

Bottleneck Transformers for Visual Recognition

Aravind Srinivas¹ Tsung-Yi Lin² Niki Parmar² Jonathon Shlens² Pieter Abbeel¹ Ashish Vaswani²
¹UC Berkeley ²Google Research
{aravind}@cs.berkeley.edu

Abstract

We present BoTNet, a conceptually simple yet powerful backbone architecture that incorporates self-attention for multiple computer vision tasks including image classification, object detection and instance segmentation. By just replacing the spatial convolutions with global self-attention in the final three bottleneck blocks of a ResNet and no other changes, our approach improves upon the baselines significantly on instance segmentation and object detection while also reducing the parameters, with minimal overhead in latency. Through the design of BoTNet, we also point out how ResNet bottleneck blocks with self-attention can be viewed as Transformer blocks. Without any bells and whistles, BoTNet achieves 44.4% Mask AP and 49.7% Box AP on the COCO Instance Segmentation benchmark using the Mask R-CNN framework; surpassing the previous best published single model and single scale results of ResNeSt [72] evaluated on the COCO validation set. Finally, we present a simple adaptation of the BoTNet design for image classification, resulting in models that achieve a strong performance of 84.7% top-1 accuracy on the ImageNet benchmark while being up to 2.33x faster in “compute”¹ time than the popular EfficientNet models on TPU-v3 hardware. We hope our simple and effective approach will serve as a strong baseline for future research in self-attention models for vision.²

1. Introduction

Deep convolutional backbone architectures [40, 56, 29, 70, 59] have enabled significant progress in image classification [54], object detection [17, 42, 21, 20, 53], instance segmentation [25, 12, 28]. Most landmark backbone architectures [40, 56, 29] use multiple layers of 3×3 convolutions.

While the convolution operation can effectively capture local information, vision tasks such as object detection, instance segmentation, keypoint detection require modeling long range dependencies. For example, in instance segmentation, being able to collect and associate scene information

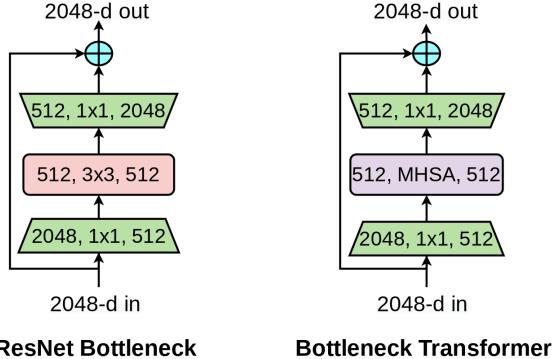


Figure 1: **Left:** A ResNet Bottleneck Block, **Right:** A Bottleneck Transformer (BoT) block. The only difference is the replacement of the spatial 3×3 convolution layer with Multi-Head Self-Attention (MHSA). The structure of the self-attention layer is described in Figure 4.

from a large neighborhood can be useful in learning relationships across objects [34]. In order to globally aggregate the locally captured filter responses, convolution based architectures require stacking multiple layers [56, 29]. Although stacking more layers indeed improves the performance of these backbones [72], an explicit mechanism to model global (non-local) dependencies could be a more powerful and scalable solution without requiring as many layers.

Modeling long-range dependencies is critical to natural language processing (NLP) tasks as well. Self-attention is a computational primitive [64] that implements pairwise entity interactions with a content-based addressing mechanism, thereby learning a rich hierarchy of associative features across long sequences. This has now become a standard tool in the form of Transformer [64] blocks in NLP with prominent examples being GPT [48, 3] and BERT [13, 44] models.

A simple approach to using self-attention in vision is to replace spatial convolutional layers with the multi-head self-attention (MHSA) layer proposed in the Transformer [64] (Figure 1). This approach has seen progress on two seemingly different approaches in the recent past. On the one hand, we have models such as SASA [51], SANet [73], Axial-SASA [66], etc that propose to replace spatial convo-

¹Forward and backward propagation for batch size 32

²Code and pre-trained models will be made available.

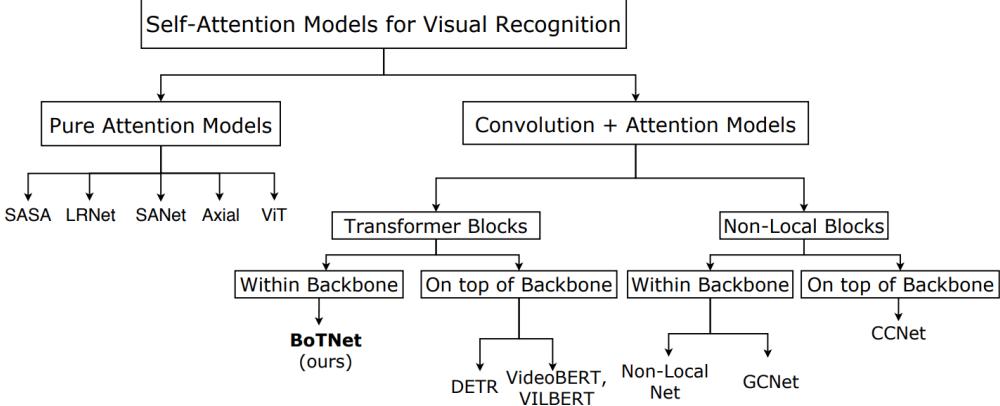


Figure 2: A taxonomy of deep learning architectures using self-attention for visual recognition. Our proposed architecture BoTNet is a hybrid model that uses both convolutions and self-attention. The specific implementation of self-attention could either resemble a Transformer block [64] or a Non-Local block [67] (difference highlighted in Figure 4). BoTNet is different from architectures such as DETR [8], VideoBERT [58], VILBERT [46], CCNet [37], etc by employing self-attention within the backbone architecture, in contrast to using them outside the backbone architecture. Being a hybrid model, BoTNet differs from pure attention models such as SASA [51], LRNet [35], SANet [73], Axial-SASA [33, 66] and ViT [14].

lutions in ResNet bottleneck blocks [29] with different forms of self-attention (local, global, vector, axial, etc). On the other hand, we have the Vision Transformer (ViT) [14], that proposes to stack Transformer blocks [64] operating on linear projections of non-overlapping patches. It may appear that these approaches present two different classes of architectures. We point out that it is *not the case*. Rather, ResNet bottleneck blocks with the MHSA layer can be viewed as Transformer blocks with a bottleneck structure, modulo minor differences such as the residual connections, choice of normalization layers, etc. (Figure 3). Given this equivalence, we call ResNet bottleneck blocks with the MHSA layer as *Bottleneck Transformer* (BoT) blocks.

Here are a few challenges when using self-attention in vision: (1) Image sizes are much larger (1024×1024) in object detection and instance segmentation compared to image classification (224×224). (2) The memory and computation for self-attention scale quadratically with spatial dimensions [61], causing overheads for training and inference.

To overcome these challenges, we consider the following design: (1) Use convolutions to *efficiently* learn *abstract* and *low resolution* featuremaps from large images; (2) Use global (*all2all*) self-attention to process and aggregate the information contained in the featuremaps captured by convolutions. Such a hybrid design (1) uses existing and well optimized primitives for both convolutions and all2all self-attention; (2) can deal with large images efficiently by having convolutions do the spatial downsampling and letting attention work on smaller resolutions. Here is a simple practical instantiation of this hybrid design: Replace *only* the final three bottleneck blocks of a ResNet with BoT blocks *without any other*

changes. Or in other words, take a ResNet and only replace the final three 3×3 convolutions with MHSA layers (Fig 1, Table 1). This simple change improves the mask AP by 1.2% on the COCO instance segmentation benchmark [42] over our canonical baseline that uses ResNet-50 in the Mask R-CNN framework [28] with *no hyperparameter differences* and minimal overheads for training and inference. Moving forward, we call this simple instantiation as BoTNet given its connections to the Transformer through the BoT blocks. While we note that there is no novelty in its construction, we believe the simplicity and performance make it a useful reference backbone architecture that is worth studying.

Using BoTNet, we demonstrate significantly improved results on instance segmentation and object detection *without any bells and whistles* such as Cascade R-CNN [5], FPN changes [43, 19, 45, 60], hyperparameter changes [59], etc. A few key results from BoTNet are listed below: (1) Performance gains across various training configurations (Section 4.1), data augmentations (Section 4.2) and ResNet family backbones (Section 4.5); (2) Significant boost from BoTNet on small objects (+2.4 Mask AP and +2.6 Box AP) (Section 4.4, Section A.3); (3) Performance gains over Non-Local layers (Section 4.7); (4) Gains that scale well with larger images resulting in **44.4%** mask AP, competitive with state-of-the-art performance among entries that only study backbone architectures with modest training schedules (up to 72 epochs) and no extra data or augmentations.³.

We finally offer a simple adaptation of the BoTNet ar-

³State-of-the-art notion is based only on the entries in <https://paperswithcode.com/sota/instance-segmentation-on-coco-minival>. We also note *unlike* the ResNeSt-200 result, we do not use Cascade R-CNN which could give us additional gains.

chitecture for good gains on the ImageNet benchmark [54] after noting that a straightforward use of BoTNet does not provide substantial gains. Using this adaptation, we design a family of BoTNet models that achieve up to **84.7%** top-1 accuracy on the ImageNet validation set, while being upto **2.33x** faster than the popular EfficientNet models in terms of *compute* time on TPU-v3 hardware. By providing *strong results* through BoTNet, we hope that self-attention becomes a widely used primitive in future vision architectures.

2. Related Work

A taxonomy of deep learning architectures that employ self-attention for vision is presented in Figure 2. In this section, we focus on: (1) Transformer vs BoTNet; (2) DETR vs BoTNet; (3) Non-Local vs BoTNet.

