between these two architectures, the self-attention layers they rely on induce a strong computational bottleneck, especially at large spatial resolutions. In this work, we explore the idea of reducing the time spent training these layers by initializing them as convolutional layers. This enables us to transition smoothly from any pre-trained CNN to its functionally identical hybrid model, called Transformed CNN (T-CNN). With only 50 epochs of fine-tuning, the resulting T-CNNs demonstrate significant performance gains over the CNN (+2.2% top-1 on ImageNet-1k for a ResNet50-RS) as well as substantially improved robustness (+11% top-1 on ImageNet-C). We analyze the representations learnt by the T-CNN, providing deeper insights into the fruitful interplay between convolutions and self-attention. Finally, we experiment initializing the T-CNN from a partially trained CNN, and find that it reaches better performance than the corresponding hybrid model trained from scratch, while reducing training time.

Introduction

Since the success of AlexNet in 2012 [1], the field of Computer Vision has been dominated by Convolutional Neural Networks (CNNs) [2, 3]. Their local receptive fields give them a strong inductive bias to exploit the spatial structure of natural images [4–6], while allowing them to scale to large resolutions seamlessly. Yet, this inductive bias limits their ability to capture long-range interactions.

In this regard, self-attention (SA) layers, originally introduced in language models [7–9], have gained interest as a building block for vision Ramachandran *et al.* [10] and Zhao *et al.* [11]. Recently, they gave rise to a plethora of Vision Transformer (ViT) models, able to compete with state-of-the-art CNNs in various tasks Dosovitskiy *et al.* [12], Touvron *et al.* [13], Wu *et al.* [14], Touvron *et al.* [15], Liu *et al.* [16] and Heo *et al.* [17] while demonstrating better robustness [18, 19]. However, capturing long-range dependencies necessarily comes at the cost of quadratic complexity in input size, a computational burden which many recent directions have tried to alleviate [20–23]. Additionally, ViTs are generally harder to train [24, 25], and require vast amounts of pre-training [12] or distillation from a convolutional teacher [26–28] to match the performance of CNNs.

Faced with the dilemma between efficient CNNs and powerful ViTs, several approaches have aimed to bridge the gap between these architectures. On one side, hybrid models append SA layers onto convolutional backbones [28–32], and have already fueled successful results in a variety of tasks [33–37]. Conversely, a line of research has studied the benefit of introducing convolutional biases in Transformer architectures to ease learning [38–40]. Despite these interesting compromises, modelling long-range dependencies at low computational cost remains a challenge for practitioners.

^{*}stephane.dascoli@ens.fr



Figure 1: **Transformed ResNets strike a strong accuracy-robustness balance.** Our models (red) significantly outperform the original ResNet-RS models (dark blue) they were initialized from when evaluated on ImageNet-1k. On various robustness benchmarks (ImageNet-C, A and R, from left to right), they narrow or close the gap with Transformer architectures.

Contributions At a time when pre-training on vast datasets has become common practice, we ask the following question: does one need to train the SA layers during the whole learning process? Could one instead learn cheap components such as convolutions first, leaving the SA layers to be learnt at the end? In this paper, we take a step in this direction by presenting a method to fully reparameterize a pre-trained convolutional layer as a *Gated Positional Self-Attention* (GPSA) layer [38]. The latter is initialized to reproduce the mapping of the convolutional layer, but is then encouraged to learn more general mappings which are not accessible to the CNN by adjusting positional gating parameters.

We leverage this method to reparametrize pre-trained CNNs as functionally equivalent hybrid models. After only 50 epochs of fine-tuning, the resulting Transformed CNNs (T-CNNs) boast significant performance and robustness improvements as shown in Fig. 1, demonstrating the practical relevance of our method. We analyze the inner workings of the T-CNNs, showing how they learn more robust representations by combining convolutional heads and SA heads in a complementary way. Finally, we investigate how performance gains depend on the reparametrization epoch. Results suggest that reparametrizing at intermediate times is optimal in terms of speed-performance trade-offs.

Related work Our work mainly builds on two pillars. First, the idea that SA layers can express any convolution, introduced by Cordonnier *et al.* [41]. This idea was recently leveraged in d'Ascoli *et al.* [38], which initialize the SA layers of the ViT as *random* convolutions and observe performance gains compared to the standard initialization, especially in the low-data regime where inductive biases are most useful. Our approach is a natural follow-up of this idea: what happens if the SA layers are instead initialized as *trained* convolutions?

Second, we exploit the following learning paradigm: train a simple and fast model, then reparameterize it as a more complex model for the final stages of learning. This approach was studied from a scientific point of view in d'Ascoli *et al.* [42], which shows that reparameterizing a CNN as a fullyconnected network (FCN) halfway through training can lead the FCN to outperform the CNN. Yet, the practical relevance of this method is limited by the vast increase in number of parameters required by the FCN to functionally represent the CNN. In contrast, our reparameterization hardly increases the parameter count of the CNN, making it easily applicable to any state-of-the-art CNN. Note that these reparameterization methods can be viewed an informed version of dynamic architecture growing algorithms such as AutoGrow [43].

In the context of hybrid models, various works have studied the performance gains obtained by introducing MHSA layers in ResNets with minimal architectural changes [28, 31, 32]. However, the MHSA layers used in these works are initialized randomly and need to be trained from scratch. Our approach is different, as it makes use of GPSA layers, which can be initialized to represent the same function as the convolutional layer it replaces. We emphasize that the novelty in our work is not in the architectures used, but in the unusual way they are blended together.

