# Co-Scale Conv-Attentional Image Transformers

Weijian Xu[*],    Yifan Xu[*],    Tyler Chang,    Zhuowen Tu
University of California San Diego
{wex041, yix081, tachang, ztu}@ucsd.edu

## Abstract

*In this paper, we present Co-scale conv-attentional image Transformers (CoaT), a Transformer-based image classifier equipped with co-scale and conv-attentional mechanisms. First, the co-scale mechanism maintains the integrity of Transformers' encoder branches at individual scales, while allowing representations learned at different scales to effectively communicate with each other; we design a series of serial and parallel blocks to realize the co-scale attention mechanism. Second, we devise a conv-attentional mechanism by realizing a relative position embedding formulation in the factorized attention module with an efficient convolution-like implementation. CoaT empowers image Transformers with enriched multi-scale and contextual modeling capabilities. On ImageNet, relatively small CoaT models attain superior classification results compared with the similar-sized convolutional neural networks and image/vision Transformers. The effectiveness of CoaT's backbone is also illustrated on object detection and instance segmentation, demonstrating its applicability to the downstream computer vision tasks.*

## 1. Introduction

A notable recent development in artificial intelligence is the creation of attention mechanisms [44] and Transformers [36], which have made a profound impact in a range of fields including natural language processing [6, 23], document analysis [45], speech recognition [7], and computer vision [8, 2]. In the past, state-of-the-art image classifiers have been built primarily on convolutional neural networks (CNNs) [16, 15, 30, 29, 10, 42] that operate on layers of filtering processes. Recent developments [34, 8] however begin to show encouraging results for Transformer-based image classifiers.

In essence, both the convolution [16] and attention [44] operations address the fundamental representation problem for structured data (e.g. images and text) by modeling the

---

[*] indicates equal contribution.
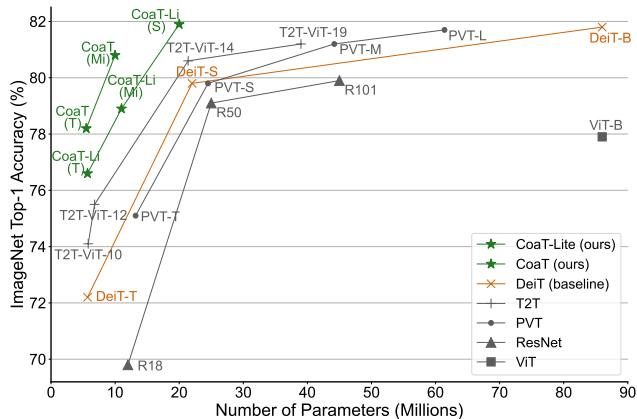    Code at `https://github.com/mlpc-ucsd/CoaT`.



Figure 1. **Model Size vs. ImageNet Accuracy.** Our CoaT model significantly outperforms other image Transformers. In particular, CoaT achieves competitive accuracies (78.2% and 80.8% top-1) compared to EfficientNet B0 and B2 with similar model size. Details are in Table 2.

local contents, as well as the contexts. The receptive fields in CNNs are gradually expanded through a series of convolution operations. The attention mechanism [44, 36] is, however, different from the convolution operations: (1) the receptive field at each location or token in self-attention [36] readily covers the entire input space since each token is "matched" with all tokens including itself; (2) the self-attention operation for each pair of tokens computes a dot product between the "query" (the token in consideration) and the "key" (the token being matched with) to weight the "value" (of the token being matched with). Moreover, although the convolution and the self-attention operations both perform a weighted sum, their weights are computed differently: in CNNs, the weights are learned during training but fixed (or gated [5]) during testing; in the self-attention mechanism, the weights are dynamically computed based on the similarity or affinity between every pair of tokens. As a consequence, the self-similarity operation in the self-attention mechanism provides modeling means that are potentially more adaptive and general than convolution operations. In addition, the introduction of position encodings and embeddings [36] provides Transformers with additional flexibility

1

to model spatial configurations beyond fixed input structures.

Of course, the advantages of the attention mechanism are not given for free, since the self-attention operation computes an affinity/similarity that is more computationally demanding than linear filtering in convolution. The early development of Transformers has been mainly focused on natural language processing tasks [36, 6, 23] since text is "shorter" than an image, and text is easier to tokenize. In computer vision, self-attention has been adopted to provide added modeling capability for various applications [39, 43, 50]. With the underlying framework increasingly developed [8, 34], Transformers start to bear fruit in computer vision [2, 8] by demonstrating their enriched modeling capabilities.

In the seminal DEtection TRansformer (DETR) [2] algorithm, Transformers are adopted to perform object detection and panoptic segmentation, but DETR still uses CNN backbones to extract the basic image features. Efforts have recently been made to build image classifiers from scratch, all based on Transformers [8, 34, 38]. While Transformer-based image classifiers have reported encouraging results, performance and design gaps to the well-developed CNN models still exist. For example, in [8, 34], an input image is divided into a single grid of fixed patch size. In this paper, we develop Co-scale conv-attentional image Transformers (CoaT) by introducing two mechanisms of practical significance to Transformer-based image classifiers. The contributions of our work are summarized as follows:

- We introduce a co-scale mechanism to image Transformers by maintaining encoder branches at separate scales while engaging cross-layer attention. Two types of co-scale blocks are developed, namely a serial and a parallel block, realizing **fine-to-coarse**, **coarse-to-fine**, and **cross-scale attention** image modeling.

- We design a conv-attention module to realize **relative position embeddings** in the **factorized attention** module using a convolution-like attention operation that achieves significantly enhanced computation efficiency when compared with vanilla self-attention layers in Transformers.

Our resulting Co-scale conv-attentional image Transformers (CoaT) learn effective representations under a modularized architecture. On ImageNet, CoaT achieves state-of-the-art classification results when compared with the competitive convolutional neural networks (e.g. Efficient-Net [33]), while significantly outperforming the competing Transformer-based image classifiers by a large margin [8, 34, 38].

## 2. Related Works

Our work is inspired by the recent efforts [8, 34] to realize Transformer-based image classifiers. ViT [8] demonstrates the feasibility of building Transformer-based image classifiers from scratch, but its performance on ImageNet [26] is not achieved without including additional training data; DeiT [34] attains results comparable to convolution-based classifiers by using an effective training strategy together with model distillation, removing the data requirement in [8]. Both ViT [8] and DeiT [34] are however based on a single image grid of fixed patch size.

The development of our co-scale conv-attentional transformers (CoaT) is motivated by two observations: (1) multi-scale modeling typically brings enhanced capability to representation learning [10, 25, 37]; (2) the intrinsic connection between relative position encoding and convolution makes it possible to carry out efficient self-attention using conv-like operations. As a consequence, the superior performance of the CoaT classifier shown in the experiments comes from two of our new designs in Transformers: (1) a co-scale mechanism that allows cross-layer attention; (2) a conv-attention module to realize an efficient self-attention operation. Next, we highlight the differences of the two proposed modules with the standard operations and concepts.