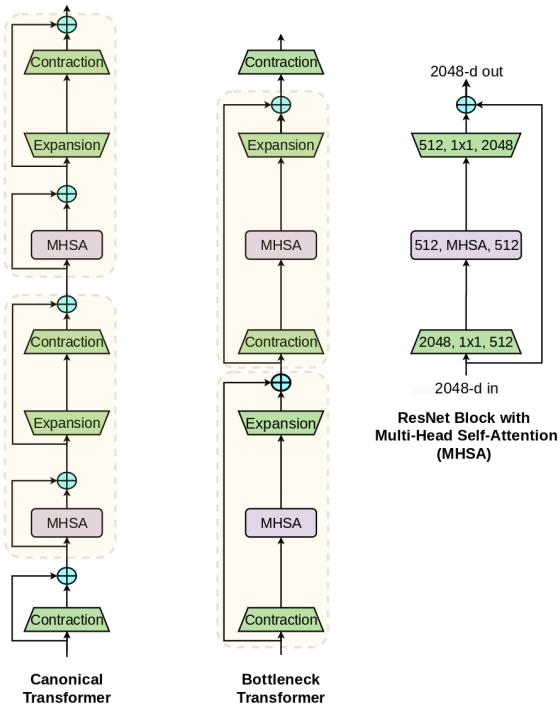


Figure 3: Left: Canonical view of the Transformer with the boundaries depicting the definition of a Transformer block as described in Vaswani et. al [64]. **Middle:** Bottleneck view of the Transformer with boundaries depicting what we define as the Bottleneck Transformer (BoT) block in this work. The architectural structure that already exists in the Transformer can be interpreted a ResNet bottleneck block [29] with Multi-Head Self-Attention (MHSA) [64] with a different notion of block boundary as illustrated. **Right:** An instantiation of the Bottleneck Transformer as a ResNet bottleneck block [29] with the difference from a canonical ResNet block being the replacement of 3×3 convolution with MHSA.

Connection to the Transformer: As the title of the paper suggests, one key message in this paper is that ResNet bottleneck blocks with Multi-Head Self-Attention (MHSA) layers can be viewed as Transformer blocks with a bottleneck structure. This is visually explained in Figure 3 and we name this block as Bottleneck Transformer (BoT). We note that the architectural design of the BoT block is not our contribution. Rather, we point out the relationship between MHSA ResNet bottleneck blocks and the Transformer with the hope that it improves our understanding of architecture design spaces [49, 50] for self-attention in computer vision. There are still a few differences aside from the ones already visible in the figure (residual connections and block boundaries): (1) Normalization: Transformers use Layer Normalization [1] while BoT blocks use Batch Normalization [38] as is typical in ResNet bottleneck blocks [29]; (2) Non-Linearities: Transformers use one non-linearity in the FFN block, while the ResNet structure allows BoT block to use three non-linearities; (3) Output projections: The MHSA block in a Transformer contains an output projection while the MHSA layer (Fig 4) in a BoT block (Fig 1) does not; (4) We use the SGD with momentum optimizer typically used in computer vision [29, 28, 22] while Transformers are generally trained with the Adam optimizer [39, 64, 8, 14].

Connection to DETR: Detection Transformer (DETR) is a detection framework that uses a Transformer to implicitly perform region proposals and localization of objects instead of using an R-CNN [21, 20, 53, 28]. Both DETR and BoT-Net attempt to use self-attention to improve the performance on object detection and instance (or panoptic) segmentation. The difference lies in the fact that DETR uses Transformer blocks outside the backbone architecture with the motivation to get rid of region proposals and non-maximal suppression for simplicity. On the other hand, the goal in BoTNet is to provide a backbone architecture that uses Transformer-like blocks for detection and instance segmentation. We are agnostic to the detection framework (be it DETR or R-CNN). We perform our experiments with the Mask [28] and Faster R-CNN [53] systems and leave it for future work to integrate BoTNet as the backbone in the DETR framework. With visibly good gains on small objects in BoTNet, we believe there maybe an opportunity to address the lack of gain on small objects found in DETR, in future (Section A.3).

Connection to Non-Local Neural Nets: Non-Local (NL) Nets [67] make a connection between the Transformer and the Non-Local-Means algorithm [4]. They insert NL blocks into the final one (or) two blockgroups (c_4, c_5) in a ResNet and improve the performance on video recognition and instance segmentation. Like NL-Nets [67, 6], BoTNet is a hybrid design using convolutions and global self-attention. (1) Three differences between a NL layer and a MHSA layer (illustrated in Figure 4): use of multiple heads, value projection and po-

sition encodings in MHSA (NL block doesn't use them as per the implementation in https://github.com/facebookresearch/ImageNet-Adversarial-Training/blob/master/resnet_model.py#L92); (2) NL blocks use a bottleneck with channel factor reduction of 2 (instead of 4 in BoT blocks which adopt the ResNet structure); (3) NL blocks are *inserted* as *additional* blocks into a ResNet backbone as opposed to *replacing* existing convolutional blocks as done by BoTNet. Section 4.7 offers a comparison between BoTNet, NLNet as well as a NL-like version of BoTNet where we *insert* BoT blocks in the same manner as NL blocks instead of replacing.

3. Method

stage	output	ResNet-50	BoTNet-50
c1	512×512	$7 \times 7, 64$, stride 2	$7 \times 7, 64$, stride 2
		3×3 max pool, stride 2	3×3 max pool, stride 2
c2	256×256	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
c3	128×128	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
c4	64×64	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
c5	32×32	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ \text{MHSA, } 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
# params.		25.5×10^6	20.8×10^6
M.Adds		85.4×10^9	102.98×10^9
TPU steptime		786.5 ms	1032.66 ms

Table 1: Architecture of BoTNet-50 (BoT50): The only difference in BoT50 from ResNet-50 (R50) is the use of MHSA layer (Figure 4) in c5. For an input resolution of 1024×1024 , the MHSA layer in the first block of c5 operates on 64×64 while the remaining two operate on 32×32 . We also report the parameters, multiply-adds (m. adds) and training time throughput (TPU-v3 steptime on a v3-8 Cloud-TPU). BoT50 has only 1.2x more m.adds. than R50. The overhead in training throughout is 1.3x. BoT50 also has 1.2x fewer parameters than R50. While it may appear that it is simply the aspect of performing slightly more computations that might help BoT50 over the baseline, we show that it is not the case in Section 4.5.

BoTNet by design is simple: replace the final three spatial (3×3) convolutions in a ResNet with Multi-Head Self-Attention (MHSA) layers that implement global (*all2all*) self-attention over a 2D featuremap (Fig 4). A ResNet typically has 4 stages (or blockgroups) commonly referred to

as $[c2, c3, c4, c5]$ with strides $[4, 8, 16, 32]$ relative to the input image, respectively. Stacks $[c2, c3, c4, c5]$ consist of multiple *bottleneck* blocks with residual connections (e.g, R50 has $[3, 4, 6, 3]$ bottleneck blocks).

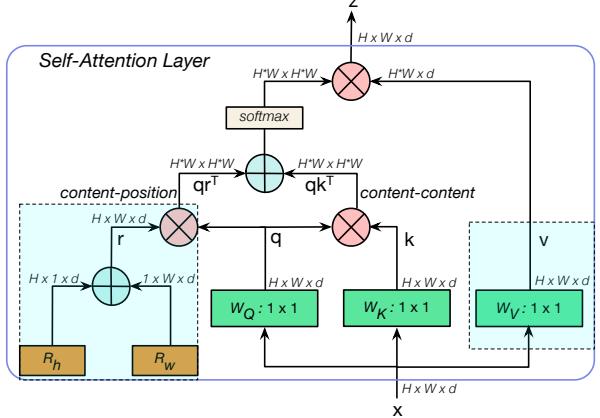


Figure 4: Multi-Head Self-Attention (MHSA) layer used in the BoT block. While we use 4 heads, we do not show them on the figure for simplicity. *all2all* attention is performed on a 2D featuremap with *split* relative position encodings R_h and R_w for height and width respectively. The attention logits are $qk^T + qr^T$ where q, k, r represent query, key and position encodings respectively (we use relative distance encodings [55, 2, 51]). \oplus and \otimes represent element wise sum and matrix multiplication respectively, while 1×1 represents a pointwise convolution. Along with the use of multiple heads, the highlighted blue boxes (**position encodings** and the **value projection**) are the *only* three elements that are not present in the Non-Local Layer [67, 69].

Approaches that use self-attention throughout the backbone [51, 2, 73, 14] are feasible for input resolutions (224×224 (for classification) and 640×640 (for detection experiments in SASA [51])) considered in these papers. Our goal is to use attention in more realistic settings of high performance instance segmentation models, where typically images of larger resolution (1024×1024) are used. Considering that self-attention when performed globally across n entities requires $O(n^2d)$ memory and computation [64], we believe that the simplest setting that adheres to the above factors would be to incorporate self-attention at the lowest resolution featuremaps in the backbone, ie, the residual blocks in the c5 stack. The c5 stack in a ResNet backbone typically uses 3 blocks with one spatial 3×3 convolution in each. Replacing them with MHSA layers forms the basis of the BoTNet architecture. The first block in c5 uses a 3×3 convolution of stride 2 while the other two use a stride of 1. Since *all2all* attention is not a strided operation, we use a 2×2 average-pooling with a stride 2 for the first BoT block. The BoTNet architecture is described in Table 1

Backbone	epochs	AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅ ^{bb}	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
R50	12	39.0	59.4	42.7	35.0	56.4	37.4
BoT50	12	39.4 (+ 0.4)	60.3 (+ 0.9)	43 (+ 0.3)	35.3 (+ 0.3)	57 (+ 0.6)	37.5 (+ 0.1)
R50	24	41.2	61.8	45.1	36.9	58.8	39.8
BoT50	24	42.8 (+ 1.6)	64.1 (+ 2.3)	46.9 (+ 1.8)	38.0 (+ 1.1)	60.9 (+ 2.1)	40.5 (+ 0.7)
R50	36	42.1	62.5	46.3	37.7	59.6	40.5
BoT50	36	43.6 (+ 1.5)	65.3 (+ 2.8)	47.6 (+ 1.3)	38.9 (+ 1.2)	62.5 (+ 2.9)	41.3 (+ 0.8)
R50	72	42.8	63.0	46.9	37.9	60.0	40.7
BoT50	72	43.7 (+ 0.9)	64.7 (+ 1.7)	47.9 (+ 1.0)	38.7 (+ 0.8)	61.8 (+ 1.8)	41.1 (+ 0.4)

Table 2: Comparing R50 and BoT50 under the 1x (12 epochs), 3x (36 epochs) and 6x (72 epochs) settings, trained with image resolution 1024×1024 and multi-scale jitter of $[0.8, 1.25]$.

and the MHSA layer is presented in Figure 4. The strided version of the BoT block is presented in Figure 12.