1 Background

Multi-head self-attention The SA mechanism is based on a trainable associative memory with (key, query) vector pairs. To extract the semantic interpendencies between the L elements of a sequence $X \in \mathbb{R}^{L \times D_{in}}$, a sequence of "query" embeddings $Q = W_{qry}X \in \mathbb{R}^{L \times D_h}$ is matched against another sequence of "key" embeddings $K = W_{key}X \in \mathbb{R}^{L \times D_h}$ using inner products. The result is an attention matrix whose entry (ij) quantifies how semantically relevant Q_i is to K_j :

$$\boldsymbol{A} = \operatorname{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^{\top}}{\sqrt{D_{h}}}\right) \in \mathbb{R}^{L \times L}.$$
 (1)

Multi-head SA layers use several SA heads in parallel to allow the learning of different kinds of dependencies:

$$MSA(\boldsymbol{X}) := \sum_{h=1}^{N_h} [SA_h(\boldsymbol{X})] \boldsymbol{W}_{out}^h, \qquad SA_h(\boldsymbol{X}) := \boldsymbol{A}^h \boldsymbol{X} \boldsymbol{W}_{val}^h,$$
(2)

where $\boldsymbol{W}_{\mathrm{val}}^h \in R^{D_{in} \times D_v}$ and $\boldsymbol{W}_{\mathrm{out}}^h \in R^{D_v \times D_{out}}$ are two learnable projections.

To incorporate positional information, ViTs usually add absolute position information to the input at embedding time, before propagating it through the SA layers. Another possibility is to replace the vanilla SA with positional SA (PSA), including a position-dependent term in the softmax [10, 44]. Although there are several way to parametrize the positional attention, we use encodings r_{ij} of the relative position of pixels *i* and *j* as in [38, 41]:

$$\boldsymbol{A}_{ij}^{h} := \operatorname{softmax} \left(\boldsymbol{Q}_{i}^{h} \boldsymbol{K}_{j}^{h\top} + \boldsymbol{v}_{pos}^{h\top} \boldsymbol{r}_{ij} \right).$$
(3)

Each attention head learns an embedding $v_{pos}^h \in \mathbb{R}^{D_{pos}}$, and the relative positional encodings $r_{ij} \in \mathbb{R}^{D_{pos}}$ only depend on the distance between pixels i and j, denoted denoted as a two-dimensional vector δ_{ij} .

Self-attention as a generalized convolution Cordonnier *et al.* [41] shows that a multi-head PSA layer (Eq. 3) with N_h heads and dimension $D_{pos} \ge 3$ can express any convolutional layer of filter size $\sqrt{N_h}$, with D_{in} input channels and min (D_v, D_{out}) output channels, by setting the following:

In the above, the *center of attention* $\Delta^h \in \mathbb{R}^2$ is the position to which head h pays most attention to, relative to the query pixel, whereas the *locality strength* $\alpha^h > 0$ determines how focused the attention is around its center Δ^h . When α^h is large, the attention is focused only on the pixel located at Δ^h ; when α^h is small, the attention is spread out into a larger area. Thus, the PSA layer can achieve a convolutional attention map by setting the centers of attention Δ^h to each of the possible positional offsets of a $\sqrt{N_h} \times \sqrt{N_h}$ convolutional kernel, and sending the locality strengths α^h to some large value.

2 Approach

In this section, we introduce our method for mapping a convolutional layer to a functionally equivalent PSA layer with minimal increase in parameter count. To do this, we leverage the GPSA layers introduced in d'Ascoli *et al.* [38].

Loading the filters We want each head h of the PSA layer to functionally mimic the pixel h of a convolutional filter $W_{\text{filter}} \in \mathbb{R}^{N_h \times D_{in} \times D_{out}}$, where we typically have $D_{out} \ge D_{in}$. Rewriting the action of the MHSA operator in a more explicit form, we have

$$MHSA(\boldsymbol{X}) = \sum_{h=1}^{N_h} \boldsymbol{A}^h \boldsymbol{X} \underbrace{\boldsymbol{W}_{val}^h \boldsymbol{W}_{out}^h}_{\boldsymbol{W}^h \in \mathbb{R}^{D_{in} \times D_{out}}}$$
(5)

In the convolutional configuration of Eq. 4, $A^h X$ selects pixel h of X. Hence, we need to set $W^h = W^h_{\text{filter}}$. However, as a product of matrices, the rank of W_h is bottlenecked by D_v . To avoid this being a limitation, we need $D_v \ge D_{in}$ (since $D_{out} \ge D_{in}$). To achieve this with a minimal number of parameters, we choose $D_v = D_{in}$, and simply set the following initialization:

$$\boldsymbol{W}_{\mathrm{val}}^{h} = \boldsymbol{I}, \qquad \boldsymbol{W}_{\mathrm{out}}^{h} = \boldsymbol{W}_{\mathrm{filter}}^{h}.$$
 (6)

Note that this differs from the usual choice made in SA layers, where $D_v = \lfloor D_{in}/N_h \rfloor$. However, to keep the parameter count the same, we share the same W_{val}^h across different heads h, since it plays a symmetric role at initialization.

Note that this reparameterization introduces three additional matrices compared to the convolutional filter: W_{qry} , W_{key} , W_{val} , each containing $D_{in} \times D_{in}$ parameters. However, since the convolutional filter contains $N_h \times D_{in} \times D_{out}$ parameters, where we typically have $N_h = 9$ and $D_{out} \in \{D_{in}, 2D_{in}\}$, these additional matrices are much smaller than the filters and hardly increase the parameter count. This can be seen from the model sizes in Tab. 2.