- **Co-Scale vs. Multi-Scale**. Multi-scale approaches have a long history in computer vision [40, 21]. Convolutional neural networks [16, 15, 10] naturally implement a fine-to-coarse strategy. U-Net [25] enforces an extra coarse-to-fine route in addition to the standard fine-to-coarse path; HRNet [37] provides a further enhanced modeling capability by keeping simultaneous fine and coarse scales throughout the convolution layers. In a parallel development [38] to ours, layers of different scales are in tandem for the image Transformers but [38] merely carries out a fine-to-coarse strategy. The co-scale mechanism proposed here differs from the existing methods in how the responses are computed and interact with each other: CoaT consists of a series of highly modularized serial and parallel blocks to enable fine-to-coarse, coarse-to-fine, and cross-scale attention on tokenized representations. The joint attention mechanism across different scales in our co-scale module provides enhanced modeling power beyond the standard linear fusion in existing multi-scale approaches [10, 25, 37].

- **Conv-Attention vs. Attention.** Pure attention-based models [24, 13, 50, 8, 34] have been introduced to the vision domain. [24, 13, 50] replace convolutions in ResNet-like architecture by self-attention modules for better local and non-local relation modeling. In contrast, [8, 34] directly adapt the Transformer [36] for image recognition. Recently, there have been works [1, 4] enhancing the attention mechanism by introducing convolution. LambdaNets [1] introduces an efficient self-attention alternative for global context modeling and employs a 3-D convolution to realize the relative posi-

tion embeddings in local context modeling. CPVT [4] designs 2-D depthwise convolutions as the conditional positional encoding after self-attention. In our conv-attention: (1) we adopt an efficient factorized attention following [1]; (2) we design a depthwise convolution-based relative position encoding, and (3) extend it to be an alternative case in convolutional position encoding, related to CPVT [4]. Detailed discussion of our network design and its relation with LambdaNets [1] and CPVT [4] can be found in Section 4.1 and 4.2.

## 3. Revisit Scaled Dot-Product Attention

Transformers take as input a sequence of vector representations (i.e. tokens) $\mathbf{x}_1, ..., \mathbf{x}_N$, or equivalently $X \in \mathbb{R}^{N \times C}$. The self-attention mechanism in Transformers projects each $\mathbf{x}_i$ into corresponding query, key, and value vectors, using learned linear transformations $W^Q, W^K$, and $W^V \in \mathbb{R}^{C \times C}$. Thus, the projection of the whole sequence generates representations $Q, K, V \in \mathbb{R}^{N \times C}$. In the *scaled dot-product attention* from original Transformers [36]:

$$\text{Att}(X) = \text{softmax}\left(\frac{QK^\top}{\sqrt{C}}\right)V \qquad (1)$$

The softmax outputs can be seen as an attention map from queries to keys. Note that we hide the batch size $B$ and the number of attention heads $H$ in the actual multi-head self-attention [36] for simplified derivation. Self-attention outputs are then passed through a linear layer and a simple feed-forward network, with residual connections and layer-norm applied before and after the feed-forward network.

In vision transformers [8, 34], the input sequence of vectors is formulated by the concatenation of a class token CLS and the flattened feature vectors $\mathbf{x}_1, ..., \mathbf{x}_{HW}$ as image tokens from the feature maps $F \in \mathbb{R}^{H \times W \times C}$, for a total length of $N = HW + 1$. Due to the high dimensions of natural images in pixels (i.e. $N \gg C$), we cannot afford the computation of applying the self-attention with high resolution because the softmax logits in Equation 1 have $O(N^2)$ space complexity and $O(N^2C)$ time complexity for the whole sequence. To tackle this computation issue, [8, 34] tokenize the image by patches instead of pixels to reduce the length of sequence. However, the coarse splitting (e.g. 16×16 patches) limits the capability of the model to represent details within each patch, which is essential in many vision tasks. To address this dilemma between computation and representation capability, we propose a *co-scale* mechanism that can enable interaction between different scales to produce a rich representation, with the help of an efficient *conv-attentional* mechanism that lowers the computation complexity for high-resolution feature maps.
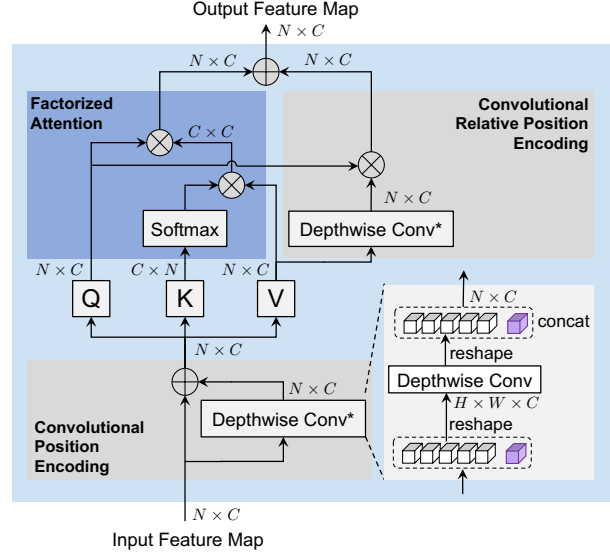


Figure 2. **Illustration of the conv-attentional module.** We apply a convolutional position encoding on the image tokens from the input. The resulting features are fed into a factorized attention with a convolutional relative position encoding.
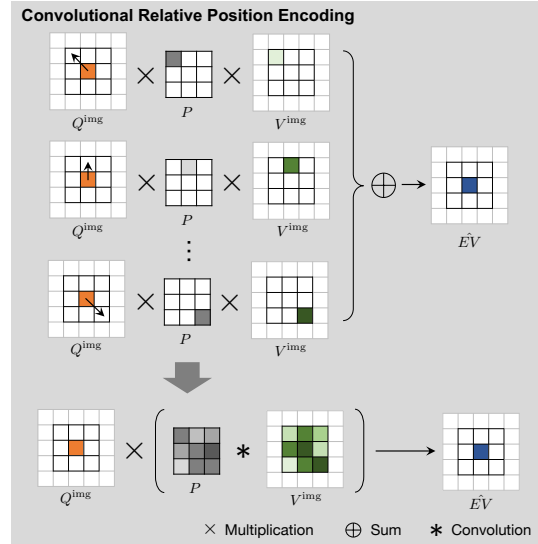


Figure 3. **Illustration of the convolutional relative position encoding.** (Upper) Each query (orange), value (green) and the corresponding relative position encoding (gray) are multiplied together, where the sum of the resulting products leads to the output (blue) in $\hat{E}V$. (Lower) The equivalent form to generate the output, which is the product of the query and the convolution between the position encoding map and value tokens.