Relative Position Encodings: In order to make the attention operation *position aware*, Transformer based architectures typically make use of a position encoding [64]. It has been observed lately that *relative-distance-aware* position encodings [55] are better suited for vision tasks [2, 51, 73]. This can be attributed to attention not only taking into account the content information but also relative distances between features at different locations, thereby, being able to effectively associate information across objects with positional awareness. In BoTNet, we adopt the 2D relative position self-attention implementation from [51, 2].

4. Experiments

We study the benefits of BoTNet for instance segmentation and object detection. We perform a thorough ablation study of various design choices through experiments on the COCO dataset [42]. We report the standard COCO metrics including the AP^{bb} (averaged over IoU thresholds), AP₅₀^{bb}, AP₇₅^{bb}; AP₅₀^{mk}, AP₇₅^{mk} for box and mask respectively. As is common practice these days, we train using the COCO train set and report results on the COCO val (or minival) set as followed in Detectron [22]⁴. Our experiments are based on the Google Cloud TPU detection codebase⁵. We run all the baselines and ablations with the same codebase. Unless explicitly specified, our training infrastructure uses v3-8 Cloud-TPU which contains 8 cores with 16 GB memory per core. We train with the bfloat16 precision and cross-replica batch normalization [38, 68, 28, 22, 47] using a batch size of 64.

⁴train - 118K images, val - 5K images

⁵<https://github.com/tensorflow/tpu/tree/master/models/official/detection>

4.1. BoTNet improves over ResNet on COCO Instance Segmentation with Mask R-CNN

We consider the simplest and most widely used setting: ResNet-50⁶ backbone with FPN⁷. We use images of resolution 1024×1024 with a multi-scale jitter of $[0.8, 1.25]$ (scaling the image dimension between 820 and 1280, in order to be consistent with the Detectron setting of using 800×1300). In this setting, we benchmark both the ResNet-50 (R50) and BoT ResNet-50 (BoT50) as the backbone architectures for multiple training schedules: **1x**: 12 epochs, **2x**: 24 epochs, **3x**: 36 epochs, **6x**: 72 epochs⁸, all using the same hyperparameters for both the backbones across all the training schedules (Table 2). We clearly see that BoT50 is a significant improvement on top of R50 barring the 1x schedule (12 epochs). This suggests that BoT50 warrants longer training in order to show significant improvement over R50. We also see that the improvement from BoT50 in the 6x schedule (72 epochs) is worse than its improvement in the 3x schedule (32 epochs). This suggests that training much longer with the default scale jitter hurts. We address this by using a more aggressive scale jitter (Section 4.2).

4.2. Scale Jitter helps BoTNet more than ResNet

In Section 4.1, we saw that training much longer (72 epochs) reduced the gains for BoT50. One way to address this is to increase the amount of multi-scale jitter which has been known to improve the performance of detection and segmentation systems [15, 18]. Table 3 shows that BoT50 is significantly better than R50 (+ 2.1% on AP^{bb} and + 1.7% on AP^{mk}) for multi-scale jitter of $[0.5, 2.0]$, while also showing significant gains (+ 2.2% on AP^{bb} and + 1.6% on

⁶We use the ResNet backbones pre-trained on ImageNet classification as is common practice. For BoTNet, the replacement layers are **not** pre-trained but randomly initialized for simplicity; the remaining layers are initialized from a pre-trained ResNet.

⁷FPN refers to Feature Pyramid Network [41]. We use it in every experiment we report results on, and our FPN levels from 2 to 6 (p2 to p6) similar to Detectron [22].

⁸1x, 2x, 3x and 6x convention is adopted from MoCo [26].

AP^{mk}) for scale jitter of $[0.1, 2.0]$, suggesting that BoTNet (self-attention) benefits more from extra augmentations such as multi-scale jitter compared to ResNet (pure convolutions).

Backbone	jitter	AP^{bb}	AP^{mk}
R50	$[0.8, 1.25]$	42.8	37.9
BoT50	$[0.8, 1.25]$	43.7 (+ 0.9)	38.7 (+ 0.8)
R50	$[0.5, 2.0]$	43.7	39.1
BoT50	$[0.5, 2.0]$	45.3 (+ 1.8)	40.5 (+ 1.4)
R50	$[0.1, 2.0]$	43.8	39.2
BoT50	$[0.1, 2.0]$	45.9 (+ 2.1)	40.7 (+ 1.5)

Table 3: Comparing R50 and BoT50 under three settings of multi-scale jitter, all trained with image resolution 1024×1024 for 72 epochs (6x training schedule).

4.3. Relative Position Encodings Boost Performance

BoTNet uses relative position encodings [55]. We present an ablation for the use of relative position encodings by benchmarking the individual gains from content-content interaction (qk^T) and content-position interaction (qr^T) where q, k, r represent the query, key and relative position encodings respectively. The ablations (Table 4) are performed with the canonical setting⁹. We see that the gains from qr^T and qk^T are complementary with qr^T more important, ie, qk^T standalone contributes to 0.6% AP^{bb} and 0.6% AP^{mk} improvement over the R50 baseline, while qr^T standalone contributes to 1.0% AP^{bb} and 0.7 % AP^{mk} improvement. When combined together ($qk^T + qr^T$), the gains on both AP^{bb} and AP^{mk} are additive (1.5% and 1.2% respectively). We also see that using absolute position encodings (qr_{abs}^T) does not provide as much gain as relative. This suggests that introducing relative position encodings into architectures like DETR [8] is an interesting direction for future work.

Backbone	Att. Type	AP^{bb}	AP^{mk}
R50	-	42.1	37.7
BoT50	qk^T	42.7 (+ 0.6)	38.3 (+ 0.6)
BoT50	qr^T_{relative}	43.1 (+ 1.0)	38.4 (+ 0.7)
BoT50	$qk^T + qr^T_{\text{relative}}$	43.6 (+ 1.5)	38.9 (+ 1.2)
BoT50	$qk^T + qr^T_{\text{abs}}$	42.5 (+ 0.4)	38.1 (+ 0.4)

Table 4: Ablation for Relative Position Encoding: Gains from the two types of interactions in the MHSA layers, content-content (qk^T) and content-position (qr^T).

4.4. Why replace all three c5 spatial convolutions?

Is replacing all the 3 spatial convolutions in c5 the *minimum effective* change or could it be simplified even further?

⁹res:1024x1024, 36 epochs (3x schedule), multi-scale jitter:[0.8, 1.25]

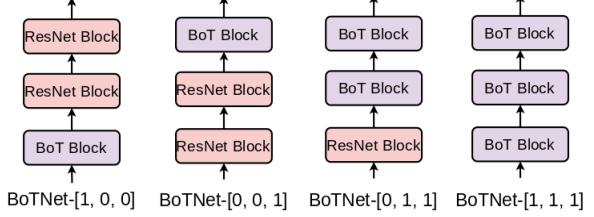


Figure 5: Replacement configs for BoTNet in the c5 block-group of a ResNet

We perform an ablation study on the replacement design in order to answer this question (Please refer to Table 5, Figure 5). The baseline R50 and BoT50 go by the notation $[0, 0, 0]$ and $[1, 1, 1]$ since the former does not replace anything while the latter replaces the spatial convolution in all three c5 blocks. As mentioned already, the first replacement in c5 operates on 64×64 feature map while the remaining two operate on 32×32 . We ablate for the configs: $[0, 0, 1]$, $[0, 1, 1]$ and $[1, 0, 0]$. The first two configs test how useful is the replacement when performed only on the smaller 32×32 featuremap(s) once and twice respectively, while the last tests how useful is the replacement when performed only on the larger 64×64 featuremap.

First, we see that in terms of aggregate measures such as AP^{bb} and AP^{mk} , each of the configs for BoT50 is a strict improvement over R50, with similar performance. Config. $[1, 0, 0]$ is closer to the performance of BoT50 ($[1, 1, 1]$) compared to the other configurations, however a difference of 0.2 AP^{bb} is within the noise typically observed in COCO experiments. It is indeed surprising that just a single self-attention replacement layer right at the end ($[0, 0, 1]$) provides a visible gain of 1.3 AP^{bb} . When contrasted with the performance of R101 (43.2 AP^{bb} and 38.4 AP^{mk}), the config. $[0, 0, 1]$ is very much competitive with 43.4 AP^{bb} and 38.6 AP^{mk} , with more efficient compute step-time on the TPU (1.2x faster). Nevertheless, the gains on large objects for the $[0, 0, 1]$ config (+ 0.6 AP_L^{bb}) are not as significant as those in R101 (+ 1.6 AP_L^{bb}).