Gated Positional self-attention Recent work [38] has highlighted an issue with standard PSA: the fact that the content and positional terms in Eq. 3 are potentially of very different magnitudes, in which case the softmax ignores the smallest of the two. This can typically lead the PSA to adopt a greedy attitude: choosing the form of attention (content or positional) which is easiest at a given time then sticking to it.

To avoid this, the ConViT d'Ascoli *et al.* [38] uses GPSA layers which sum the content and positional terms *after* the softmax, with their relative importances governed by a learnable *gating* parameter λ_h (one for each attention head). In GPSA layers, the attention is parametrized as follows:

$$\boldsymbol{A}_{ij}^{h} := (1 - \sigma(\lambda_{h})) \operatorname{softmax} \left(\boldsymbol{Q}_{i}^{h} \boldsymbol{K}_{j}^{h\top} \right) + \sigma(\lambda_{h}) \operatorname{softmax} \left(\boldsymbol{v}_{pos}^{h\top} \boldsymbol{r}_{ij} \right),$$
(7)

where $\sigma : x \mapsto 1/(1+e^{-x})$ is the sigmoid function. In the positional part, the encodings r_{ij} are fixed rather than learnt (see Eq. 4), which makes changing input resolution straightforward (see SM. B) and leaves only 3 learnable parameters per head: Δ_1, Δ_2 and α^1 .

How convolutional should the initialization be? The convolutional initialization of GPSA layers involves two parameters, determining how strictly convolutional the behavior is: the initial value of the *locality strength* α , which determines how focused each attention head is on its dedicated pixel, and the initial value of the *gating parameters* λ , which determines the importance of the positional information versus content. If $\lambda_h \gg 0$ and $\alpha \gg 1$, the T-CNN will perfectly reproduce the input-output function of the CNN, but may stay stuck in the convolutional configuration. Conversely, if $\lambda_h \ll 0$ and $\alpha \ll 1$, the T-CNN will poorly reproduce the input-output function of the CNN. Hence, we choose $\alpha = 1$ and $\lambda = 1$ to lie in between these two extremes. This puts the T-CNN "on the verge of locality", enabling it to escape locality effectively throughout training.

¹Since α represents the temperature of the softmax, its value must stay positive at all times. To ensure this, we instead learn a rectified parameter $\tilde{\alpha}$ using the softplus function: $\alpha = \frac{1}{\beta} \log(1 + e^{-\beta \tilde{\alpha}})$, with $\beta = 5$.



Figure 2: **T-CNNs present better speed-accuracy trade-offs than the CNNs they stem from.** Total training time (original training + finetuning) is normalized by the total training time of the ResNet50-RS. Inference throughput is the number of images processed per second on a V100 GPU at batch size 32.

Architectural details To make our setup as canonical as possible, we focus on ResNet architectures [45], which contain 5 stages, with spatial resolution halfed and number of channels doubled at each stage. Our method involves reparameterizing 3×3 convolutions as GPSA layers with 9 attention heads. However, global SA is too costly in the first layers, where the spatial resolution is large. We therefore only reparameterize the last stage of the architecture, while replacing the first stride-2 convolution by a stride-1 convolution, exactly as in [32]. We also add explicit padding layers to account for the padding of the original convolutions.

3 Performance of the Transformed CNNs

In this section, we apply our reparametrization to state-of-the-art CNNs, then fine-tune the resulting T-CNNs to learn better representations. This method allows to fully disentangle the training of the SA layers from that of the convolutional backbone, which is of practical interest for two reasons. First, it minimizes the time spent training the SA layers, which typically have a slower throughput. Second, it separates the algorithmic choices of the CNN backbone from those of the SA layers, which are typically different; for example, CNNs are typically trained with SGD whereas SA layers perform much better with adaptive optimizers such as Adam [24], an incompatibility which may limit the performance of usual hybrid models.

Training details To minimize computational cost, we restrict the fine-tuning to 50 epochs². Following [24], we use the AdamW optimizer, with a batch size of 1024^3 . The learning rate is warmed up to 10^{-4} then annealed using a cosine decay. To encourage the T-CNN to escape the convolutional configuration and learn content-based attention, we use a larger learning rate of 0.1 for the gating parameters of Eq. 7 (one could equivalently decrease the temperature of the sigmoid function).

We use the same data augmentation scheme as the DeiT [13], as well as rather large stochastic depth coefficients d_r reported in Tab. 1. Hoping that our method could be used as an alternative to the commonly used practice of fine-tuning models at higher resolution, we also increase the resolution during fine-tuning [46]. In this setting, a ResNet50 requires only 6 hours of fine-tuning on 16 V100 GPUs, compared to 33 hours for the original training. For our largest model (ResNet350-RS), the fine-tuning lasts 50 hours.

²We study how performance depends on the number of fine-tuning epochs in SM. C.

