
Not All Images are Worth 16x16 Words: Dynamic Transformers for Efficient Image Recognition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Vision Transformers (ViT) have achieved remarkable success in large-scale image
2 recognition. They split each 2D image into a fixed number of patches, each of
3 which is treated as a token. Generally, representing an image with more tokens
4 would lead to higher prediction accuracy, while it also results in drastically in-
5 creased computational cost. To achieve a decent trade-off between accuracy and
6 speed, the number of tokens is empirically set to 16x16 or 14x14. In this paper,
7 we argue that each image has its own characteristics, and ideally the token number
8 should be conditioned on the individual input. In fact, we have observed that there
9 are some “easy” images which can be accurately predicted with a mere number of
10 4x4 tokens, while some “hard” ones may need a finer representation. Inspired by
11 this phenomenon, we propose a Dynamic Transformer to automatically configure
12 a proper number of tokens for each input image. This is achieved by cascading
13 multiple Transformers with increasing number of tokens, which are sequentially
14 activated in an adaptive fashion at test time, i.e., the inference is terminated once
15 a sufficiently confident prediction is produced. We further design efficient fea-
16 ture reuse and relationship reuse mechanisms across different components of the
17 Dynamic Transformer to reduce redundant computation. Extensive empirical
18 results on ImageNet, CIFAR-10, and CIFAR-100 demonstrate that our method
19 significantly outperforms the competitive baselines in terms of both theoretical
20 computational efficiency and practical inference speed.

21

1 Introduction

22 Transformers, the dominant self-attention-based models in natural language processing (NLP) [9,
23 34, 3], have been successfully adapted to image recognition problems [10, 49, 32, 15] recently. In
24 particular, vision Transformers achieve state-of-the-art performance on the large scale ImageNet
25 benchmark [8], while exhibit excellent scalability with the further growing dataset size (e.g., on JFT-
26 300M [10]). These models split each image into a fixed number of patches and embed them into 1D
27 tokens as inputs. Typically, representing the data using more tokens contributes to higher prediction
28 accuracy, but leads to intensive computational cost, which grows quadratically with respect to the
29 token number in self-attention blocks. For a proper trade-off between efficiency and effectiveness,
30 existing works empirically adopt 14x14/16x16 tokens [10, 49].

31 In this paper, we argue that it may not be optimal to treat all
32 samples with the same number of tokens. In fact, there exist
33 considerable variations among different images (e.g., contents,
34 scales of objects, backgrounds, etc.). Therefore, the number of
35 representative tokens should ideally be configured specifically
36 for each input. This issue is critical for the computational
37 efficiency of the models. For example, we train a T2T-ViT-12 [49] with varying token numbers, and

Table 1: Accuracy and computa-
tional cost of T2T-ViT-12 with dif-
ferent token numbers on ImageNet.

# of Tokens	14x14	7x7	4x4
Accuracy	76.7%	70.3%	60.8%
FLOPs	1.78G	0.47G	0.21G

38 report the corresponding accuracy and FLOPs in Table 1. One can observe that adopting the officially
 39 recommended 14x14 tokens only correctly recognizes $\sim 15.9\%$ (76.7% v.s. 60.8%) more test samples
 40 compared to that of using 4x4 tokens, while increases the computational cost by 8.5 times (1.78G
 41 v.s. 0.21G). In other words, computational resources are wasted on applying the unnecessary 14x14
 42 tokens to many “easy” images for which 4x4 tokens are sufficient.

43 Motivated by this observation, we propose a novel *Dynamic Vision Transformer* (DVT) framework, aiming to
 44 automatically configure a decent token number conditioned
 45 on each image for high computational efficiency. In spe-
 46 cific, a cascade of Transformers are trained using increas-
 47 ing number of tokens. At test time, these models are se-
 48 quentially activated starting with less tokens. Once a pre-
 49 diction with sufficient confidence has been produced, the
 50 inference procedure will be terminated immediately. As a
 51 consequence, the computation is unevenly allocated among
 52 “easy” and “hard” samples by adjusting the token number,
 53 yielding a considerable improvement in efficiency. Impor-
 54 tantly, we further develop *feature-wise* and *relationship-
 55 wise* reuse mechanisms to reduce redundant computation.
 56 The former allows the downstream models to be trained on
 57 the basis of previously extracted deep features, while the
 58 later enables leveraging existing upstream self-attention
 59 relationships to learn more accurate attention maps. Illus-
 60 trative examples of our method are given in Figure 1.
 61