## 4. Conv-Attention Module

### 4.1. Factorized Attention Mechanism

In Equation 1, the materialization of the softmax logits and attention maps leads to the $O(N^2)$ space complexity
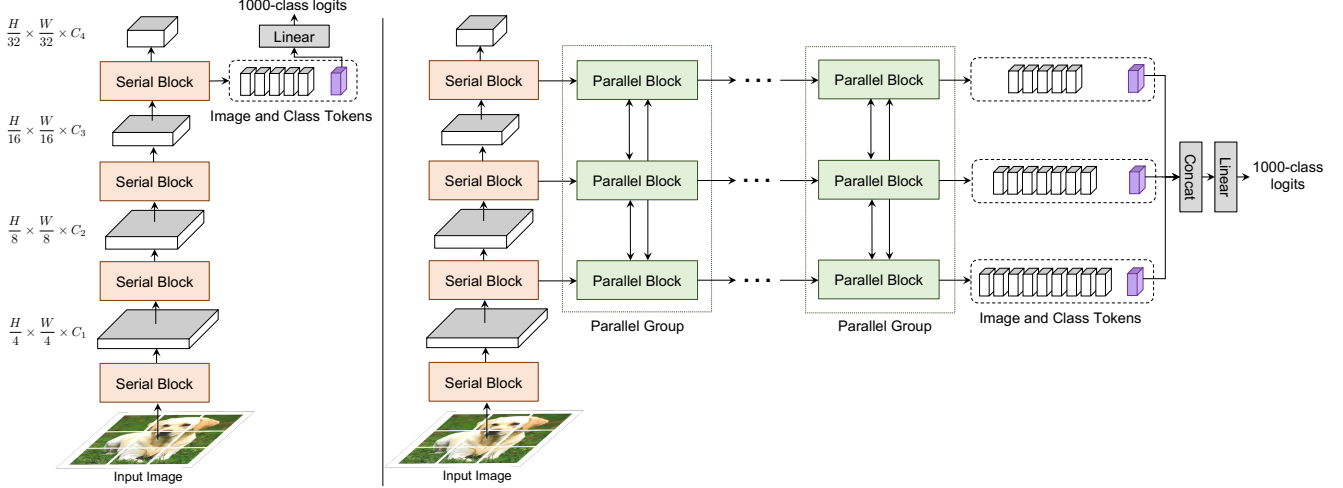
Figure 4. **CoaT model architecture.** (Left) The overall network architecture of **CoaT-Lite**. CoaT-Lite consists of serial blocks only, where image features are down-sampled and processed in a sequential order. (Right) The overall network architecture of **CoaT**. CoaT consists of serial blocks and parallel blocks that interact via co-scale attention.

and $O(N^2C)$ time complexity. Inspired by recent works [3, 28, 1] on linearization of self-attention, we approximate the softmax attention map by factorizing it using two functions $\phi(\cdot), \psi(\cdot) : \mathbb{R}^{N \times C'} \to \mathbb{R}^{N \times C'}$ and compute the second matrix multiplication (keys and values) together:

$$\text{FactorAtt}(X) = \phi(Q)\Big(\psi(K)^\top V\Big) \qquad (2)$$

The factorization leads to a $O(NC')$ space complexity and $O(NC'^2)$ time complexity, where both are linear functions of the sequence length $N$. Performer [3] uses random projections in $\phi$ and $\psi$ for a provable approximation, but with the cost of relatively large $C'$. Efficient-Attention [28] applies the softmax function for both $\phi$ and $\psi$, which is efficient but causes a significant performance drop on the vision tasks in our experiments. Here, we develop our factorized attention mechanism following LambdaNets [1] with $\phi$ as the identity function and $\psi$ as the softmax:

$$\text{FactorAtt}(X) = \frac{Q}{\sqrt{C}}\Big(\text{softmax}(K)^\top V\Big) \qquad (3)$$

where $\text{softmax}(\cdot)$ is applied across the tokens in the sequence in an element-wise manner and the projected channels $C' = C$. Different from LambdaNets [1], we also add the scaling factor $1/\sqrt{C}$ back in due to its normalizing effect, bringing better performance. This factorized attention takes $O(NC)$ space complexity and $O(NC^2)$ time complexity. It is noteworthy that the proposed factorized attention following [1] is not a direct approximation of the scaled dot-product attention, but it can still be regarded as a generalized attention mechanism modeling the feature interactions using query, key and value vectors.

## 4.2. Convolution as Position Encoding

Our factorized attention module mitigates the computational burden from the original scaled dot-product attention. However, because we compute $L = \text{softmax}(K)^\top V \in \mathbb{R}^{C \times C}$ first, $L$ can be seen as a global data-dependent linear transformation for every feature vector in the query map $Q$. This indicates if we have two query vectors $\mathbf{q}_1, \mathbf{q}_2 \in \mathbb{R}^C$ from $Q$ and $\mathbf{q}_1 = \mathbf{q}_2$, then their corresponding self-attention outputs will be the same:

$$\text{FactorAtt}(X)_1 = \frac{\mathbf{q}_1}{\sqrt{C}}L = \frac{\mathbf{q}_2}{\sqrt{C}}L = \text{FactorAtt}(X)_2 \quad (4)$$

Without the position encoding, the Transformer is only composed of linear layers and self-attention modules. Thus, the output of a token is dependent on the corresponding input without awareness of any difference in its locally nearby features. This property is unfavorable for vision tasks such as semantic segmentation (e.g. the same blue patches in the sky and the sea are segmented as the same category).

**Convolutional Relative Position Encoding.** To enable vision tasks, ViT and DeiT [8, 34] insert absolute position embeddings into the input, which may have limitations in modeling the relative relations between local tokens. Instead, following [27], we can integrate a relative position encoding $P = \{\mathbf{p}_i, i = -\frac{M-1}{2}, ..., \frac{M-1}{2}\}$ with window size $M$ to obtain the relative attention map $EV \in \mathbb{R}^{N \times C}$ in attention formulation if tokens are regarded as a 1-D sequence:

$$\text{RelFactorAtt}(X) = \frac{Q}{\sqrt{C}}\Big(\text{softmax}(K)^\top V\Big) + EV \quad (5)$$

where the encoding matrix $E \in \mathbb{R}^{N \times N}$ has elements:

$$E_{ij} = \mathbb{1}(i,j)\mathbf{q}_i \cdot \mathbf{p}_{j-i}, \ \ 1 \le i, j \le N \qquad (6)$$

4

in which $\mathbb{1}(i,j) = \mathbb{1}_{\{|j-i|\leq(M-1)/2\}}(i,j)$ is an indicator function. Each element $E_{ij}$ represents the relation from query $\mathbf{q}_i$ to the value $\mathbf{v}_j$ within window $M$, and $(EV)_i$ aggregates all related value vectors with respective to query $\mathbf{q}_i$. Unfortunately, the $EV$ term still requires $O(N^2)$ space complexity and $O(N^2C)$ time complexity. In CoaT, we propose to simplify the $EV$ term to $\hat{EV}$ by considering each channel in the query, key and value vectors as *internal heads*. Thus, for each internal head $l$, we have:

$$E_{ij}^{(l)} = \mathbb{1}(i,j)q_i^{(l)}p_{j-i}^{(l)}, \;\; \hat{EV}_i^{(l)} = \sum_j E_{ij}^{(l)}v_j^{(l)} \quad (7)$$