³Confirming the results of [24], we obtained worse results with SGD.

| | Training | | | | Fine-tuning | | | | | |
|--------------|----------|-------|----------|-------|-------------|-------|------------|-------|---------|-------|
| Backbone | | | | | | | Without SA | | With SA | |
| | Res. | d_r | TTT | Top-1 | Res. | d_r | TTT | Top-1 | TTT | Top-1 |
| ResNet50-RS | 160 | 0.0 | 1 (ref.) | 78.8 | 224 | 0.1 | 1.16 | 80.4 | 1.30 | 81.0 |
| ResNet101-RS | 192 | 0.0 | 1.39 | 80.3 | 224 | 0.1 | 1.65 | 81.9 | 1.79 | 82.4 |
| ResNet152-RS | 256 | 0.0 | 3.08 | 81.2 | 320 | 0.2 | 3.75 | 83.4 | 4.13 | 83.7 |
| ResNet200-RS | 256 | 0.1 | 4.15 | 82.8 | 320 | 0.2 | 5.04 | 83.7 | 5.42 | 84.0 |
| ResNet270-RS | 256 | 0.1 | 6.19 | 83.8 | 320 | 0.2 | 7.49 | 83.9 | 7.98 | 84.3 |
| ResNet350-RS | 288 | 0.1 | 10.49 | 84.0 | 320 | 0.2 | 12.17 | 84.1 | 12.69 | 84.5 |

Table 1: Statistics of the models considered, trained from scratch on ImageNet. Top-1 accuracy is measured on ImageNet-1k validation set. "TTT" stands for total training time (including fine-tuning), normalized by the total training time of the ResNet50-RS. d_r is the stochastic depth coefficient used for the various models.



Figure 3: **Robustness is most improved for strong and blurry corruption categories.** We report the relative improvement between the top-1 accuracy of the T-ResNet50-RS and that of the ResNet50-RS on ImageNet-C, averaging over the different corruption categories (left) and corruption severities (right).

Performance gains We applied our method to pre-trained ResNet-RS [47] models, using the weights provided by the timm package [48]. These models are derived from the original ResNet [45], but use improved architectural features and training strategies, enabling them to reach better speed-accuracy trade-offs than EfficientNets. Results are presented in Tab. 1, where we also report the baseline improvement of fine-tuning in the same setting but without SA. In all cases, our fine-tuning improves top-1 accuracy, with a significant gap over the baseline. To demonstrate the wide applicability of our method, we report similar improvements for ResNet-D architectures in SM. D.

Despite the extra fine-tuning epochs and their slower throughput, the resulting T-CNNs match the performance of the original CNNs at equal throughput, while significantly outperforming them at equal total training time, as shown in the Pareto curves of Fig. $2(a)^4$. However, the major benefit of the reparametrization is in terms of robustness, as shown in Fig. 2(b) and explained below.

Robustness gains Recent work [18, 19] has shown that Transformer-based architectures are more robust to input perturbations than convolutional architectures. We therefore investigate whether our fine-

⁴We estimated the training times of the original ResNet-RS models based on their throughput, for the same hardware as used for the T-ResNet-RS.

| Model | Res. | Params | Speed | Flops | ImNet-1k | ImNet-C | ImNet-A | ImNet-R | | |
|----------------------|------|--------|-------|-------|----------|---------|---------|---------|--|--|
| Transformers | | | | | | | | | | |
| ViT-B/16 | 224 | 86 M | 182 | 16.9 | 77.9 | 52.2 | 7.0 | 21.9 | | |
| ViT-L/16 | 224 | 307 M | 55 | 59.7 | 76.5 | 49.3 | 6.1 | 17.9 | | |
| DeiT-S | 224 | 22 M | 544 | 4.6 | 79.9 | 55.4 | 18.9 | 31.0 | | |
| DeiT-B | 224 | 87 M | 182 | 17.6 | 82.0 | 60.7 | 27.4 | 34.6 | | |
| ConViT-S | 224 | 28 M | 296 | 5.4 | 81.5 | 59.5 | 24.5 | 34.0 | | |
| ConViT-B | 224 | 87 M | 139 | 17.7 | 82.4 | 61.9 | 29.0 | 36.9 | | |
| CNNs | | | | | | | | | | |
| ResNet50 | 224 | 25 M | 736 | 4.1 | 76.8 | 46.1 | 4.2 | 21.5 | | |
| ResNet101 | 224 | 45 M | 435 | 7.85 | 78.0 | 50.2 | 6.3 | 23.0 | | |
| ResNet101x3 | 224 | 207 M | 62 | 69.6 | 80.3 | 53.4 | 9.1 | 24.5 | | |
| ResNet152x4 | 224 | 965 M | 18 | 183.1 | 80.4 | 54.5 | 11.6 | 25.8 | | |
| ResNet50-RS | 160 | 36 M | 938 | 4.6 | 78.8 | 36.8 | 5.7 | 39.1 | | |
| ResNet101-RS | 192 | 64 M | 674 | 12.1 | 80.3 | 44.1 | 11.8 | 44.8 | | |
| ResNet152-RS | 256 | 87 M | 304 | 31.2 | 81.2 | 49.9 | 23.4 | 45.9 | | |
| ResNet200-RS | 256 | 93 M | 225 | 40.4 | 82.8 | 49.3 | 25.4 | 48.1 | | |
| ResNet270-RS | 256 | 130 M | 152 | 54.2 | 83.8 | 53.6 | 26.6 | 48.7 | | |
| ResNet350-RS | 288 | 164 M | 89 | 87.5 | 84.0 | 53.9 | 34.9 | 49.7 | | |
| Our transformed CNNs | | | | | | | | | | |
| T-ResNet50-RS | 224 | 38 M | 447 | 17.6 | 81.0 | 48.0 | 18.7 | 42.9 | | |
| T-ResNet101-RS | 224 | 66 M | 334 | 25.1 | 82.4 | 52.9 | 27.7 | 47.8 | | |
| T-ResNet152-RS | 320 | 89 M | 128 | 65.8 | 83.7 | 54.5 | 39.8 | 50.6 | | |
| T-ResNet200-RS | 320 | 96 M | 105 | 80.2 | 84.0 | 57.0 | 41.2 | 51.1 | | |
| T-ResNet270-RS | 320 | 133 M | 75 | 107.2 | 84.3 | 58.6 | 43.7 | 51.4 | | |
| T-ResNet350-RS | 320 | 167 M | 61 | 130.5 | 84.5 | 59.2 | 44.8 | 53.8 | | |