62 Notably, DVT is designed as a general framework. Most of the state-of-the-art image recognition
 63 Transformers, such as ViT [10], DeiT [32], and T2T-ViT [49], can be straightforwardly deployed as
 64 its backbones for higher efficiency. Our method is also appealing in its flexibility. The computational
 65 cost of DVT is able to be adjusted online by simply adapting the early-termination criterion. This
 66 characteristic makes DVT suitable for the cases where the available computational resources fluctuate
 67 dynamically or a minimal power consumption is required to achieve a given performance. Both
 68 situations are ubiquitous in real-world applications (e.g., searching engines and mobile apps).

69 The performance of DVT is evaluated on ImageNet [8] and CIFAR [23] with T2T-ViT [49] and DeiT
 70 [32]. Experimental results show that DVT significantly improves the efficiency of the backbones. For
 71 examples, DVT reduces the computational cost of T2T-ViT by 1.6-3.6x without sacrificing accuracy.
 72 The real inference speed on a NVIDIA 2080Ti GPU is consistent with our theoretical results.

73 2 Related Work

74 **Vision Transformers.** Inspired by the success of Transformers on NLP tasks [9, 34, 3, 37], vision
 75 Transformers (ViT) have recently been developed for image recognition [10]. Although ViT by itself
 76 is not comparable with state-of-the-art convolutional networks (CNN) on the standard ImageNet
 77 benchmark, it attains excellent results when pre-trained on the larger JFT-300M dataset. DeiT [32]
 78 studies the training strategy of ViT and proposes a knowledge distilling-based approach, surpassing
 79 the performance of ResNet [17]. Some following works such as T2T-ViT [49], TNT [15], CaiT [33],
 80 DeepViT [53], CPVT [6], LocalViT [25] and CrossViT [5] focus on improving the architecture design
 81 of ViT. Another line of research proposes to integrate the inductive bias of CNN into Transformers
 82 [43, 7, 48, 14]. There are also attempts to adapt ViT for other vision tasks (e.g., object detection,
 83 semantic segmentation, etc.) [27, 38, 11, 18, 50, 52]. The most majority of these concurrent works
 84 represent each image with a fix number of tokens. To the best of our knowledge, we are the first to
 85 consider configuring token numbers conditioned on the inputs.

86 **Efficient deep networks.** Computational efficiency plays a critical role in real-world scenarios, where
 87 the executed computation translates into power consumption, carbon emission or latency. A number
 88 of works have been done on reducing the computational cost of CNNs [20, 29, 19, 47, 51, 28, 30].
 89 However, designing efficient vision Transformers is still an under-explored topic. T2T-ViT [49]
 90 proposes a light-weighted tokens-to-token module and obtains a competitive accuracy-parameter
 91 trade-off compared to MobileNetV2 [29]. LeViT [14] accelerates the inference of Transformer
 92 models by involving convolutional layers. Swin Transformer [27] introduces an efficient shifted

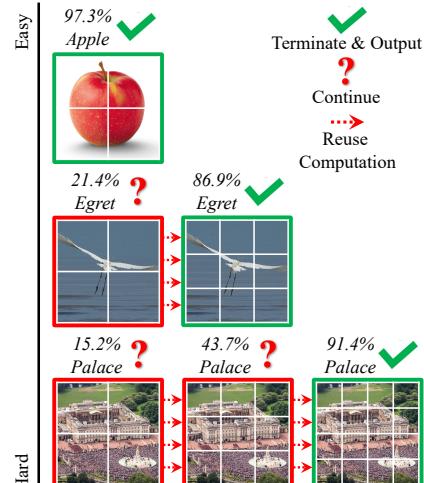


Figure 1: Examples for DVT.