In practice, we can use a 1-D depthwise convolution to compute $\hat{EV}$:

$$\hat{EV}^{(l)} = Q^{(l)} \circ \text{Conv1D}(P^{(l)}, V^{(l)}), \quad (8)$$

$$\hat{EV} = Q \circ \text{DepthwiseConv1D}(P, V) \quad (9)$$

where $\circ$ is the Hadamard product. It is noteworthy that in vision Transformers, we have two types of tokens, the class (`CLS`) token and image tokens. Thus, we use a 2-D depthwise convolution (with window size $M$ and kernel $P$) and apply it only to the reshaped image tokens (i.e. $Q^{\text{img}}, V^{\text{img}} \in \mathbb{R}^{H \times W \times C}$ from $Q, V$ respectively):

$$\hat{EV}^{\text{img}} = Q^{\text{img}} \circ \text{DepthwiseConv2D}(P, V^{\text{img}}) \quad (10)$$

$$\hat{EV} = \text{concat}(\hat{EV}^{\text{img}}, \mathbf{0}) \quad (11)$$

$$\text{ConvAtt}(X) = \frac{Q}{\sqrt{C}}\Big(\text{softmax}(K)^{\top}V\Big) + \hat{EV} \quad (12)$$

Based on our derivation, the depthwise convolution can be seen as a special case of relative position encoding.
*Convolutional Relative Position Encoding vs Other Relative Position Encoding.* The commonly referred relative position encoding [27] works in the standard scaled dot-product attention settings since the encoding matrix $E$ is combined with the softmax logits in the attention maps, which are not materialized in our factorized attention. Related to our work, LambdaNets [1] attempts to use a 3-D convolution to compute $EV$ directly, but it costs $O(NC^2)$ space complexity and $O(NC^2M^2)$ time complexity, which leads to heavy computation when channel size $C$ is large. In contrast, our factorized attention computes $\hat{EV}$ that only takes $O(NC)$ space complexity and $O(NCM^2)$ time complexity, achieving better efficiency than LambdaNets.

**Convolutional Position Encoding.** We then extend the idea of convolutional relative position encoding to a general convolutional position encoding case. Convolutional relative position encoding models local position-based relationship between queries and values. Similar to the absolute position encoding used in most image Transformers [8, 34], we would like to insert the position relationship into the input
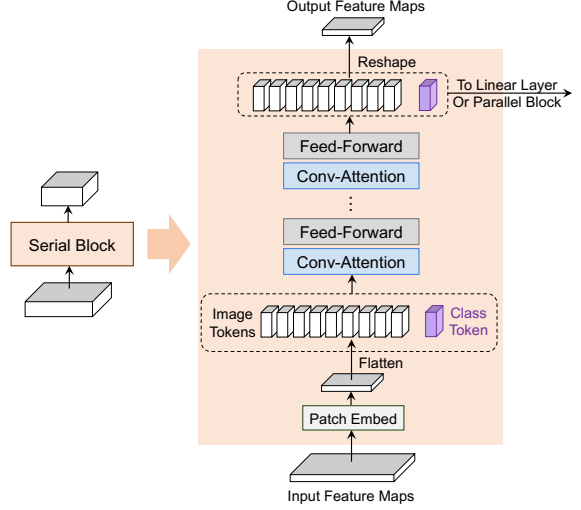


Figure 5. **Schematic illustration of the serial block in CoaT.** Input feature maps are first down-sampled by a patch embedding layer, and then tokenized features (along with a class token) are processed by multiple conv-attention and feed-forward layers.

image features directly to enrich the effects of relative position encoding. In each conv-attentional module, we insert a depthwise convolution into the input features $X$ and concatenate the resulting position-aware features back to the input features following the standard absolute position encoding scheme (see Figure 2 lower part), which resembles the realization of conditional position encoding in CPVT [4].

CoaT and CoaT-Lite share the convolutional position encoding weights and convolutional relative position encoding weights for the serial and parallel modules within the same scale. We set convolution kernel size to 3 for the convolutional position encoding. We set convolution kernel size to 3, 5 and 7 for image features from different attention heads for convolutional relative position encoding.

The work of CPVT [4] explores the use of convolution as conditional position encodings by inserting it after the feed-forward network under a single scale ($\frac{H}{16} \times \frac{W}{16}$). Our work focuses on applying convolution as relative position encoding and a general position encoding with factorized attention in a co-scale setting.

### 4.3. Conv-Attentional Mechanism

The final conv-attentional module is shown in Figure 2: We apply the first convolutional position encoding on the image tokens from the input. Then, we feed it into $\text{ConvAtt}(\cdot)$ including factorized attention and the convolutional relative position encoding. The resulting map is used for the subsequent feed-forward networks.

We show an illustration of the proposed convolutional relative position encoding strategy in Figure 3. Each query in the image tokens $Q^{\text{img}}$ seeks all its nearby values within the window from value image tokens $V^{\text{img}}$. The product of
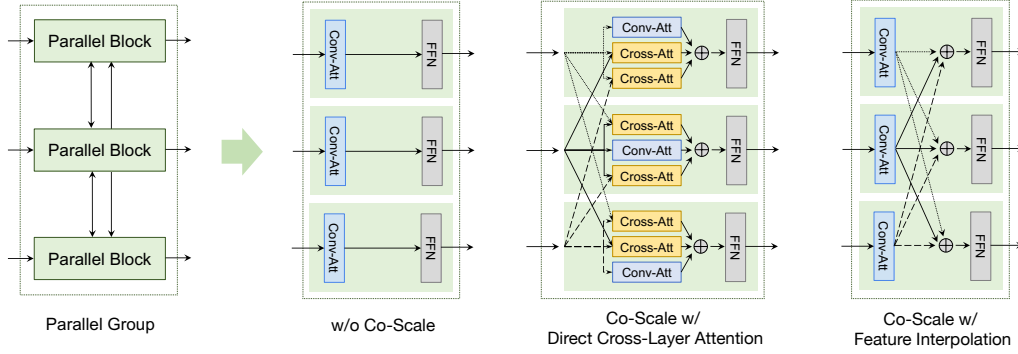
Figure 6. **Schematic illustration of the parallel group in CoaT.** For "w/o Co-Scale", tokens learned at the individual scales are combined to perform the classification but absent intermediate co-scale interaction for the individual paths of the parallel blocks. We propose two "Co-Scale" variants, namely Direct Cross-Layer Attention and Feature Interpolation. Co-Scale with Feature Interpolation is adopted in the final CoaT-Lite and CoaT models reported in the ImageNet benchmark.

the query, each obtained value, and corresponding relative position encoding in $P$ is then summed to output $\hat{E}V$. To reduce the computation, we consider the relative position encoding map $P$ as a convolutional kernel and convolve it with the value image tokens $V^{\text{img}}$ first. Then, we multiply the query with the result of convolution and generate the output $\hat{E}V$. Note that Figure 3 shows a simple case where query, position encoding, and value have a single channel (corresponding to one internal head in Equation 8). For the multi-channel version (corresponding to Equation 9, 10), the multiplication and the convolution in Figure 3 should be replaced by a Hadamard product and a depthwise convolution operation.

## 5. Co-Scale Conv-Attentional Transformers

### 5.1. Co-Scale Mechanism

The proposed co-scale mechanism is designed to introduce cross-scale attention to image transformers. Here, we describe two types of co-scale blocks in the CoaT architecture, namely serial and parallel blocks.