Table 2: Accuracy of our models on various benchmarks. Throughput is the number of images processed per second on a V100 GPU at batch size 32. The ViT and ResNet results are reported in [18]. For ImageNet-C, we keep a resolution of 224 at test time to avoid distorting the corruptions.

tuning procedure brings robustness gains to the original CNNs. To do so, we consider three benchmarks. First, ImageNet-C [49], a dataset containing 15 sets of randomly generated corruptions, grouped into 4 categories: 'noise', 'blur', 'weather', and 'digital'. Each corruption type has five levels of severity, resulting in 75 distinct corruptions. Second, ImageNet-A [50], a dataset containing naturally "adversarial" examples from ImageNet. Finally, we evaluate robustness to distribution shifts with ImageNet-R [51], a dataset with various stylized "renditions" of ImageNet images ranging from paintings to embroidery, which strongly modify the local image statistics.

As shown in Tab. 2 and illustrated in Fig. 1, the T-ResNet-RS substantially outperforms the ResNet-RS on all three benchmarks. For example, our T-ResNet101-RS reaches similar or higher top-1 accuracy than the ResNet200-RS on each task, despite its lower top-1 accuracy on ImageNet-1k. This demonstrates that SA improves robustness more than it improves classification accuracy.

To better understand where the benefits come from, we decompose the improvement of the T-ResNet50-RS over the various corruption severeties and categories of ImageNet-C in Fig. 3. We observe that improvement increases almost linearly with corruption severity. Although performance is higher in all corruption categories, there is a strong variability: the T-CNN shines particularly in tasks where the objects in the image are less sharp due to lack of contrast, bad weather or blurriness. We attribute this to the ability of SA to distinguish shapes in the image, as investigated in Sec 4.

4 Dissecting the Transformed CNNs

In this section, we analyze various observables to understand how the representations of a T-ResNet270-RS evolve from those of the ResNet270-RS throughout training.



Figure 4: The later layers effectively escape the convolutional configuration. A: top-1 accuracy throughout the 50 epochs of fine-tuning of a T-ResNet270-RS. B: size of the receptive field of the various heads h (thin lines), calculated as α_h^{-1} (see Eq. 3). Thick lines represent the average over the heads. C: depicts how much attention the various heads h (thin lines) pay to positional information, through the value of $\sigma(\lambda_h)$ (see Eq. 7). Thick lines represent the average over the heads.

Unlearn to better relearn In Fig. 4A, we display the train and test accuracy throughout training⁵. The dynamics decompose into two distinct phases: accuracy dips down during the learning rate warmup phase (first 5 epochs of training), then increases back up as the learning rate is decayed.

Interestingly, as shown in SM. A, the depth of the dip depends on the learning rate. For too small learning rates, the dip is small, but the test accuracy increases too slowly after the dip; for too large learning rates, the test accuracy increases rapidly after the dip, but the dip is too deep to be compensated for. This suggests that the T-CNN needs to "unlearn" to some extent, a phenomenon reminiscent of the

⁵The train accuracy is lower than the test accuracy due to the heavy data augmentation used during fine-tuning.



(a) Input image

(b) Attention maps

Figure 5: **GPSA layers combine local and global attention in a complementary way.** We depicted the attention maps of the four GPSA layers of the T-ResNet270-RS, obtained by feeding the image on the left through the convolutional backbone, then selecting a query pixel in the center of the image (red box). For each head *h*, we indicate the value of the gating parameter $\sigma(\lambda_h)$ in red (see Eq. 7). In each layer, at least one of the heads learns to perform content-based attention ($\sigma(\lambda_h) = 0$).

"catapult" mechanism of Lewkowycz *et al.* [52] which propels models out of sharp minima to land in wider minima.

Escaping the convolutional representation In Fig. 4B, we show the evolution of the "attention span" $1/\alpha_h$ (see Eq. 4), which reflects the size of the receptive field of attention head h. On average (thick lines), this quantity increases in the first three layers, showing that the attention span widens, but variability exists among different attention heads (thin lines): some broaden their receptive field, whereas others contract it.

In Fig. 4C, we show the evolution of the gating parameters λ^h of Eq. 7, which reflect how much attention head *h* pays to position versus content. Interestingly, the first layer stays strongly convolutional on average, as $\mathbb{E}_h \sigma(\lambda_h)$ rapidly becomes close to one (thick blue line). The other layers strongly escape locality, with most attention heads focusing on content information at the end of fine-tuning.

In Fig. 5, we display the attention maps after fine-tuning. A clear divide appears between the "convolutional" attention heads, which remain close to their initialization, and the "content-based" attention heads, which learn more complex dependencies. Notice that the attention head initially focusing on the query pixel (head 5) stays convolutional in all layers. Throughout the layers, the shape of the central object is more and more clearly visible, as observed in [53]. This supports the hypothesis that robustness gains obtained for blurry corruptions (see Fig. 3) are partly due to the ability of the SA layers to isolate objects from the background.

5 When should one start learning the self-attention layers?

Previous sections have demonstrated the benefits of initializing T-CNNs from pre-trained CNNs, a very compelling procedure given the wide availability of pretrained models. But one may ask: how does this compare to training a hybrid model from scratch? More generally, given a computational budget, how long should the SA layers be trained compared to the convolutional backbone?