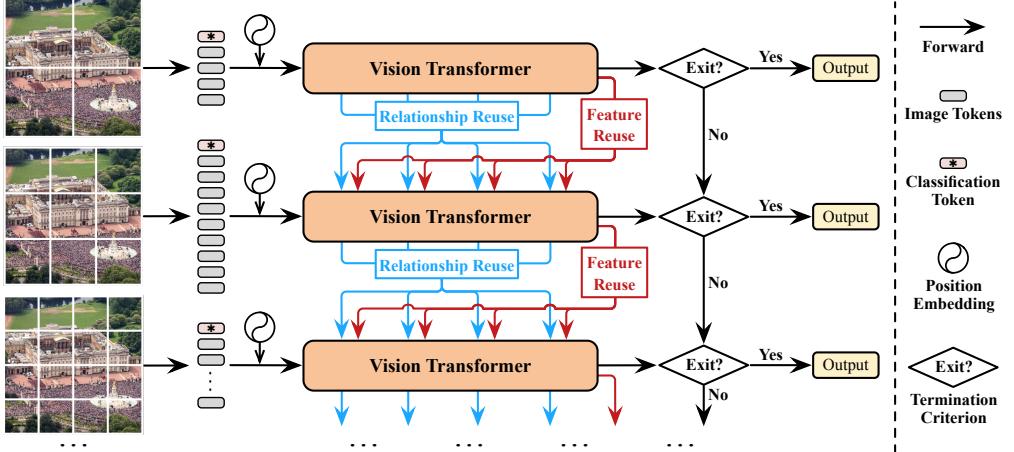


Figure 2: An overview of *Dynamic Vision Transformers* (DVT). Under the objective of configuring proper token numbers conditioned on the inputs, we cascade multiple Transformers with increasing number of tokens. At test time, they are sequentially activated until a convincing prediction (e.g. sufficiently confident) has been obtained or the final model has been inferred. The feature and relationship reuse mechanisms allow reusing computation across different Transformers.

93 window-based approach in multi-stage vision Transformers. Compared to these models with fixed
 94 computational graphs, the proposed DVT framework improves the efficiency by adaptively changing
 95 the architecture of the network on a per-sample basis.

96 **Dynamic models.** Designing dynamic architectures is an effective approach for efficient deep
 97 learning [16]. In the context of recognition tasks, MSDNet and its variants [21, 46, 24] develop a
 98 multi-classifier CNN architecture to perform early exiting for easy samples. Another type of dynamic
 99 CNNs skips redundant layers [35, 39, 45] or channels [26] conditioned on the inputs. Besides,
 100 the spatial adaptive paradigm [13, 4, 41, 36, 40] has been proposed for efficient image and video
 101 recognition. Although these works are related to DVT on the spirit of adaptive computation, they are
 102 developed based on CNN, while DVT is tailored for vision Transformers.

103 3 Dynamic Vision Transformer

104 Vision Transformers [10, 15, 32, 49] split each 2D image into 1D tokens, while model their long
 105 range interaction with the self-attention mechanism [34]. As aforementioned, to correctly recognize
 106 some ‘‘hard’’ images and achieve high accuracy, the number of tokens usually needs to be large,
 107 leading to the quadratically grown computational cost. However, ‘‘easier’’ images that make up the
 108 bulk of the datasets typically require far fewer tokens and much less costs (as shown in Table 1).
 109 Inspired by this observation, we propose a *Dynamic Vision Transformers* (DVT), aiming to improve
 110 the computational efficiency of Transformers via adaptively reducing the number of representative
 111 tokens for each input.

112 In specific, we propose to deploy multiple Transformers trained with increasing number of tokens,
 113 such that one can sequentially activate them for each test image until obtaining a convincing prediction
 114 (e.g., with sufficient confidence). The computation is allocated unevenly across different samples for
 115 improving the overall efficiency. It is worth noting that, if all the Transformers are learned separately,
 116 the computation performed by upstream models will simply be abandoned once a downstream
 117 Transformer is activated, resulting in considerable inefficiency. To alleviate this problem, we introduce
 118 the efficient feature and relationship reuse mechanisms.

119 3.1 Overview

120 **Inference.** We start by describing the inference procedure of DVT, which is shown in Figure 2. For
 121 each test sample, we first coarsely represent it using a small number of 1D token embeddings. This
 122 can be achieved by either straightforwardly flattening the split image patches [10, 15] or leveraging
 123 techniques like the tokens-to-token module [49]. We infer a vision Transformer with these few
 124 tokens to obtain a quick prediction. This process enjoys high efficiency since the computation cost of
 125 Transformers grows quadratically with respect to token number. Then the prediction will be evaluated

126 with certain criterion to determine whether it is reliable enough to be retrieved immediately. In this
 127 paper, early-termination is performed when the model is sufficiently confident (details in Section 3.3).