**CoaT Serial Block.** A serial block (shown in Figure 5) models image representations in a reduced resolution. In a typical serial block, we first down-sample input feature maps by a certain ratio using a patch embedding layer (2D convolution layer), and flatten the reduced feature maps into a sequence of image tokens. We then concatenate image tokens with an additional CLS token, a specialized vector to perform image classification, and apply multiple conv-attentional modules as described in Section 4 to learn internal relationships among image tokens and the CLS token. Finally, we separate the CLS token from the image tokens and reshape the image tokens to 2D feature maps for the next serial block.

**CoaT Parallel Block.** We realize cross-scale attention between parallel blocks in each parallel group (shown in Figure 6). In a typical parallel group, we have sequences of input

features (image tokens and CLS token) from serial blocks with different scales. To realize fine-to-coarse, coarse-to-fine, and cross-scale attention in the parallel group, we develop two strategies: (1) direct cross-layer attention; (2) attention with feature interpolation. In this paper, we adopt attention with feature interpolation for better empirical performance. The effectiveness of both strategies are shown in Section 6.4.
*Direct cross-layer attention.* In direct cross-layer attention, we form query, key, and value vectors from input features for each scale. For attention within the same layer, we use the conv-attention (Figure 2) with the query, key and value vectors from current scale. For attention across different layers, we down-sample or up-sample the key and value vectors to match the resolution of other scales. We then perform a cross-attention, which extends the conv-attention with query from current scale and key and value from another scale. Finally, we sum outputs of conv-attention and cross-attention together and applies a shared feed-forward layer. With direct cross-layer attention, the cross-layer information is fused in a cross-attention fashion.
*Attention with feature interpolation.* Instead of performing cross-layer attention directly, we also present attention with feature interpolation. First, the input image features from different scales are processed by independent conv-attention modules. Then, we down-sample or up-sample image features from each scale to match the other scales' dimensions using bilinear interpolation or keep the same for its own scale. The features belong to the same scale are summed in the parallel group, and they are further passed into a shared feed-forward layer. In this way, the conv-attentional module in the next step can learn cross-layer information based on the feature interpolation in the current step.

### 5.2. Model Architecture

**CoaT-Lite.** CoaT-Lite, Figure 4 (Left), processes input images with a series of serial blocks following a fine-to-coarse pyramid structure. Given an input image $I \in \mathbb{R}^{H \times W \times C}$,

Table 1. **Architecture details of CoaT-Lite and CoaT models.** $C_i$ represents the hidden dimension of the attention layers in block $i$; $H_i$ represents the number of attention heads in the attention layers in block $i$; $R_i$ represents the expansion ratio for the feed-forward hidden layer dimension between attention layers in block $i$. Multipliers indicate the number of conv-attentional modules in block $i$.

| Blocks | Output | CoaT-Lite | | | CoaT | |
|---|---|---|---|---|---|---|
| | | Tiny | Mini | Small | Tiny | Mini |
| Serial Block ($S_1$) | $56 \times 56$ | $\begin{bmatrix} C_1 = 64 \\ H_1 = 8 \\ R_1 = 8 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_1 = 64 \\ H_1 = 8 \\ R_1 = 8 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_1 = 64 \\ H_1 = 8 \\ R_1 = 8 \end{bmatrix} \times 3$ | $\begin{bmatrix} C_1 = 152 \\ H_1 = 8 \\ R_1 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_1 = 152 \\ H_1 = 8 \\ R_1 = 4 \end{bmatrix} \times 2$ |
| Serial Block ($S_2$) | $28 \times 28$ | $\begin{bmatrix} C_2 = 128 \\ H_2 = 8 \\ R_2 = 8 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_2 = 128 \\ H_2 = 8 \\ R_2 = 8 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_2 = 128 \\ H_2 = 8 \\ R_2 = 8 \end{bmatrix} \times 4$ | $\begin{bmatrix} C_2 = 152 \\ H_2 = 8 \\ R_2 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_2 = 216 \\ H_2 = 8 \\ R_2 = 4 \end{bmatrix} \times 2$ |
| Serial Block ($S_3$) | $14 \times 14$ | $\begin{bmatrix} C_3 = 256 \\ H_3 = 8 \\ R_3 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_3 = 320 \\ H_3 = 8 \\ R_3 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_3 = 320 \\ H_3 = 8 \\ R_3 = 4 \end{bmatrix} \times 6$ | $\begin{bmatrix} C_3 = 152 \\ H_3 = 8 \\ R_3 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_3 = 216 \\ H_3 = 8 \\ R_3 = 4 \end{bmatrix} \times 2$ |
| Serial Block ($S_4$) | $7 \times 7$ | $\begin{bmatrix} C_4 = 320 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_4 = 512 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_4 = 512 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 3$ | $\begin{bmatrix} C_4 = 152 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} C_4 = 216 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 2$ |
| Parallel Group | $\begin{bmatrix} 28 \times 28 \\ 14 \times 14 \\ 7 \times 7 \end{bmatrix}$ | | | | $\begin{bmatrix} C_4 = 152 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 6$ | $\begin{bmatrix} C_4 = 216 \\ H_4 = 8 \\ R_4 = 4 \end{bmatrix} \times 6$ |
| #Params | | 5.7M | 11M | 20M | 5.5M | 10M |

each serial block down-samples the image features into lower resolution, resulting in a sequence of four resolutions: $F_1 \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_1}$, $F_2 \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times C_2}$, $F_3 \in \mathbb{R}^{\frac{H}{16} \times \frac{W}{16} \times C_3}$, $F_4 \in \mathbb{R}^{\frac{H}{32} \times \frac{W}{32} \times C_4}$. In CoaT-Lite, we obtain the CLS token in the last serial block, and perform image classification via a linear projection layer based on the CLS token.

**CoaT.** Our CoaT model, shown in Figure 4 (Right), consists of both serial and parallel blocks. Once we obtain multi-scale feature maps $\{F_1, F_2, F_3, F_4\}$ from the serial blocks, we pass $F_2, F_3, F_4$ into the parallel group with three separate parallel blocks. In CoaT, we concatenate the independent CLS tokens for each feature scale, and reduce the channel dimension to perform image classification with the same procedure as CoaT-Lite.

**Model Variants.** In this paper, we explore CoaT and CoaT-Lite with three different model sizes, namely Tiny, Mini, and Small. Architecture details are shown in Table 1. For example, tiny models represent those with a 5M parameter budget constraint. Specifically, these tiny models have four serial blocks, each with two conv-attentional modules. In CoaT-Lite Tiny architectures, the hidden dimensions of the attention layers increase for later blocks. CoaT Tiny sets the hidden dimensions of the attention layers in the parallel group to be equal, and performs the co-scale attention within the parallel group for six steps. Mini and small models follow the same architecture design but with increased embedding dimensions and increased numbers of conv-attentional modules within blocks.

## 6. Experiments

### 6.1. Experiment Details

**Image Classification.** We perform image classification on the standard ILSVRC-2012 ImageNet dataset [26]. The standard ImageNet benchmark contains 1.3 million images in the training set and 50K images in the validation set, covering 1000 object classes. Image cropping sizes are set to $224 \times 224$. For fair comparison, we perform data augmentation such as MixUp [48], CutMix [47], random erasing [51], repeated augmentation [11], and label smoothing [31], following identical procedures in DeiT [34] with the exception that stochastic depth is not used [14].