Transformed CNN versus hybrid models To answer the first question, we consider a ResNet-50 trained on ImageNet for 400 epochs. We use SGD with momentum 0.9 and a batch size of 1024, warming up the learning rate for 5 epochs before a cosine decay. To achieve a strong baseline, we use the same augmentation scheme as in [13] for the DeiT. Results are reported in Tab. 3. In this modern training setting, the vanilla ResNet50 reaches a solid performance of 79.04% on ImageNet, well above the 77% usually reported in litterature.

| Name | t_1 | t_2 | Train time | Top-1 |
|-----------------|-------|-------|------------|-------|
| Vanilla CNN | 400 | 0 | 2.0k mn | 79.04 |
| Vanilla CNN↑320 | 450 | 0 | 2.4k mn | 79.78 |
| T-CNN | 400 | 50 | 2.3k mn | 79.88 |
| T-CNN↑320 | 400 | 50 | 2.7k mn | 80.84 |
| Vanilla hybrid | 0 | 400 | 2.8k mn | 79.95 |
| T-CNN* | 100 | 300 | 2.6k mn | 80.44 |
| T-CNN* | 200 | 200 | 2.4k mn | 80.28 |
| T-CNN* | 300 | 100 | 2.2k mn | 79.28 |

Table 3: The benefit of late reparametrization. We report the top-1 accuracy of a ResNet-50 on ImageNet reparameterized at various times t_1 during training. \uparrow 320 stands for fine-tuning at resolution 320. The models with a \star keep the same optimizer after reparametrization, in contrast with the usual T-CNNs.

The T-CNN obtained by fine-tuning the ResNet for 50 epochs at same resolution obtains a top-1 accuracy of 79.88%, with a 15% increase in training time, and 80.84 as resolution 320, with a 35% increase in training time. In comparison, the hybrid model trained for 400 epochs in the same setting only reaches 79.95%, in spite of a 40% increase in training time. Hence, fine-tuning yields better results than training the hybrid model from scratch.

What is the best time to reparametrize? We now study a scenario between the two extreme cases: what happens if we reparametrize halfway through training? To investigate this question in a systematic way, we train the ResNet50 for t_1 epochs, then reparametrize and resume training for another t_2 epochs, ensuring that $t_1 + t_2 = 400$ in all cases. Hence, $t_1 = 400$, amounts to the vanilla ResNet50, whereas $t_1 = 0$ corresponds to the hybrid model trained from scratch. To study how final performance depends on t_1 in a fair setting, we keep the same optimizer and learning rate after the reparametrization, in contrast with the fine-tuning procedure which uses fresh optimizer.

Results are presented in Tab. 3. Interestingly, the final performance evolves non-monotonically, reaching a maximum of 80.44 for $t_1 = 100$, then decreasing back down as the SA layers have less and less time to learn. This non-monotonicity is remarkably similar to that observed in [42], where reparameterizing a CNN as a FCN in the early stages of training enables the FCN to outperform the CNN. Crucially, this result suggests that reparametrizing during training not only saves time, but also helps the T-CNN find better solutions.

Discussion

In this work, we showed that complex building blocks such as self-attention layers need not be trained from start. Instead, one can save in compute time while gaining in performance and robustness by initializing them from pre-trained convolutional layers. At a time where energy savings and robustness are key stakes, we believe this finding is important.

On the practical side, our fine-tuning method offers an interesting new direction for practitioners. One clear limitation of our method is the prohibitive cost of reparametrizing the early stages of CNNs. This cost could however be alleviated by using linear attention methods [21], an important direction for future work. Note also that while our T-CNNs significantly improve the robustness of CNNs, they do not systematically reach the performance of end-to-end Transformers such as the DeiT (for example on ImageNet-C, see Fig. 1). Bridging this gap is an important next step for hybrid models.

On the theoretical side, our results spark several interesting questions. First, why is it better to reparametrize at intermediate times? One natural hypothesis, which will be explored in future work, is

that SA layers benefit from capturing meaningful dependencies between the features learnt by the CNN, rather than the random correlations which exist at initialization. Second, why are the representations learnt by the SA layers more robust? By inspecting the attention maps and the most improved corruption categories of ImageNet-C, we hypothesized that SA helps isolating objects from the background, but a more thorough analysis is yet to come.

Acknowledgements We thank Matthew Leavitt, Hugo Touvron, Hervé Jégou and Francisco Massa for helpful discussions. SD and GB acknowledge funding from the French government under management of Agence Nationale de la Recherche as part of the "Investissements d'avenir" program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute).

References

- 1. Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. *Communications of the ACM* **60**, 84–90 (2017).
- 2. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**, 2278–2324 (1998).
- 3. LeCun, Y. *et al.* Backpropagation applied to handwritten zip code recognition. *Neural computation* **1**, 541–551 (1989).
- Scherer, D., Müller, A. & Behnke, S. Evaluation of Pooling Operations in Convolutional Architectures for Object Recognition en. in Artificial Neural Networks – ICANN 2010 (eds Diamantaras, K., Duch, W. & Iliadis, L. S.) (Springer, Berlin, Heidelberg, 2010), 92–101.
- Schmidhuber, J. Deep learning in neural networks: An overview. en. Neural Networks 61, 85-117. http://www.sciencedirect.com/science/article/pii/S0893608014002135 (2021) (Jan. 2015).
- 6. Goodfellow, I., Bengio, Y. & Courville, A. Deep Learning (MIT Press, 2016).
- 7. Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- 8. Vaswani, A. et al. Attention is all you need in Advances in neural information processing systems (2017), 5998–6008.
- 9. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- 10. Ramachandran, P. *et al.* Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909* (2019).
- 11. Zhao, H., Jia, J. & Koltun, V. Exploring self-attention for image recognition in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020), 10076–10085.
- 12. Dosovitskiy, A. *et al.* An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv preprint arXiv:2010.11929* (2020).
- 13. Touvron, H. *et al.* Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877* (2020).
- 14. Wu, B. *et al.* Visual Transformers: Token-based Image Representation and Processing for Computer Vision. *arXiv:2006.03677 [cs, eess]*. arXiv: 2006.03677. http://arxiv.org/abs/2006.03677 (2020) (July 2020).
- 15. Touvron, H., Cord, M., Sablayrolles, A., Synnaeve, G. & Jégou, H. Going deeper with Image Transformers. *arXiv preprint arXiv:2103.17239* (2021).
- 16. Liu, Z. *et al.* Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030* (2021).