128 Once the prediction fails to meet the termination criterion, the original input image will be split
 129 into more tokens for more accurate but computationally more expensive inference. Note that, here
 130 the dimension of each token embedding remains unchanged, while the number of tokens increases,
 131 enabling more fine-grained representation. An additional Transformer with the same architecture
 132 as the previous one but different parameters will be activated. By design, this stage trades off
 133 computation for higher accuracy on some “difficult” test samples. To improve the efficiency, the
 134 new model can reuse the previously learned features and relationships, which will be introduced in
 135 Section 3.2. Similarly, after obtaining a new prediction, the termination criterion will be applied, and
 136 the above procedure will proceed until the sample exits or the final Transformer has been inferred.

137 **Training.** For training, we simply train DVT to produce correct predictions at all exits (i.e., each
 138 with the corresponding number of tokens). Formally, the optimization objective is

$$\text{minimize } \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{train}}} \left[\sum_i L_{\text{CE}}(\mathbf{p}_i, y) \right], \quad (1)$$

139 where (\mathbf{x}, y) denote a sample in the training set $\mathcal{D}_{\text{train}}$ and its corresponding label. We adopt the
 140 standard cross-entropy loss function $L_{\text{CE}}(\cdot)$, while \mathbf{p}_i denotes the softmax prediction probability
 141 output by the i^{th} exit. We find that such a simple training objective works well in practice.

142 **Transformer backbone.** DVT is proposed as a general and flexible framework. It can be built on
 143 top of most existing vision Transformers like ViT [10], DeiT [32] and T2T-ViT [49] to improve their
 144 efficiency. The architecture of Transformers simply follows the implementation of these backbones.

145 3.2 Feature and Relationship Reuse

146 An important challenge to develop our DVT approach is how to facilitate the *reuse* of computation.
 147 That is, once a downstream Transformer with more tokens is inferred, it is obviously inefficient if
 148 the computation performed in previous models is abandoned. The upstream models, although being
 149 based on smaller number of input tokens, are trained with the same objective, and have extracted
 150 valuable information for fulfilling the task. Therefore, we propose two mechanisms to reuse the
 151 learned deep features and self-attention relationships. Both of them are able to improve the test
 152 accuracy significantly by involving minimal extra computational cost.

153 **Background.** For the ease of introduction, we first revisit the basic formulation of vision Trans-
 154 formers. The Transformer encoders consist of alternatively stacked multi-head self-attention (MSA)
 155 and multi-layer perceptron (MLP) blocks [34, 10]. The layer normalization (LN) [2] and residual
 156 connection [17] are applied before and after each block, respectively. Let $\mathbf{z}_l \in \mathbb{R}^{N \times D}$ denote the
 157 output of the l^{th} Transformer layer, where N is the number of tokens for each sample, and D is the
 158 dimension of each token. Note that $N = HW + 1$, which corresponds to $H \times W$ patches of the
 159 original image and a single learnable classification token. Formally, we have

$$\mathbf{z}'_l = \text{MSA}(\text{LN}(\mathbf{z}_{l-1})) + \mathbf{z}_{l-1}, \quad l \in \{1, \dots, L\}, \quad (2)$$

$$\mathbf{z}_l = \text{MLP}(\text{LN}(\mathbf{z}'_l)) + \mathbf{z}'_l, \quad l \in \{1, \dots, L\}, \quad (3)$$

160 where L is the total number of layers in the Transformer. The classification token in \mathbf{z}_L will be fed
 161 into a LN layer followed by a fully-connected layer for the final prediction. For simplicity, here we
 162 omit the details on the position embedding, which is unrelated to our main idea. No modification is
 163 performed on it in addition to the configurations of backbones.