Among the different configs for BoT50, we see that $[1, 0, 0]$ and $[0, 1, 1]$ are the best in terms of good performance on both small and large objects. Surprisingly, the actual BoTNet config ($[1, 1, 1]$) shows significant boost on small objects (2.6 AP_S^{bb}), but does not show substantial gain on large objects, even relative to other ablation configs. We suspect this could be due to poor optimization and leave it for future work to carefully understand how self-attention affects the performance on small and large objects.

Based on these ablations, consider the question: is it better to replace convolutions with self-attention in c5 (BoT50) vs stacking more convolution layers (R101)? An argument in favor of R101 is that the gains are clear on both small and large objects unlike BoT50 where the gains are much

Backbone	c5-attn.	TC Time	AP ^{bb}	AP _S ^{bb}	AP _L ^{bb}	AP ^{mk}	AP _S ^{mk}	AP _L ^{mk}
R50	[0, 0, 0]	786.5	42.1	22.5	59.1	37.7	18.3	54.9
BoT 50	[0, 0, 1]	813.7	43.4 (+ 1.3)	23.7 (+ 1.2)	59.7 (+ 0.6)	38.6 (+ 0.9)	19.0 (+ 0.7)	55.6 (+ 0.7)
BoT 50	[0, 1, 1]	843.87	43.4 (+ 1.3)	24.0 (+ 1.5)	60.2 (+ 1.1)	38.6 (+ 0.9)	19.4 (+ 1.1)	55.9 (+ 1.0)
BoT R50	[1, 0, 0]	983.2	43.7 (+ 1.6)	38.9 (+ 1.2)	23.9 (+ 1.4)	60.6 (+ 1.5)	19.3 (+ 1.0)	55.9 (+ 1.0)
BoT 50	[1, 1, 1]	1032.7	43.6 (+ 1.5)	25.1 (+ 2.6)	59.4 (+ 0.3)	38.9 (+ 1.2)	20.7 (+ 2.4)	55.5 (+ 0.6)
R101	[0, 0, 0]	928.7	43.3 (+ 1.2)	24.2 (+ 1.7)	60.7 (+ 1.6)	38.4 (+ 0.7)	19.6 (+ 1.3)	56.8 (+ 1.9)
BoT101	[1, 1, 1]	1174.9	45.5 (+ 2.2)	26.0 (+ 1.8)	62.3 (+ 1.6)	40.4 (+ 2.0)	21.1 (+ 1.5)	58.0 (+ 1.2)

Table 5: Ablation study on the replacement design in BoTNet: All models are trained with the canonical config with image size 1024×1024 , jitter [0.8, 1.25], 36 epochs. TC Time refers to TPU-v3 Compute Step Time (in milliseconds) during training.

more on small objects. However, there does exist a BoT50 config that can strictly improve upon R101, ie the [0, 1, 1] config. It has similar properties to R101 (gain on both small and large objects), similar performance on aggregate measures like AP^{bb} and AP^{mk}, with a more efficient compute steptime. Hence, we can affirmatively say that *self-attention replacement is more efficient than stacking convolutions*.

4.5. BoTNet improves backbones in ResNet Family

How well does the replacement setup of BoTNet work for other backbones in the ResNet family? Table 6 presents the results for BoTNet with R50, R101, and R152. All these experiments use the canonical training setting (refer to footnote in 4.3). These results demonstrate that BoTNet is applicable as a drop-in replacement for any ResNet backbone. Note that BoT50 is better than R101 (+ 0.3% AP^{bb}, + 0.5% AP^{mk}) while it is competitive with R152 on AP^{mk}. Replacing 3 spatial convolutions with all2all attention gives more improvement in the metrics compared to stacking 50 more layers of convolutions (R101), and is competitive with stacking 100 more layers (R152), supporting our initial hypothesis that long-range dependencies are better captured through attention than stacking convolution layers.¹⁰

4.6. BoTNet scales well with larger images

We benchmark BoTNet as well as baseline ResNet when trained on 1280×1280 images in comparison to 1024×1024 using the best config: multi-scale jitter of [0.1, 2.0] and training for 72 epochs. Results are presented in Tables 7 and 9. Results in Table 7 suggest that BoTNet benefits from training on larger images for all of R50, R101 and R152. BoTNet trained on 1024×1024 (leave alone 1280×1280) is significantly better than baseline ResNet trained on 1280×1280 . Further, BoT200 trained with 1280×1280 achieves a AP^{bb} of 49.7% and AP^{mk} of 44.4%. We believe this result

¹⁰Note that while one may argue that the improvements of BoT50 over R50 could be attributed to having 1.2x more M. Adds, BoT50 (121×10^9 M.Adds) is also better than R101 (162.99×10^9 B M. Adds) and is competitive with R152 (240.56×10^9 M. Adds) despite performing significantly less computation.

Backbone	AP ^{bb}	AP ^{mk}
R50	42.1	37.7
BoT50	43.6 (+ 1.5)	38.9 (+ 1.2)
R101	43.3	38.4
BoT101	45.5 (+ 2.2)	40.4 (+ 2.0)
R152	44.2	39.1
BoT152	46.0 (+ 1.8)	40.6 (+ 1.5)

Table 6: Comparing R50, R101, R152, BoT50, BoT101 and BoT152; all 6 setups using the canonical training schedule of 36 epochs, 1024×1024 images, multi-scale jitter [0.8, 1.25].

highlights the power of self-attention, in particular, because it has been achieved without any bells and whistles such as modified FPN [43, 19, 15, 60], cascade RCNN [5], etc. This result surpasses the previous best published single model single scale instance segmentation result from ResNeSt [72] evaluated on the COCO minival (44.2% AP^{mk}).

Backbone	res	AP ^{bb}	AP ^{mk}
R50	1280	44.0	39.5
BoT50	1024	45.9 (+ 1.9)	40.7 (+ 1.2)
BoT50	1280	46.1 (+ 2.1)	41.2 (+ 1.8)
R101	1280	46.4	41.2
BoT101	1024	47.4 (+ 1.0)	42.0 (+ 0.8)
BoT101	1280	47.9 (+ 1.5)	42.4 (+ 1.2)

Table 7: All the models are trained for 72 epochs with a multi-scale jitter of [0.1, 2.0].

4.7. Comparison with Non-Local Neural Networks

How does BoTNet compare to Non-Local Neural Networks? NL ops are *inserted* into the c4 stack of a ResNet backbone between the pre-final and final bottleneck blocks. This *adds* more parameters to the model, whereas BoTNet ends up reducing the model parameters (Table 6). In the

Backbone	Change in backbone	AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅ ^{bb}	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
R50	-	42.1	62.5	46.3	37.7	59.6	40.5
R50 + NL [67]	+ 1 NL block in c4	43.1 (+ 1.0)	64.0 (+ 1.5)	47.4 (+ 1.1)	38.4 (+ 0.7)	61.1 (+ 1.5)	40.9 (+ 0.4)
R50 + BoT (c4)	+ 1 BoT block in c4	43.7 (+ 1.6)	64.7 (+ 2.2)	47.8 (+ 1.5)	38.9 (+ 1.2)	61.6 (+ 2.0)	41.8 (+ 1.3)
R50 + BoT (c4, c5)	+ 2 BoT blocks in c4, c5	44.9 (+ 2.8)	66.0 (+ 3.5)	49.0 (+ 2.7)	39.7 (+ 2.0)	62.9 (+ 3.3)	43.5 (+ 2.5)
BoT50	Replacement in c5	43.6 (+ 1.5)	65.3 (+ 2.8)	47.6 (+ 1.3)	38.9 (+ 1.2)	62.5 (+ 2.9)	41.3 (+ 0.8)

Table 8: Comparison between BoTNet and Non-Local (NL) Nets: All models trained for 36 epochs with image size 1024×1024 , jitter $[0.8, 1.25]$.

Backbone	AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅ ^{bb}	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
BoT152	49.5	71.0	54.2	43.7	68.2	47.4
BoT200	49.7	71.3	54.6	44.4	68.9	48.2

Table 9: BoT152 and BoT200 trained for 72 epochs with a multi-scale jitter of $[0.1, 2.0]$.

NL mould, we add ablations where we introduce BoT block in the exact same manner as the NL block. We also run an ablation with the insertion of two BoT blocks, one each in the c4, c5 stacks. Results are presented in Table 8. Adding a NL improves AP^{bb} by 1.0 and AP^{bb} by 0.7, while adding a BoT block gives +1.6 AP^{bb} and +1.2 AP^{mk} showing that BoT block design is better than NL. Further, BoT-R50 (which replaces instead of adding new blocks) provides +1.5 AP^{bb} and + 1.2 AP^{mk}, as good as adding another BoT block and better than adding one additional NL block.

4.8. Image Classification on ImageNet

4.8.1 BoTNet-S1 architecture

While we motivated the design of BoTNet for detection and segmentation, it is a natural question to ask whether the BoTNet architecture design also helps improve the image classification performance on the ImageNet [54] benchmark. Prior work [69] has shown that *adding* Non-Local blocks to ResNets and training them using canonical settings does *not* provide substantial gains. We observe a similar finding for BoTNet-50 when contrasted with ResNet-50, with both models trained with the canonical hyperparameters for ImageNet [50]: 100 epochs, batch size 1024, weight decay 1e-4, standard ResNet data augmentation, cosine learning rate schedule (Table 10). BoT50 does *not* provide significant gains over R50 on ImageNet though it does provide the benefit of reducing the parameters while maintaining comparable computation (M.Adds).