- 17. Heo, B. *et al.* Rethinking spatial dimensions of vision transformers. *arXiv preprint arXiv:2103.16302* (2021).
- 18. Bhojanapalli, S. *et al.* Understanding robustness of transformers for image classification. *arXiv preprint arXiv:2103.14586* (2021).
- 19. Mao, X. *et al.* Rethinking the Design Principles of Robust Vision Transformer. *arXiv preprint arXiv:2105.07926* (2021).
- 20. Bello, I. Lambdanetworks: Modeling long-range interactions without attention. *arXiv preprint arXiv:2102.08602* (2021).
- 21. Wang, S., Li, B., Khabsa, M., Fang, H. & Ma, H. L. Self-Attention with Linear Complexity. *arXiv* preprint arXiv:2006.04768 (2020).
- 22. Choromanski, K. et al. Rethinking attention with performers. arXiv preprint arXiv:2009.14794 (2020).
- 23. Katharopoulos, A., Vyas, A., Pappas, N. & Fleuret, F. Transformers are rnns: Fast autoregressive transformers with linear attention in International Conference on Machine Learning (2020), 5156–5165.
- 24. Zhang, J. *et al.* Why are Adaptive Methods Good for Attention Models? *arXiv preprint arXiv:1912.03194* (2019).
- 25. Liu, L., Liu, X., Gao, J., Chen, W. & Han, J. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249* (2020).
- 26. Hinton, G., Vinyals, O. & Dean, J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
- 27. Jiang, Z. et al. Token Labeling: Training a 85.4% Top-1 Accuracy Vision Transformer with 56M Parameters on ImageNet 2021. arXiv: 2104.10858 [cs.CV].
- 28. Graham, B. *et al.* LeViT: a Vision Transformer in ConvNet's Clothing for Faster Inference. *arXiv* preprint arXiv:2104.01136 (2021).
- 29. Chen, Y., Kalantidis, Y., Li, J., Yan, S. & Feng, J. A2-Nets: Double Attention Networks. *arXiv preprint arXiv:1810.11579* (2018).
- 30. Bello, I., Zoph, B., Vaswani, A., Shlens, J. & Le, Q. V. Attention augmented convolutional networks in Proceedings of the IEEE International Conference on Computer Vision (2019), 3286–3295.
- 31. Chen, Z. et al. Visformer: The Vision-friendly Transformer 2021. arXiv: 2104.12533 [cs.CV].
- 32. Srinivas, A. *et al.* Bottleneck Transformers for Visual Recognition. *arXiv e-prints*, arXiv:2101.11605. arXiv: 2101.11605 [cs.CV] (Jan. 2021).
- 33. Carion, N. *et al.* End-to-End Object Detection with Transformers. *arXiv preprint arXiv:2005.12872* (2020).
- 34. Hu, H., Gu, J., Zhang, Z., Dai, J. & Wei, Y. Relation networks for object detection in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018), 3588–3597.
- 35. Chen, Y.-C. et al. Uniter: Universal image-text representation learning in European Conference on Computer Vision (2020), 104–120.
- 36. Locatello, F. et al. Object-centric learning with slot attention. arXiv preprint arXiv:2006.15055 (2020).
- 37. Sun, C., Myers, A., Vondrick, C., Murphy, K. & Schmid, C. Videobert: A joint model for video and language representation learning in Proceedings of the IEEE International Conference on Computer Vision (2019), 7464–7473.
- 38. d'Ascoli, S. *et al.* Convit: Improving vision transformers with soft convolutional inductive biases. *arXiv preprint arXiv:2103.10697* (2021).
- 39. Wu, H. *et al.* Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808* (2021).

- 40. Yuan, K. *et al.* Incorporating Convolution Designs into Visual Transformers. *arXiv preprint arXiv:2103.11816* (2021).
- 41. Cordonnier, J.-B., Loukas, A. & Jaggi, M. On the relationship between self-attention and convolutional layers. *arXiv preprint arXiv:1911.03584* (2019).
- 42. d'Ascoli, S., Sagun, L., Biroli, G. & Bruna, J. Finding the Needle in the Haystack with Convolutions: on the benefits of architectural bias in Advances in Neural Information Processing Systems (2019), 9334–9345.
- 43. Wen, W., Yan, F., Chen, Y. & Li, H. Autogrow: Automatic layer growing in deep convolutional networks in Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (2020), 833–841.
- 44. Shaw, P., Uszkoreit, J. & Vaswani, A. Self-attention with relative position representations. *arXiv* preprint arXiv:1803.02155 (2018).
- 45. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition in Proceedings of the IEEE conference on computer vision and pattern recognition (2016), 770–778.
- 46. Touvron, H., Vedaldi, A., Douze, M. & Jégou, H. Fixing the train-test resolution discrepancy. *arXiv preprint arXiv:1906.06423* (2019).
- 47. Bello, I. *et al.* Revisiting ResNets: Improved Training and Scaling Strategies. *arXiv preprint arXiv:2103.07579* (2021).
- 48. Wightman, R. *PyTorch Image Models* https://github.com/rwightman/pytorch-image-models. 2019.
- 49. Hendrycks, D. & Dietterich, T. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. *Proceedings of the International Conference on Learning Representations* (2019).
- 50. Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J. & Song, D. Natural Adversarial Examples. *CVPR* (2021).
- 51. Hendrycks, D. *et al.* The Many Faces of Robustness: A Critical Analysis of Out-of-Distribution Generalization. *arXiv preprint arXiv:2006.16241* (2020).
- 52. Lewkowycz, A., Bahri, Y., Dyer, E., Sohl-Dickstein, J. & Gur-Ari, G. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218* (2020).
- 53. Caron, M. *et al.* Emerging properties in self-supervised vision transformers. *arXiv preprint arXiv:2104.14294* (2021).
- 54. He, T. et al. Bag of tricks for image classification with convolutional neural networks in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2019), 558–567.