All experiment results of our models in Table 2 are reported at 300 epochs, consistent with previous methods [34]. We train all models with a global batch size of 2048 with the NVIDIA Automatic Mixed Precision (AMP) enabled. We adopt the AdamW [20] optimizer with cosine learning rate decay, five warm-up epochs, and weight decay 0.05. The learning rate is scaled as $5 \times 10^{-4} \times \frac{\text{global batch size}}{512}$.

**Object Detection and Instance Segmentation.** We conduct object detection and instance segmentation experiments on the Common Objects in Context (COCO2017) dataset [19]. The COCO2017 benchmark contains 118K training images and 5K validation images. We evaluate the generalization ability of CoaT in object detection and instance segmentation with the Mask R-CNN [9] framework. We follow the identical data processing settings in Mask R-CNN and enable the feature pyramid network (FPN) [18] to utilize multi-scale features. In addition, we perform object detection based on Deformable DETR [52] following its data processing settings, using random horizontal flips, resizing, and cropping as augmentation techniques.

Table 2. **CoaT performance on ImageNet-1K validation set.** Our CoaT models consistently outperform other methods while being parameter efficient. ConvNets, attention-based models and Transformers with similar model size are grouped together for comparison. "#GFLOPs" is calculated with the input size 224×224. "*" indicates the improved performance after retraining with a better training scheme.

| Arch. | Model | #Params | Input | #GFLOPs | Top-1 Acc. @input |
|---|---|---|---|---|---|
| ConvNets | EfficientNet-B0 [33] | 5.3M | $224^2$ | 0.4 | 77.1% |
| | ShuffleNet [49] | 5.4M | $224^2$ | 0.5 | 73.7% |
| | MobileNet [12] | 6.9M | $224^2$ | 0.6 | 74.7% |
| | MnasNet-A3 [32] | 5.2M | $224^2$ | 0.4 | 76.7% |
| Transformers | DeiT-Tiny [34] | 5.7M | $224^2$ | 1.3 | 72.2% |
| | CPVT-Tiny [4] | 5.7M | $224^2$ | - | 73.4% |
| | T2T-ViT-10 [46] | 5.8M | $224^2$ | 1.2 | 74.1% |
| | **CoaT-Lite Tiny** (Ours) | 5.7M | $224^2$ | 1.6 | 76.6% |
| | **CoaT Tiny** (Ours) | 5.5M | $224^2$ | 4.4 | **78.2%** |
| ConvNets | EfficientNet-B2 [33] | 9M | $260^2$ | 1.0 | 80.1% |
| | ResNet-18 [10] | 12M | $224^2$ | 1.8 | 69.8% |
| | REDNet-26 [17] | 9M | $224^2$ | 1.7 | 75.9% |
| Attention | SAN10 [50] | 11M | $224^2$ | 2.2 | 77.1% |
| | LambdaNets [1] | 15M | $224^2$ | - | 78.4% |
| Transformers | PVT-Tiny [38] | 13M | $224^2$ | 1.9 | 75.1% |
| | **CoaT-Lite Mini** (Ours) | 11M | $224^2$ | 2.0 | 78.9% |
| | **CoaT Mini** (Ours) | 10M | $224^2$ | 6.8 | **80.8%** |
| ConvNets | EfficientNet-B4 [33] | 19M | $380^2$ | 4.2 | **82.9%** |
| | ResNet-50 [10] | 25M | $224^2$ | 4.1 | 76.2% |
| | ResNet-50* [10] | 25M | $224^2$ | 4.1 | 79.1% |
| | ResNeXt50-32x4d [42] | 25M | $224^2$ | 4.3 | 77.6% |
| | ResNeXt50-32x4d* [42] | 25M | $224^2$ | 4.3 | 79.5% |
| | REDNet-101 [17] | 25M | $224^2$ | 4.7 | 79.1% |
| | REDNet-152 [17] | 34M | $224^2$ | 6.8 | 79.3% |
| | ResNet-101 [10] | 45M | $224^2$ | 7.9 | 77.4% |
| | ResNet-101* [10] | 45M | $224^2$ | 7.9 | 79.9% |
| | ResNeXt101-32x4d [42] | 45M | $224^2$ | 8.0 | 78.8% |
| | ResNeXt101-32x4d* [42] | 45M | $224^2$ | 8.0 | 80.6% |
| | ResNet-152 [10] | 60M | $224^2$ | 11.6 | 78.3% |
| | ResNet-152* [10] | 60M | $224^2$ | 11.6 | 80.8% |
| | ResNeXt101-64x4d [42] | 84M | $224^2$ | 15.6 | 79.6% |
| | ResNeXt101-64x4d* [42] | 84M | $224^2$ | 15.6 | 81.5% |
| Transformers | DeiT-Small [34] | 22M | $224^2$ | 4.6 | 79.8% |
| | PVT-Small [38] | 24M | $224^2$ | 3.8 | 79.8% |
| | CPVT-Small [4] | 22M | $224^2$ | - | 80.5% |
| | T2T-ViT$_t$-14 [46] | 22M | $224^2$ | 5.2 | 80.7% |
| | T2T-ViT$_t$-19 [46] | 39M | $224^2$ | 8.4 | 81.4% |
| | PVT-Medium [38] | 44M | $224^2$ | 6.7 | 81.2% |
| | PVT-Large [38] | 61M | $224^2$ | 9.8 | 81.7% |
| | DeiT-Base [34] | 86M | $224^2$ | 17.6 | 81.8% |
| | CPVT-Base [4] | 86M | $224^2$ | - | **81.9%** |
| | ViT-B/16 [8] | 86M | $384^2$ | 55.4 | 77.9% |
| | ViT-L/16 [8] | 307M | $384^2$ | 190.7 | 76.5% |
| | **CoaT-Lite Small** (Ours) | 20M | $224^2$ | 4.0 | **81.9%** |

For Mask R-CNN optimization, we train the model with the ImageNet-pretrained backbone on 8 GPUs via SGD with momentum with learning rate 0.02. The training period contains 90K steps in 1× setting and 270K steps in 3× setting following Detectron2 [41]. For Deformable DETR optimization, we train the model with the pretrained backbone for 50 epochs, using an AdamW optimizer with initial learning rate $2 \times 10^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We reduce the learning rate by a factor of 10 at epoch 40.