A simple method to fix this lack of gain is to take advantage of the image sizes typically used for image classification. In image classification, we often deal with much smaller image sizes (224×224) compared to those used in object detection and segmentation (1024×1024). The featuremaps

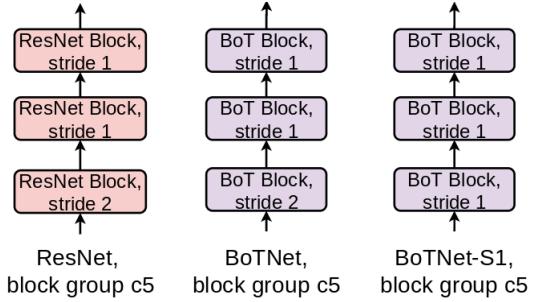


Figure 6: The c5 (final) block groups for ResNet (left), BoTNet (middle) and BoTNet-S1 (right).

on which the BoT blocks operate are hence much smaller (e.g 14×14 , 7×7) compared to those in instance segmentation and detection (e.g 64×64 , 32×32). With the same number of parameters, and, without a significant increase in computation, the BoTNet design in the c5 blockgroup can be changed to uniformly use a stride of 1 in all the final MHSA layers. We call this design as BoTNet-S1 (S1 to depict stride 1 in the final blockgroup). We note that this architecture is similar in design to the hybrid models explored in Vision Transformer (ViT) [14] that use a ResNet up to stage c4 prior to stacking Transformer blocks. The main difference between BoTNet-S1 and the hybrid ViT models lies in the use of BoT blocks as opposed to regular Transformer blocks (other differences being normalization layer, optimizer, etc as mentioned in the contrast to Transformer in Related Work (Sec. 2)). The architectural distinction amongst ResNet, BoTNet and BoTNet-S1, in the final blockgroup, is visually explained in Figure 6). The strided BoT block is visually explained in Figure 12.

4.8.2 Evaluation in the standard training setting

We first evaluate this design for the 100 epoch setting along with R50 and BoT50. We see that BoT-S1-50 improves on top of R50 by 0.9% in the regular setting (Table 10). This improvement does however come at the cost of more computation (m.adds). Nevertheless, the improvement is a promising signal for us to design models that scale well with

larger images and improved training conditions that have become more commonly used since EfficientNets [59].

Backbone	M.Adds	Params	top-1 acc.
R50	3.86G	25.5M	76.8
BoT50	3.79G	20.8M	77.0 (+0.2)
BoT-S1-50	4.27G	20.8M	77.7 (+0.9)

Table 10: ImageNet results in regular training setting: 100 epochs, batch size 1024, weight decay 1e-4, standard ResNet augmentation, for all three models.

4.8.3 Effect of data augmentation and longer training

We saw from our instance segmentation experiments that BoTNet and self-attention benefit more from regularization such as data augmentation (in the case of segmentation, increased multi-scale jitter) and longer training. It is natural to expect that the gains from BoT and BoT-S1 could improve when training under an improved setting: 200 epochs, batch size 4096, weight decay 8e-5, RandAugment (2 layers, magnitude 10), and label smoothing of 0.1. In line with our intuition, the gains are much more significant in this setting for both BoT50 (+ 0.6%) and BoT-S1-50 (+ 1.4%) compared to the baseline R50 (Table 11).

Backbone	top-1 acc.	top-5 acc.
R50	77.7	93.9
BoT50	78.3 (+0.6)	94.2 (+0.3)
BoT-S1-50	79.1 (+1.4)	94.4 (+0.5)

Table 11: ImageNet results in an improved training setting: 200 epochs, batch size 4096, weight decay 8e-5, RandAugment (2 layers, magnitude 10), and label smoothing of 0.1

4.8.4 Effect of SE blocks, SiLU and lower weight decay

Other aspects of improved training of backbone architectures for image classification has been the use of Squeeze-Excitation (SE) blocks [36], SiLU non-linearity [16, 52, 32] and further lowering the weight decay (for eg., EfficientNet uses 1e-6). When employing SE blocks in BoTNet and BoTNet-S1, we only do so for the ResNet blocks that employ 3×3 convolutions since self-attention is already designed for contextual global pooling. As expected, these changes lead to further improvements in the accuracy for all the models, with the gains from BoTNet remaining intact over the baseline SENet (ResNet with SE blocks) (Table 12).

The improvements from Tables 10, 11, 12 suggest that BoTNets are a good replacement for ResNets and SENets, especially when trained with data augmentations and longer

Backbone	top-1 acc.	top-5 acc.
SE50	79.2	94.6
BoT50	79.6 (+0.4)	94.6
BoT-S1-50	80.4 (+1.2)	95.0 (+0.4)

Table 12: ImageNet results in a further improved training setting with SE blocks and SiLU non-linearity: 200 epochs, batch size 4096, weight decay 4e-5, RandAugment (2 layers, magnitude 10), and label smoothing of 0.1. R50 with SE blocks is referred to as SE50.

training schedules. We have also made sure to present strong baselines (eg. 79.2% top-1 acc. SENet-50).

4.8.5 Comparison to EfficientNets

The previous ablations confirm that BotNets and ResNets both benefit from improved training settings such as lower weight decay, longer training, better augmentations, and the use of SiLU non-linearity. What if we go further in terms of improvement in training and use all the elements from the EfficientNet training regime? More precisely, training for 350 epochs with SE blocks, SiLU, RandAugment, lower weight decay, increased input image resolutions, DropConnect [65] and Dropout [57] (final linear layer for the logits). As in Table 12, we will refer to ResNets trained with SE blocks as SENets [36]. We benchmark a variety of models under this improved setting in Table 14. For the block configuration of BoTNet-S1 and SENet-350 presented in Table 14, please refer to Table 13. ResNets and SENets (50, 101, 152) use the standard block configuration and strides. Table 14 reveals the following key findings:

- BoTNets are significantly more efficient in terms of compute step time than EfficientNets, in the higher accuracy regime (B4 and beyond).
- ResNets and SENets perform really well in the lower accuracy regime, outperforming both EfficientNets and BoTNets (until and including B4).
- EfficientNets may be better in terms of M.Adds, but do not map as well as BoTNets, onto the latest hardware accelerators such as TPUs. For pretty much every single EfficientNet model, there exists a better ResNet, SENet or a BoTNet (similar accuracy with more efficient compute time).
- The input image sizes required for BoTNets are much smaller compared to EfficientNets suggesting that self-attention is a more efficient operation for pooling context compared to depthwise separable convolutions.

Model	Block Groups
ResNet-50	[3, 4, 6, 3]
ResNet-101	[3, 4, 23, 3]
ResNet-152	[3, 8, 36, 3]
SENet-50	[3, 4, 6, 3]
SENet-101	[3, 4, 23, 3]
SENet-152	[3, 8, 36, 3]
SENet-350	[4, 40, 60, 12]
BoTNet-S1-59	[3, 4, 6, 6]
BoTNet-S1-77	[3, 4, 6, 12]
BoTNet-S1-110	[3, 4, 23, 6]
BoTNet-S1-128	[3, 4, 23, 12]

Table 13: Backbones and their block group configurations. All of them use the standard convolutional stem of the ResNet [29]. [3, 4, 6, 6], refers to the use of 3, 4, 6 and 6 blocks respectively in stages c₂, c₃, c₄, c₅. For BoTNet-S1 design, the final blockgroup c₅ uses a stride of 1 for all blocks. In order to reflect the improved training setting used in EfficientNets, the four BoTNet models (above) make use of SE blocks for blocks in the groups c₂, c₃, c₄.

4.8.6 ResNets and SENets are strong baselines until 83% top-1 accuracy

It is still not a widely established fact that ResNets and SENets (*without* additional tweaks such as ResNet-D [30], width scaling [71]) can be strong baselines for ImageNet. This was brought to attention through RegNets by Radavosovic et al. [50] but only in a standard training setting of 100 epochs. In Figure 7, we show that ResNets and SENets achieve strong performance in the improved EfficientNet training setting. They are strong enough that they can outperform all the EfficientNets from B0-B5 as revealed by the Pareto curve. *Particularly worth noting is the fact that ViT-384 (Vision Transformer finetuned to image size 384) is worse than a well tuned ResNet-50*, while the improved DeiT-384¹¹ is much worse than a SENet (in particular, S3 on the plot). Similarly, S2 is a strictly superior model compared to DeiT-224. BoTNets T3 and T4 *do not* outperform SENets, while T5 does perform better than S4. This suggests that pure convolutional models such as ResNets and SENets are still the best performing models until an accuracy regime of 83% top-1 accuracy. It is only beyond that accuracy regime where BoTNets begin to shine. We hope that by presenting these strong baselines, future research on attention models and backbone architectures take extra effort in comparing to strong convolutional baselines.

¹¹We report the accuracy of DeiT models without the distillation component for fair comparisons to baselines that do not use distillation.

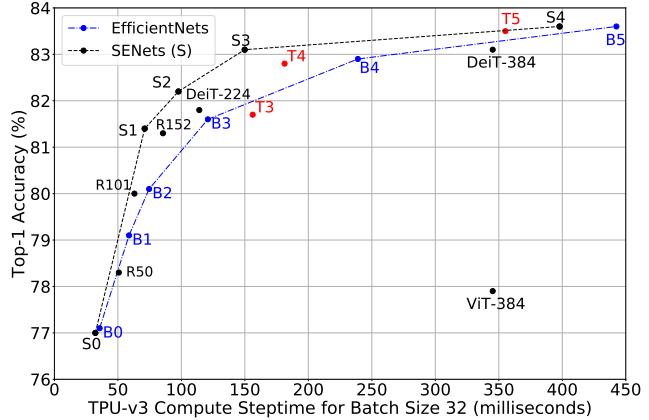


Figure 7: Presenting the results from Table 14 until the 83% accuracy regime (until EfficientNet-B5) in the form of a scatter plot. S0-S4 are SENets. T3, T4, T5 are BoTNets. B0-B5 are EfficientNets. ViT [14] and DeiT [63] refer to the recently popular Vision Transformer models, both the original and the improved versions.