A Changing of learning rate

As shown in Fig. 4 of the main text, the learning dynamics decompose into two phases: the learning rate warmup phase, where the test loss drops, then the learning rate decay phase, where the test loss increases again. This could lead one to think that the maximal learning rate is too high, and the dip could be avoided by choosing a lower learning rate. Yet this is not the case, as shown in Fig. 6. Reducing the maximal learning rate indeed reduces the dip, but it also slows down the increase in the second phase of learning. This confirms that the model needs to "unlearn" the right amount to find better solutions.



Figure 6: **The larger the learning rate, the lower the test accuracy dips, but the faster it climbs back up.** We show the dynamics of the ResNet50, fine-tuned for 50 epochs at resolution 224, for three different values of the maximal learning rate.

B Changing the test resolution

One advantage of the GPSA layers introduced by [38] is how easily they adapt to different image resolutions. Indeed, the positional embeddings they use are fixed rather than learnt. They simply consist in 3 values for each pair of pixels: their euclidean distance $||\delta||$, as well as their coordinate distance δ_1 , δ_2 (see Eq. 4). Our implementation automatically adjusts these embeddings to the input image, allowing us to change the test resolution seamlessly.

In Fig. 7, we show how the top-1 accuracies of our T-ResNet-RS models compares to those of the ResNet-RS models finetuned at same resolution but without SA. At test resolution 416, our T-ResNetRS-350 reaches an impressive top-1 accuracy of 84.9%, beyond those of the best EfficientNets and BotNets [32].

C Changing the number of epochs

In Tab. 4, we show how the top-1 accuracy of the T-ResNet-RS model changes with the number of fine-tuning epochs. As expected, performance increases significantly as we fine-tune for longer, yet we chose to set a maximum of 50 fine-tuning epochs to keep the computational cost of fine-tuning well below that of the original training.

D Changing the architecture

Our framework, which builds on the timm package, makes changing the original CNN architecture very easy. We applied our fine-tuning procedure to the ResNet-D models [54] with the exact same hyperparameters, and observed substantial performance gains, similar to the ones obtained for ResNet-RS, see Tab. 5. This suggests the wide applicability of our method.



Figure 7: **Performance at different test-time resolutions, for the finetuned models with and without SA.** The ResNet50-RS and ResNet101-RS models are finetuned at resolution 224, and all other models are finetuned at resolution 320.

| Model | Epochs | Top-1 acc |
|----------------|--------|-----------|
| ResNet50-RS | 0 | 79.91 |
| T-ResNet50-RS | 10 | 80.11 |
| T-ResNet50-RS | 20 | 80.51 |
| T-ResNet50-RS | 50 | 81.02 |
| ResNet101-RS | 0 | 81.70 |
| T-ResNet101-RS | 10 | 81.54 |
| T-ResNet101-RS | 20 | 81.90 |
| T-ResNet101-RS | 50 | 82.39 |

Table 4: **Longer fine-tuning increases final performance.** We report the top-1 accuracies of our models on ImageNet-1k at resolution 224.

| Model | Original res. | Original acc. | Fine-tune res. | Fine-tune acc. | Gain |
|----------------|---------------|---------------|----------------|----------------|------|
| T-ResNet50-D | 224 | 80.6 | 320 | 81.6 | +1.0 |
| T-ResNet101-D | 320 | 82.3 | 384 | 83.1 | +0.8 |
| T-ResNet152-D | 320 | 83.1 | 384 | 83.8 | +0.7 |
| T-ResNet200-D | 320 | 83.2 | 384 | 83.9 | +0.7 |
| T-ResNet50-RS | 160 | 78.8 | 224 | 81.0 | +2.8 |
| T-ResNet101-RS | 192 | 81.2 | 224 | 82.4 | +1.2 |
| T-ResNet152-RS | 256 | 83.0 | 320 | 83.7 | +0.7 |
| T-ResNet200-RS | 256 | 83.4 | 320 | 84.0 | +0.6 |

Table 5: **Comparing the performance gains of the ResNet-RS and ResNet-D architectures.** Top-1 accuracy is measured on ImageNet-1k validation set. The pre-trained models are all taken from the timm library [48].

E More attention maps



(b) Attention maps

Figure 8: **GPSA layers combine local and global attention in a complementary way.** We depicted the attention maps of the four GPSA layers of the T-ResNet270-RS, obtained by feeding the image on the left through the convolutional backbone, then selecting a query pixel in the center of the image (red box). For each head h, we indicate the value of the gating parameter $\sigma(\lambda_h)$ in red (see Eq. 7). ($\sigma(\lambda_h) = 0$).