## 6.2. CoaT for ImageNet Classification

Table 2 shows top-1 accuracy results of our models on the ImageNet validation set comparing with state-of-the-art methods. Except for EfficientNets, all reported models are evaluated with $224 \times 224$ image cropping. We separate model architectures into three categories: convolu-

Table 3. **Object detection and instance segmentation results based on Mask R-CNN on the COCO val2017**. Two Mask R-CNN with FPN settings (1× for 90k steps and 3× for 270k steps) are compared. We compare CoaT-Lite and CoaT results with ResNet and PVT. PVT results are taken from the reported numbers from [38] and its official repository.

| Backbone | #Params (M) | Mask R-CNN w/ FPN 1× | | | | | | Mask R-CNN w/ FPN 3× | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^m$ | $AP^m_{50}$ | $AP^m_{75}$ | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^m$ | $AP^m_{50}$ | $AP^m_{75}$ |
| ResNet-18 | 31.3 | 34.2 | 54.1 | 36.7 | 31.3 | 51.1 | 33.2 | 36.3 | 56.5 | 39.3 | 33.2 | 53.5 | 35.4 |
| PVT-Tiny [38] | 32.9 | 36.7 | 59.2 | 39.3 | 35.1 | 56.7 | 37.3 | 39.8 | 62.2 | 43.0 | 37.4 | 59.3 | 39.9 |
| **CoaT-Lite Mini** (Ours) | 30.7 | 39.2 | 61.5 | 42.1 | 36.0 | 58.0 | 38.5 | 41.6 | 63.0 | 45.1 | 37.6 | 59.7 | 40.2 |
| **CoaT Mini** (Ours) | 30.2 | **43.0** | **64.7** | **46.8** | **38.7** | **61.4** | **41.6** | **45.2** | **65.9** | **49.6** | **40.3** | **63.0** | **43.1** |
| ResNet-50 | 44.3 | 38.6 | 59.5 | 42.1 | 35.2 | 56.3 | 37.4 | 41.0 | 61.9 | 44.5 | 37.2 | 58.6 | 39.9 |
| PVT-Small [38] | 44.1 | 40.4 | 62.9 | 43.8 | 37.8 | 60.1 | 40.3 | 43.0 | 65.3 | 46.9 | 39.9 | **62.5** | 42.8 |
| **CoaT-Lite Small** (Ours) | 39.5 | **43.6** | **65.3** | **47.4** | **39.2** | **62.0** | **41.8** | **44.7** | **66.0** | **48.8** | **40.1** | 62.4 | **43.2** |

Table 4. **Object detection results based on Deformable DETR on COCO val2017**. DD ResNet-50 (rep.) is our reproduced baseline result. ResNet-50 and our CoaT-Lite are directly comparable as the DD backbone because of the similar model size.

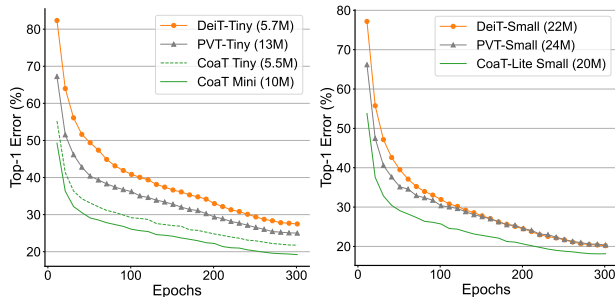| Backbone | Deformable DETR (Multi-Scale) | | | | | |
|---|---|---|---|---|---|---|
| | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
| DD ResNet-50 [52] | 44.5 | - | - | 27.6 | 47.6 | 59.6 |
| DD ResNet-50 (rep.) | 44.0 | 62.9 | 48.0 | 26.0 | 47.4 | 58.4 |
| DD **CoaT-Lite Small** (Ours) | **46.9** | **66.3** | **51.0** | **28.4** | **50.3** | **62.5** |



Figure 7. **Training curves on ImageNet-1K.** (Left) CoaT Tiny and CoaT Mini top-1 validation errors. (Right) CoaT-Lite Small top-1 validation errors.

tional networks (ConvNets), attention-based models (non-Transformer), and Transformers. Under around 5M, 10M, and 20M parameter budget constraints, CoaT and CoaT-Lite surpasses all reported Transformer-based architectures (see Table 2). In particular, our CoaT models bring a large performance gain to the baseline DeiT [34], which shows that our co-scale mechanism is essential to improve the performance of Transformer-based architectures.

Our CoaT Tiny model achieves a 21.8% error rate, a decrease of 6.0% from our DeiT-Tiny baseline. CoaT Tiny also outperforms the strongly competitive model EfficientNet-B0 by 1.1%. CoaT Mini surpasses PVT-Tiny and EfficientNet-B2 by 5.7% and 0.7% with similar model sizes. CoaT-Lite also achieves strong results under the three different model sizes while being competitively fast. Notably, our CoaT-Lite Small model (20M parameters) outperforms all reported Transformer-based architectures using similar or significantly larger model sizes.

We show training curves for CoaT and CoaT-Lite in Fig-

ure 7. CoaT converges significantly faster than competing image Transformers while achieving better generalization ability.

## 6.3. Object Detection and Instance Segmentation

Table 3 demonstrates CoaT object detection and instance segmentation results under the Mask R-CNN framework on the COCO val2017 dataset. Our CoaT and CoaT-Lite models show clear performance advantages over the ResNet and PVT backbones under both the 1× setting and the 3× setting. Our CoaT Mini obtains significant performance improvement over CoaT-Lite Mini.

We additionally perform object detection with the Deformable DETR (DD) framework in Table 4. We compare our models with the standard ResNet-50 backbone on the COCO dataset [19]. Our CoaT-Lite Small as the backbone achieves 2.9% improvement on average precision (AP) over the reproduced results of Deformable DETR with ResNet-50 [52].

## 6.4. Ablation Study

**Effectiveness of Convolutional Position Encoding.** We study the effectiveness of the combination of the convolutional position encoding in our Conv-Attention Module in Table 5. Our CoaT-Lite without any convolution position encoding results in poor performance (69.0% top-1 accuracy), indicating position encoding is essential for image transformers. We observe great improvements for CoaT-Lite variants with only either convolutional position encoding (Conv-Pos: Figure 2 bottom) or convolutional relative position encoding (Conv-Rel-Pos: Figure 2 top-right), which achieve 76.0% and 73.5% top-1 accuracy accordingly. We found CoaT-Lite with the combination of Conv-Pos and Conv-Rel-Pos learns a better classifier (76.6% top-1 accuracy), making both position encoding schemes companion rather than conflict.

**Effectiveness of Convolutional Relative Position Encoding.** For both CoaT-Lite and CoaT models, we report models with and without convolutional relative position encoding (Conv-Rel-Pos) in Table 6. We found consistent improvement for both CoaT-Lite and CoaT models that equipping

Table 5. **Effectiveness of convolutional position encoding.** All experiments are performed in the CoaT-Lite Tiny architecture. Performance is evaluated on ImageNet-1K validation set.

| Model | Conv-Pos | Conv-Rel-Pos | Top-1 Acc. |
|---|---|---|---|
| CoaT-Lite (Tiny) | ✗ | ✗ | 69.0% |
| | ✗ | ✓ | 73.5% |
| | ✓ | ✗ | 76.0% |
| | ✓ | ✓ | 76.6% |

Conv-Rel-Pos. Besides, the Conv-Rel-Pos is able to improve performance at modest increase in computational overhead.