4.8.7 BoTNets are the best performing models beyond 83% top-1 accuracy

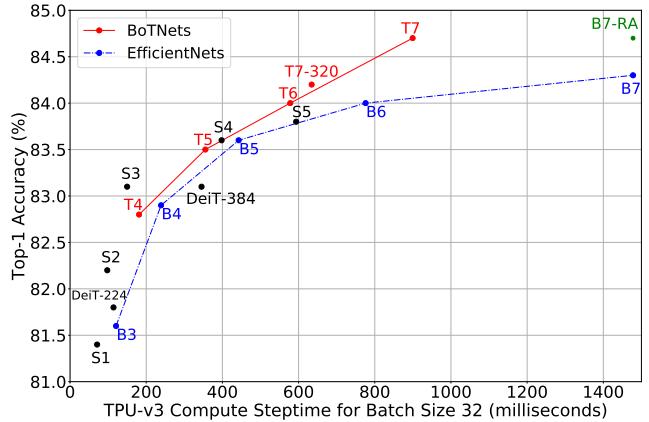


Figure 8: Presenting the results from Table 14 right until 84.7% accuracy regime in the form of a scatter plot. T4-T7 are BoTNets. S1-S5 are SENets. ViT and DeiT are Vision Transformers (baseline and improved version respectively). B7-RA is EfficientNet-B7 trained with RandAugment (instead of AutoAugment).

While SENets were a powerful model class clearly outperforming EfficientNets (up to B5) and BoTNets (up to T4), we found gains to diminish beyond SE-350 (350 layer SENet described in Table 13) trained with image size 384. This model is referred to as S5 and achieves 83.8% top-1 accuracy.¹² On

¹²While it is definitely possible to further improve this baseline through

Model	Backbone	Resolution	Top-1 Acc.	Top-5 Acc.	Params	M.Adds	Steptime
B0	EfficientNet	224	77.1	93.3	5.3M	0.39B	35.6
S0	SENet-50	160	77.0	93.5	28.02M	2.09B	32.3
B1	EfficientNet	240	79.1	94.4	7.8M	0.7B	58.8
-	ResNet-50	224	78.3	94.3	25.5M	4.09B	50.7
-	ResNet-50	256	78.8	94.5	25.5M	5.34B	62.1
-	SENet-50	224	79.4	94.6	28.02M	4.09B	64.3
B2	EfficientNet	260	80.1	94.9	9.2M	1.0B	74.6
-	ResNet-101	224	80.0	95.0	44.4M	7.8B	63.0
B3	EfficientNet	300	81.6	95.7	12M	1.8B	120.8
T3	BoTNet-S1-59	224	81.7	95.8	33.5M	7.3B	156.2
-	ResNet-152	224	81.3	95.5	60.04M	11.5B	85.6
S1	SENet-101	224	81.4	95.7	49.2M	7.8B	71.1
S2	SENet-152	224	82.2	95.9	66.6M	11.5B	97.7
B4	EfficientNet	380	82.9	96.4	19M	4.2B	238.9
T4	BoTNet-S1-110	224	82.8	96.3	54.7M	10.9B	181.3
S3	SENet-152	288	83.1	96.4	66.6M	19.04B	149.9
B5	EfficientNet	456	83.6	96.7	30M	9.9B	442.6
S4	SENet-350	320	83.6	96.6	115.18M	52.9B	397.9
T5	BoTNet-S1-128	256	83.5	96.5	75.1M	19.3B	355.2
B6	EfficientNet	528	84.0	96.8	43M	19B	776.3
S5	SENet-350	384	83.8	96.6	115.18M	52.9B	397.9
T6	BoTNet-S1-77	320	84.0	96.7	53.9M	23.3B	578.1
B7	EfficientNet	600	84.3	97.0	66M	37B	1478.4
T7-320	BoTNet-S1-128	320	84.2	96.9	75.1M	30.9B	634.3
B7-RA	EfficientNet	600	84.7	97.0	66M	37B	1478.4
T7	BoTNet-S1-128	384	84.7	97.0	75.1M	45.8B	804.5

Table 14: Various backbone architectures, evaluated under the fair setting of using the improved training elements from EfficientNet [59]. Models are grouped by the accuracy of each EfficientNet model. For ResNets, SENets and BoTNets, the best weight decay from $[2e-5, 4e-5, 8e-5, 1e-4]$ the best dropconnect from $[0.2, \text{None}]$, and the best RandAugment magnitude from $[10, 15, 24]$ are used to ensure fair comparisons to the already well tuned EfficientNets whose results we take from the latest update to the official codebase. The steptimes that we report, refer to the **compute** time on a TPU-v3 core, for a batch size of 32 for all the models. This is to ensure we do not use different batch sizes for different models and do not conflate TPU rematerialization for dense convolutions. Had we done so (i.e. picked bigger batch sizes for ResNets and BoTNets as long as they fit on memory), the speedup gains would be even higher. Further, by **compute** time, we just mean the time spent for forward and backward passes, and **not** the data loading. This is again to ensure comparisons across models do not exploit inefficient (or non-cached) data loading. B7-RA refers to EfficientNet-B7 trained with RandAugment [11].

the other hand, BoTNets scale well to larger image sizes (corroborating with our results in instance segmentation when the gains from self-attention were much more visible for larger images). In particular, T7-320 (T7 model trained with image size 320) achieves 84.2% top-1 acc., with a **2.33x**

ResNet-D and Wide ResNets, we believe those improvements are orthogonal and would apply to BoTNets as well, considering that BoTNets benefit from any improvement to the convolutional stack prior to the BoT blocks. We leave benchmarking such models for future work. Our focus is on the simplest and widely used baselines.

speedup over EfficientNet-B7 (trained with AutoAugment). T7 achieves 84.7% top-1 acc., matching the accuracy of B7-RA, with a **1.64x** speedup in efficiency. *Particularly worth noting is that BoTNets perform better than DeiT-384 (the improved ViT), showing the power of hybrid models that make use of both convolutions and self-attention compared to pure attention models such as DeiT.* While it is likely that pure attention models could outperform hybrid models such as BoTNets in future through architectural changes, we believe

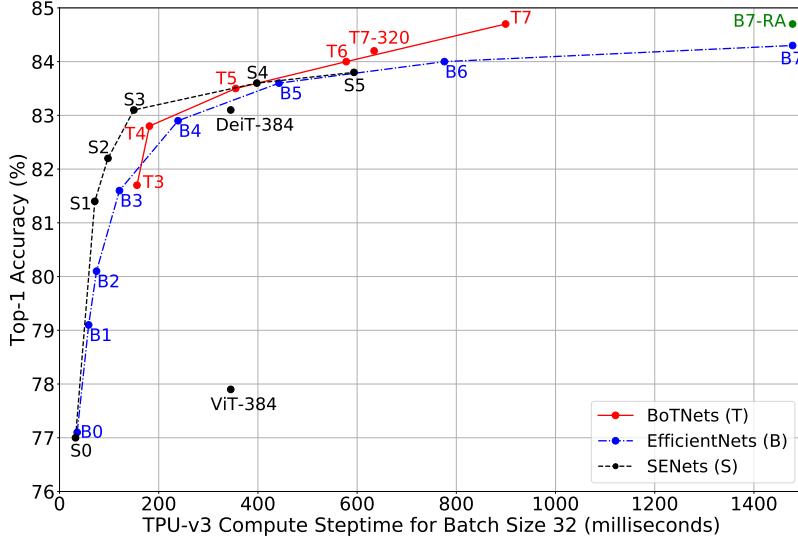


Figure 9: All backbones in Table 14 along with ViT and DeiT summarized in the form of scatter-plot and Pareto curves. SENets and BoTNets were trained while the accuracy of other models have been reported from corresponding papers.

the strong performance of the hybrid models makes them a strong baseline to keep in mind for future work.

4.9. Discussion

Figure 11 presents all the three model families (SENNets, EfficientNets and BoTNets) together in one plot. Please refer to Figures 7 and 8 for zoomed-in versions of the global figure. As discussed in the previous sections, SENNs are powerful models with strong performance that is better than EfficientNets and can be scaled all the way up to 83.8% top-1 accuracy without any bells and whistles such as ResNet-D, etc. BoTNets initially perform worse than SENNs (eg. T3) but begin to take over in terms of performance from T4 and strictly outperform EfficientNets, especially towards the end. EfficientNets do not scale well, particularly in the larger accuracy regime. This definitely suggests that the compound scaling rule in EfficientNets may not be strong enough and simple depth scaling (as in SENNs and BoTNets) works well enough. A more careful study of scaling rules for all these model classes is left for future research.

Recently, Transformer based models for visual recognition have gotten very popular since the Vision Transformer (ViT) model [14]. While our paper was developed concurrently to ViT, there has been a lot of follow-up to ViT such as DeiT [63] that have further improved the results of ViT. As seen in the results and emphasized already, DeiT-384 and ViT-384 are *clearly outperformed* by both SENNs and BoTNets currently. This suggests that convolutional and hybrid (convolution and self-attention) models are still strong models to beat as far as the ImageNet benchmark goes. The message in ViT was that pure attention models struggle in

the small data (ImageNet) regime¹³, but shine in the large data regime (JFT dataset), where inductive biases such as data augmentation and regularization tricks used in the EfficientNet training setting do not matter. Nevertheless, we think it is an interesting exercise to explore hybrid models such as BoTNet even in the large data regime simply because they seem to scale much better than SENNs and EfficientNets (Figure 11) and achieve better performance than DeiT on ImageNet. We leave such a large scale effort for future work. It is also possible to design pure attention models that scale better than ViT in the large data regime, example, those employing local attention such as SASA [51].

It is unclear what the right model class is, given that we have not yet explored the limits of hybrid and local attention models in the large data regimes, and that the pure attention models currently lag behind both convolutional and hybrid models in the small data regime. Nevertheless, with the hope that the ImageNet benchmark has been representative of the best performing models in the vision community, BoTNets are likely to be a simple and compelling baseline to always consider. While they maybe viewed as Hybrid-ViT models for simplicity, the differences have been highlighted in Section 2 to how BoT blocks are different from the regular Transformer blocks, and thereby, different from ViT.

5. Conclusion

The design of vision backbone architectures that use self-attention is an exciting topic. We hope that our work helps in improving the understanding of architecture design in this

¹³ImageNet may not be a small dataset by conventional standards, but is referred to as such here for the contrast with JFT.

space. As noted, combining our backbone with DETR [8] and studying the effect on small objects is an interesting avenue for future work. Incorporating self-attention for other computer vision tasks such as keypoint detection [7] and 3D shape prediction[23]; studying self-attention architectures for self-supervised learning in computer vision [31, 26, 62, 9, 24, 10]; and scaling to much larger datasets such as JFT, YFCC and Instagram, are ripe avenues for future research.

6. Acknowledgements

The authors thank Ilija Radosavovic for several useful discussions regarding backbone architectures and Detectron2 codebase; Pengchong Jin and Xianzhi Du for help with the TF Detection codebase; Barret Zoph, Irwan Bello and Ekin Dogus Cubuk for valuable discussions. We thank Zak Stone for extensive compute support throughout this project the through TFRC program providing Google Cloud TPUs (<https://www.tensorflow.org/tfrc>).

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 3
- [2] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3286–3295, 2019. 4, 5
- [3] 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. 1
- [4] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 2, pages 60–65. IEEE, 2005. 3
- [5] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. 2, 7
- [6] Yue Cao, Jiarui Xu, Stephen Lin, Fangyun Wei, and Han Hu. Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3
- [7] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017. 13
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. 2, 3, 6, 13, 18
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 13
- [10] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020. 13
- [11] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020. 11
- [12] Jifeng Dai, Kaiming He, and Jian Sun. Instance-aware semantic segmentation via multi-task network cascades. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3150–3158, 2016. 1
- [13] 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. 1
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020. 2, 3, 4, 8, 10, 12
- [15] Xianzhi Du, Tsung-Yi Lin, Pengchong Jin, Golnaz Ghiasi, Mingxing Tan, Yin Cui, Quoc V Le, and Xiaodan Song. Spinenet: Learning scale-permuted backbone for recognition and localization. *arXiv preprint arXiv:1912.05027*, 2019. 5, 7
- [16] Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018. 9
- [17] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 1
- [18] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D. Cubuk, Quoc V. Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation, 2020. 5
- [19] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7036–7045, 2019. 2, 7
- [20] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 1, 3
- [21] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 1, 3
- [22] Ross Girshick, Ilija Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron, 2018. 3, 5
- [23] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9785–9795, 2019. 13

- [24] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020. 13
- [25] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision*, pages 297–312. Springer, 2014. 1
- [26] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning, 2019. 5, 13
- [27] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 18
- [28] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 1, 2, 3, 5
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2, 3, 10
- [30] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2019. 10
- [31] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019. 13
- [32] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 9
- [33] Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional transformers. *arXiv preprint arXiv:1912.12180*, 2019. 2
- [34] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3588–3597, 2018. 1
- [35] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3464–3473, 2019. 2
- [36] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 9
- [37] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 603–612, 2019. 2
- [38] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 3, 5
- [39] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 3
- [40] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 1
- [41] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 5
- [42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 2, 5
- [43] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8759–8768, 2018. 2, 7
- [44] 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. 1
- [45] Yudong Liu, Yongtao Wang, Siwei Wang, TingTing Liang, Qijie Zhao, Zhi Tang, and Haibin Ling. Cbnet: A novel composite backbone network architecture for object detection. *arXiv preprint arXiv:1909.03625*, 2019. 2
- [46] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23, 2019. 2
- [47] Chao Peng, Tete Xiao, Zeming Li, Yuning Jiang, Xiangyu Zhang, Kai Jia, Gang Yu, and Jian Sun. Megdet: A large mini-batch object detector. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6181–6189, 2018. 5
- [48] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9, 2019. 1
- [49] Ilija Radosavovic, Justin Johnson, Saining Xie, Wan-Yen Lo, and Piotr Dollár. On network design spaces for visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1882–1890, 2019. 3
- [50] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10428–10436, 2020. 3, 8, 10
- [51] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Standalone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019. 1, 2, 4, 5, 12
- [52] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017. 9

- [53] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1, 3
- [54] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 1, 3, 8
- [55] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018. 4, 5, 6
- [56] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1
- [57] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. 9
- [58] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7464–7473, 2019. 2
- [59] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019. 1, 2, 9, 11, 20
- [60] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. *arXiv preprint arXiv:1911.09070*, 2019. 2, 7
- [61] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *arXiv preprint arXiv:2009.06732*, 2020. 2
- [62] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. 13
- [63] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers and distillation through attention, 2021. 10, 12
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 1, 2, 3, 4, 5, 20
- [65] Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In *International conference on machine learning*, pages 1058–1066, 2013. 9
- [66] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. *arXiv preprint arXiv:2003.07853*, 2020. 1, 2
- [67] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 2, 3, 4, 8, 19
- [68] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018. 5
- [69] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 501–509, 2019. 4, 8
- [70] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 1
- [71] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 10
- [72] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Zhi Zhang, Haibin Lin, Yue Sun, Tong He, Jonas Mueller, R. Manmatha, Mu Li, and Alexander Smola. Resnest: Split-attention networks, 2020. 1, 7
- [73] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition, 2020. 1, 2, 4, 5



Figure 10: Qualitative Results from Mask R-CNN with BoTNet-50. Images are from the COCO dataset.



Figure 11: Qualitative comparison between ResNet-50 and BoTNet-50 Mask R-CNN on images from the COCO dataset. While not very different, BoTNet is able to create more precise localizations as well as detect small objects better. A couple of examples are detecting the shadow of the kite, or the spoon next to the fork, etc.

A. Appendix

A.1. Code for the BoT block

We provide the exact code used for implementing the Multi-Head Self-Attention (MHSA) in 2D with relative-position encodings as well as the implementation of BoT block in this gist link: <https://gist.github.com/aravindsrinivas/56359b79f0ce4449bcb04ab4b56a57a2>.

A.2. Implementation: COCO Instance Segmentation and Object Detection

Our implementation is based on the Cloud TPU Detection codebase - <https://github.com/tensorflow/tpu/tree/master/models/official/detection>. Our canonical setting uses the following hyperparameters which are updated in configs/maskrcnn_config.py:

- output_size of 1024×1024 .
- aug_scale_min=0.8, aug_scale_max=1.25
- mrcnn_head:num_convs=4, num_filters=256, mask_target_size=28
- frcnn_head:num_convs=4, num_filters=256, fc_dims=1024, num_fcs=1
- rpn_head:min_level=2, max_level=6, anchors_per_location=3, num_convs=2, num_filters=256
- fpn:min_level=2, max_level=6
- anchor: num_scales=1, anchor_size=8, min_level=2, max_level=6

For all experiments, we use L2 weight decay of $4e-5$, sync-batch-norm with momentum 0.997 and epsilon 1e-4. We use batch norm in the backbone, FPN, RPN head, FRCNN head and MRCNN head. We initialize backbone weights with pre-trained ImageNet checkpoints and fine-tune all the weights (including the batch-norm parameters) as specified in MoCo [27].

Table 15 presents the hyperparameters for models that we train with batch size 64 on 8 TPU-v3 cores (equivalent to DGX-1) which applies to most of our models, in particular, all the models that we train with image size 1024×1024 .

For models that do not fit with batch size of 8 per core (for example, the ones that train with 1280×1280), we train with a global batch size of 128 on 32 TPU-v3 cores (4 images per core). The 6x schedule for these models that train with 32 cores corresponds to training with 67.5k steps (since batch size 128 is double 64), with [60k, 65k] learning

Type	Epochs	Train Steps	LR Schedule
1x	12	22.5k	[15k, 17.5k]
2x	24	45k	[37.5k, 40k]
3x	36	67.5k	[60k, 65k]
6x	72	135k	[120k, 130k]

Table 15: Learning Rate Schedules for the 1x, 2x, 3x and 6x settings for models trained with global batch size of 64, 8 TPU-v3 cores, 16 GB HBM per cores, learning rate 0.1 with schedule [0.01, 0.001]. The learning rate is initially warmed up for 500 warmup steps from 0.0067.

rate schedule. The learning rate is 0.15 with schedule [0.015, 0.0015] with 500 warmup steps started from 0.0067.

For the model that achieves best results (44.4% mask AP) with ResNet-200 backbone and BoTNet structure for blockgroup c5, we use 8 convolutions in the MRCC head.

For our object detection experiments, we just turned off the mask branch by setting the include_mask flag as False and eval_type as box instead of box_and_mask.

A.3. BoTNet improves over ResNet for Object Detection

Does BoTNet improve upon ResNet for the task of object detection as well? We verify this through experiments with the Faster R-CNN framework (Table 16). We observe that BoTNet indeed improves over ResNet significantly. BoT R50 improves upon R50 by a significant margin of + 1.6% AP_S^{bb} (AP_S^{bb} for small objects). These results suggest that *self-attention has a big effect in detecting small objects* which is considered to be an important and hard problem for deploying object detection systems in the real world. This result is in contrast with DETR [8] which observes a big improvement on large objects but not on small objects. We believe that introducing self-attention in the backbone architecture might help fix the lack of gains on small objects in DETR. Finally, we study if larger images would further benefit BoTNet for object detection. Using a multi-scale jitter of [0.1, 2.0] with 72 epoch training and image size of 1280×1280 , we see that BoT R152 achieves a strong performance of **48.4%** AP^{bb}.