Table 6. **Effectiveness of convolutional relative position encoding** are evaluated in the CoaT-Lite and CoaT architectures under Tiny, Mini and Small parameter settings. "w/o relative" indicates the model with Conv-Pos only and the other model with no indication are the model with both Conv-Pos and Conv-Rel-Pos. Performance is evaluated on ImageNet-1K validation set.

| Size | Model | #Params | Input | #GFLOPs | Top-1 Acc. @input |
|---|---|---|---|---|---|
| Tiny | CoaT-Lite w/o relative | 5.6M | $224^2$ | 1.6 | 76.0% |
| | CoaT w/o relative | 5.5M | $224^2$ | 4.3 | 77.3% |
| | CoaT-Lite | 5.7M | $224^2$ | 1.6 | 76.6% |
| | CoaT | 5.5M | $224^2$ | 4.4 | **78.2%** |
| Mini | CoaT-Lite w/o relative | 11M | $224^2$ | 1.9 | 78.2% |
| | CoaT w/o relative | 10M | $224^2$ | 6.7 | 80.4% |
| | CoaT-Lite | 11M | $224^2$ | 2.0 | 78.9% |
| | CoaT | 10M | $224^2$ | 6.8 | **80.8%** |
| Small | CoaT-Lite w/o relative | 19M | $224^2$ | 3.9 | 81.4% |
| | CoaT-Lite | 20M | $224^2$ | 4.0 | **81.9%** |

**Effectiveness of Co-Scale.**  In Table 7, we present performance results for two co-scale variants in CoaT, direct cross-layer attention and attention with feature interpolation. We also report CoaT without co-scale as a baseline. Comparing to CoaT without a co-scale mechanism, both co-scale variants show significant performance improvements. Attention with feature interpolation offers a clear advantage over direct cross-layer attention due to less computational complexity and higher accuracy.

Table 7. **Effectiveness of co-scale.** All experiments are performed in the CoaT Tiny architecture. Performance is evaluated on ImageNet-1K validation set.

| Size | Model | #Params | Input | #GFLOPs | Top-1 Acc. @input |
|---|---|---|---|---|---|
| Tiny | CoaT w/o co-scale | 5.5M | $224^2$ | 4.4 | 76.2% |
| | CoaT w/ co-scale | | | | |
| | - direct cross-layer attention | 5.5M | $224^2$ | 4.8 | 77.8% |
| | - attention with feature interpolation | 5.5M | $224^2$ | 4.4 | **78.2%** |

## 7. Visualization

We show feature and attention visualizations of our proposed CoaT model and DeiT [34] in Figure 8. For the sampled feature maps shown on the right side of the figure, we directly sample the first six feature maps after different kinds of attention blocks. The visualization of CLS attention maps is done by attending the CLS (query) to all other spatial
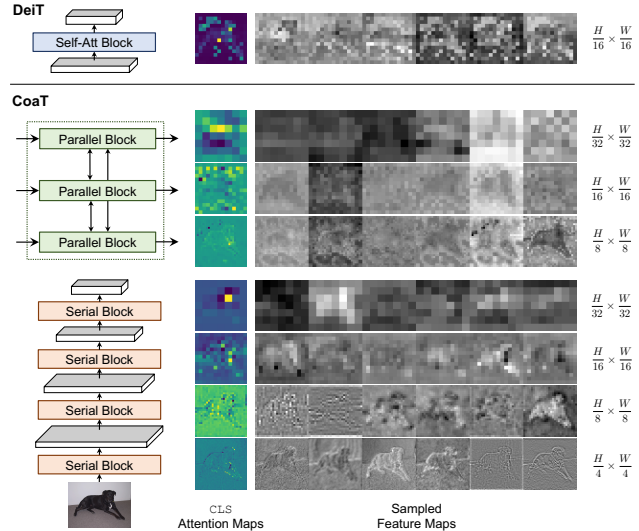


Figure 8. **Visualization of CLS attention maps and sampled feature maps.** The visualization is conducted on DeiT [34] and our proposed CoaT model trained on ImageNet. The visualization of the serial block is adopted from CoaT-Lite Mini and the visualization of the parallel block is adopted from CoaT Mini. The first six feature maps are sampled from each layer.

positions (keys) in the feature map. Note that although our factorized attention mechanism does matrix multiplication between keys and values during training, we are still able to materialize the CLS attention map that resembles the scaled dot-product attention.

Without the coarse-to-fine route, the sampled features in DeiT cannot capture low-level structure features that are essential for downstream tasks. The feature samples from DeiT also show low feature richness, resulting in poor classification performance. In our CoaT visualization, we show high-diversity multi-scale feature maps. From the visualization of serial blocks, we see both high-level abstraction and low-level parts of a dog are captured. Contexts play an important role in object recognition [22] and semantic labeling [35]. The attention visualization from parallel blocks shows that multi-scale contexts are further mixed to enhance feature richness.

## 8. Conclusion

In this paper, we have presented a Transformer based image classifier, Co-scale conv-attentional image Transformer (CoaT), in which cross-scale attention and efficient convolution-like attention operations have been developed. CoaT's small models attain strong classification results on ImageNet, and their applicability to downstream computer vision tasks have been demonstrated in object detection and instance segmentation.

# References

[1] Irwan Bello. Lambdanetworks: Modeling long-range interactions without attention. In *ICLR*, 2021.

[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.

[3] Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.

[4] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Do we really need explicit position encodings for vision transformers? *arXiv preprint arXiv:2102.10882*, 2021.

[5] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, 2017.

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.

[7] Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 5884–5888, 2018.

[8] 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. In *ICLR*, 2021.

[9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017.

[10] 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.

[11] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefler, and Daniel Soudry. Augment your batch: Improving generalization through instance repetition. In *CVPR*, 2020.

[12] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[13] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3464–3473, 2019.

[14] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016.

[15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Advances in Neural Inf. Process. Syst.*, 2012.

[16] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[17] Duo Li, Jie Hu, Changhu Wang, Xiangtai Li, Qi She, Lei Zhu, Tong Zhang, and Qifeng Chen. Involution: Inverting the inherence of convolution for visual recognition. In *CVPR*, 2021.

[18] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.

[19] 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 *ECCV*, 2014.

[20] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[21] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.

[22] Aude Oliva and Antonio Torralba. The role of context in object recognition. *Trends in cognitive sciences*, 11(12):520–527, 2007.

[23] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.

[24] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019.

[25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 2015.

[26] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

[27] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *NAACL-HLT*, 2018.

[28] Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3531–3539, 2021.

[29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[30] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.

[31] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.

[32] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *CVPR*, 2019.

[33] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.

[34] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.

[35] Zhuowen Tu. Auto-context and its application to high-level vision tasks. In *CVPR*, 2008.

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. Advances in Neural Inf. Process. Syst.*, 2017.

[37] J Wang, K Sun, T Cheng, B Jiang, C Deng, Y Zhao, D Liu, Y Mu, M Tan, X Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[38] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. *arXiv preprint arXiv:2102.12122*, 2021.

[39] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.

[40] Andrew Witkin. Scale-space filtering: A new approach to multi-scale description. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 9, pages 150–153, 1984.

[41] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. `https://github.com/facebookresearch/detectron2`, 2019.

[42] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017.

[43] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional shapecontextnet for point cloud recognition. In *CVPR*, 2018.

[44] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.

[45] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. Layoutlm: Pre-training of text and layout for document image understanding. In *SIGKDD*, 2020.

[46] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021.

[47] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019.

[48] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[49] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *CVPR*, 2018.

[50] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *CVPR*, 2020.

[51] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, pages 13001–13008, 2020.

[52] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2021.