A.4. BoT block with stride

The first block in the final block group c5 of a ResNet runs the spatial 3×3 convolution with a stride 2 on a resolution that is 2x the height and width of the other two blocks in the group. Unlike spatial convolutional layers, the MHSA layers in BoT blocks do not implement striding. In fact, implementing strided self-attention (both local and global) is an engineering challenge which we leave for future work. Our goal is to only use existing primitives. So, we adopt

Backbone	AP^{bb}	AP_S^{bb}	AP_L^{bb}
R50	41.5	24.9	54.3
BoT50	43.1 (+ 1.6)	27.6 (+ 2.7)	55.7 (+ 1.3)
R101	42.4	24.5	55.0
BoT101	44.3 (+ 1.9)	26.5 (+ 2.0)	57.5 (+ 2.5)
R152	45.1	30.1	56.6
BoT152	48.4 (+ 3.3)	32.7 (+ 2.6)	60.6 (+ 4.0)

Table 16: BoTNet for Faster R-CNN: The first four rows are trained for 36 epochs, jitter [0.8, 1.25], image size 1024×1024 . Final two rows are trained for 72 epochs with 1280×1280 image size, jitter [0.1, 2.0] using v3-32 Cloud TPU.

a very simple fix of using a local 2×2 Avg. Pooling with stride of 2 for implementing the spatial downsampling in the first BoT block. As noted in our ablation on placement of the attention (Section 4.4), we see that the self-attention in the first block is very useful for gain on small objects. At the same time, it involves running self-attention on a resolution of 64×64 , equivalent to 4096 tokens and 4096×4096 query-key matrix. We believe a well optimized strided attention implementation will drastically make this more efficient in future. The block is explained in Figure 12.

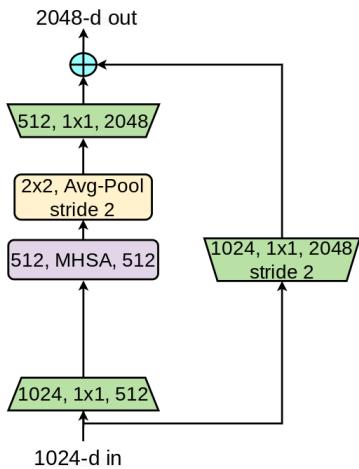


Figure 12: BoT block with stride=2, implemented in the first block of the c_5 blockgroup in a BoT-ResNet. Since the incoming tensor has 1024 channels and has a stride of 16 with respect to input resolution, there is a projection shortcut with a strided 1×1 convolution. The self-attention operation in the MHSA is global and maintains the resolution. Therefore, we use local 2×2 Avg. Pooling with a stride 2 on top of it.

A.5. Non-Local Comparison

Section 4.7 in the main paper offered a comparison between Non-Local Nets and BoTNets as well as an implemen-

tation of BoT blocks in the design of Non-Local Nets. In order to make it more clear visually, we provide an illustration of all designs in Figure 13. The figure clearly highlights the differences between inserting blocks vs replacing blocks. We note that the NL style insertion requires careful placement as prescribed in the original paper. On the other hand, simply replacing the final three blocks is a convenient design choice. Nevertheless, BoT blocks are still an improvement over vanilla NL blocks even in the insertion design, likely due to using multiple heads, relative position encodings, and value projection. As already explained in Section 4.7 in the main paper, even the *replacement* design of BotNet performs better than the *insertion* design of vanilla NL, likely due to performing more self-attention (three vs one). Adding value projections to the NL block only improved the performance by 0.2 AP^{bb}, while the gains from using 4 heads in the MHSAs in BoT block only helps by $\tilde{0.2}AP^{mk}$. Therefore, BoT blocks can be viewed as an improved implementation of NL blocks with relative position encodings being the main driver of the performance difference.

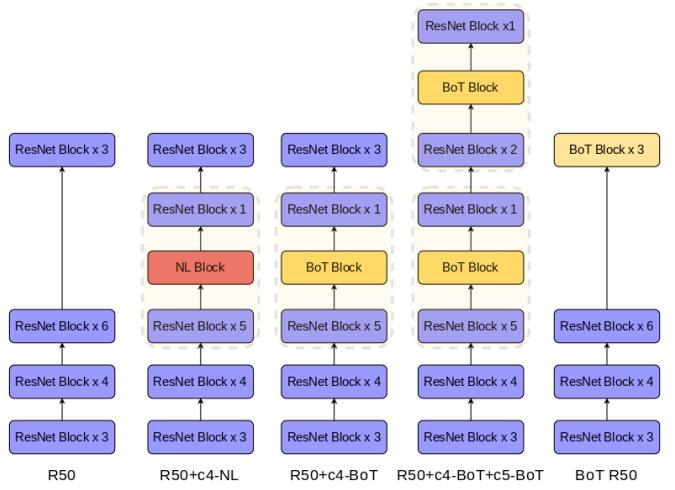


Figure 13: **First:** Regular ResNet 50 with $[3, 4, 6, 3]$ blockgroup structure. **Second:** Non-Local (NL) block *inserted* in the c_4 blockgroup of a ResNet, between the pre-final and final ResNet block, as specified by Wang et al. [67]. **Third:** BoT block *inserted* in the same manner as a NL block with the differences between a BoT and NL block highlighted in Figure 4 in the main paper. **Fourth:** 2 BoT blocks, one each in c_4, c_5 blockgroups *inserted* in the same manner (between pre-final and final block) as prescribed by Wang et al. [67]. **Fifth:** BoT50, where the final three ResNet blocks are *replaced* by BoT blocks.

A.6. Comparison to Squeeze-Excite

One may argue that squeeze-excitation blocks and self-attention blocks are very much related in the sense that both

of them perform global aggregation and provide that as a context to convolutional models. Therefore, it is natural to ask for a comparison between the two blocks especially given that Squeeze-Excite (SE) blocks are computationally cheap and easy to use compared to BoT blocks. Table 17 presents this comparison. For fair comparison to BoTNet, we only place SE blocks in c5 blockgroup and call this setup as R50+ c5-SE. We do not see any difference between R50 and R50 + c5-SE. However, we note that it is possible to get visible gains on top of R50 when placing SE blocks in all bottleneck blocks throughout the ResNet and not just in c5. Such a change (using SE in c2, c3, c4) is orthogonal to BoTNet and can be combined with BoTNet by using BoT blocks in c5 and SE blocks in c2, c3, c4. We leave exploration of such architectures for future work.

Backbone	AP ^{bb}	AP ^{mk}
R50	42.1	37.7
BoT50	43.6 (+ 1.5)	38.9 (+ 1.2)
R50 + c5-SE	42.1	37.6 (- 0.1)

Table 17: Comparing R50, BoT50 and R50 + c5-SE; all 3 setups using the canonical training schedule of 36 epochs, 1024×1024 images, multi-scale jitter [0.8, 1.25].

A.7. ImageNet Test-Set Accuracy

As has been the convention since 2017, our results in the paper only report and compare the validation set (50K images) accuracy on the ImageNet benchmark. However, the EfficientNet paper [59] presented updated results on the ImageNet test-set (100K images submitted on <http://image-net.org>). We also provide the results on the ImageNet test set for the ablation setup in Table 12, to verify that there are not any surprising differences between the validation and test set (Table 18 presents similar numbers to Table 12).

Backbone	top-1 acc.
SE50	79.3
BoT-S1-50	80.3 (+ 1.0)

Table 18: ImageNet *test set* results in a further improved training setting with SE blocks and SiLU non-linearity: 200 epochs, batch size 4096, weight decay 4e-5, RandAugment (2 layers, magnitude 10), and label smoothing of 0.1. R50 with SE blocks is referred to as SE50.

A.8. Resolution Dependency in BoTNet

Architectures such as BoTNet and ViT that use self-attention end up adopting position encodings [64]. This in turn creates a dependency at inference time to use the

same resolution that the model was trained for. For example, taking T7 trained at 384×384 , and running inference with a different resolution (say 512×512) would introduce additional positional encodings. Purely convolutional architectures such as ResNets and SENets do not face this problem. We leave it for future work to investigate various design choices in making Transformer based architectures for vision, resolution independent at inference time. Some potential ideas are multi-resolution training with bilinearly interpolated positional encodings, using a spatial convolution before every self-attention operation (which could be performed without any positional encoding), etc.

While BoTNet-like hybrid architectures benefit from the spatial dependencies implicitly learned by the convolutional stack prior to attention, the reason we still use position encodings is because of the improvements that arise from using them (Table 4). We also observe similar performance gains on the ImageNet benchmark (Table 19). Using no position encoding at all for the BoT blocks still provides a gain of + 0.7% top-1 accuracy, but lags behind the gains from using position encodings (+ 1.2%). Further, relative position encoding is not as important in ImageNet benchmark as it is in the COCO benchmark (Table 4). Absolute position encodings provide similar gains to relative position encodings. We believe this is likely due to the nature of the task being less contextual (and doesn't involve precise localization) for image classification compared to detection and segmentation. We retain relative position encodings for the ImageNet experiments for consistency with the architecture used for COCO. However, for practitioners, we recommend using absolute position encodings for image classification since it is faster (for training and inference) and simpler to implement.

Backbone	Pos-Enc	top-1 acc.
SE50	N/A	79.2
BoT-S1-50	-	79.9 (+ 0.7)
BoT-S1-50	Abs	80.2 (+ 1.0)
BoT-S1-50	Rel	80.4 (+ 1.2)

Table 19: ImageNet (val) results in a further improved training setting with SE blocks and SiLU non-linearity: 200 epochs, batch size 4096, weight decay 4e-5, RandAugment (2 layers, magnitude 10), and label smoothing of 0.1. R50 with SE blocks is referred to as SE50.

A.9. M.Adds and Params

In-built tools for computing parameters and FLOPs ($2 * M.\text{Adds}$) are often not accurate when used on architectures with einsums. We therefore explicitly calculated the M.Adds we report in this paper using this script <https://gist.github.com/aravindsrinivas/e8a9e33425e10ed0c69c1bf726b81495>.