164 **Feature reuse.** All the Transformers in DVT share the same goal of extracting discriminative
 165 representations for accurate recognition. Therefore, it is straightforward that downstream models
 166 should be learned on the basis of previously obtained deep features, rather than extracting features
 167 from scratch. The former is more efficient since the computation performed in an upstream model
 168 contributes to both itself and the successive models. To implement this idea, we propose a feature
 169 reuse mechanism (see: Figure 3). In specific, we leverage the image tokens output by the final layer
 170 of the upstream Transformer, i.e., \mathbf{z}_L^{up} , to learn a layer-wise embedding \mathbf{E}_l for the downstream model:

$$\mathbf{E}_l = f_l(\mathbf{z}_L^{\text{up}}) \in \mathbb{R}^{N \times D'}. \quad (4)$$

172 Herein, $f_l: \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times D'}$ consists of a sequence of operations starting with a LN-MLP ($\mathbb{R}^D \rightarrow \mathbb{R}^{D'}$),
 173 which introduces nonlinearity and allows more flexible transformations. Then the image tokens are

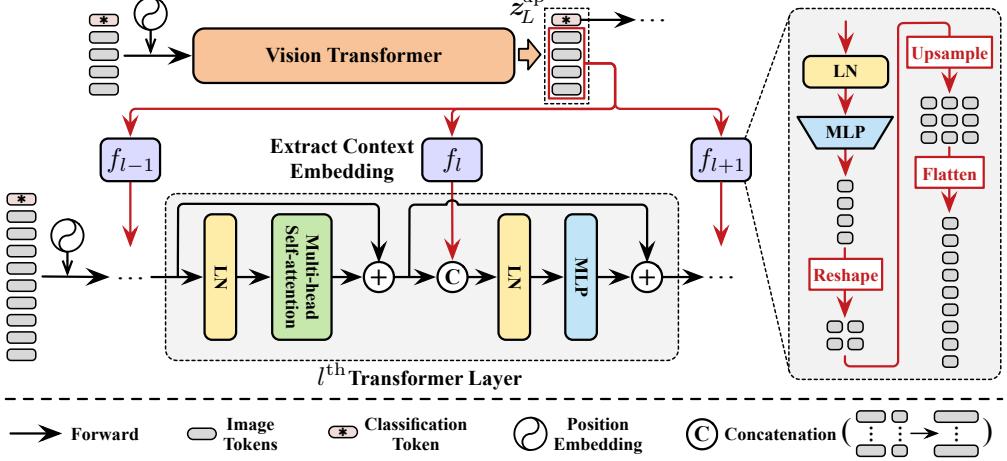


Figure 3: Illustration of the feature reuse mechanism. A layer-wise context embedding is learned based on the final representations output by the upstream model, i.e., z_L^{up} , and integrated into the MLP block of each downstream Transformer layer.

174 reshaped to the corresponding locations in the original image, upsampled and flattened to match the
 175 token number of the downstream model. Typically, we use a small D' for an efficient f_l .
 176 Consequently, the embedding E_l is injected into the downstream model, providing prior knowledge
 177 on recognizing the input image. Formally, we replace Eq. (3) by:

$$z_l = \text{MLP}(\text{LN}(\text{Concat}(z'_l, E_l))) + z'_l, \quad l \in \{1, \dots, L\}, \quad (5)$$

178 where E_l is concatenated with the intermediate tokens z'_l . We simply increase the dimension of LN
 179 and the first layer of MLP from D to $D+D'$. Since E_l is based on the upstream outputs z_L^{up} that have
 180 less tokens than z'_l , it actually concludes the context information of the input image for each token in
 181 z'_l . Therefore, we name E_l as the *context embedding*. Besides, we do not reuse the classification token
 182 and pad zero for it in Eq. (5), which we empirically find beneficial for the performance. Intuitively,
 183 Eqs. (4) and (5) allow training the downstream model to flexibly exploit the information within z_L^{up}
 184 on a per-layer basis, under the objective of minimizing the final recognition loss (Eq. (1)). This
 185 feature reuse formulation can also be interpreted as implicitly enlarging the depth of the model.

186 **Relationship reuse.** A prominent advantage of vision Transformers is that their self-attention
 187 blocks enable integrating information across the entire image, which effectively models
 188 the long-range dependencies in the data. Typically, the models need to learn a group of attention
 189 maps at each layer to describe the relationships among tokens. Apart from the deep features
 190 mentioned above, the downstream models also have access to the self-attention maps produced
 191 in previous models. We argue that these learned relationships are also capable of being reused
 192 to facilitate the learning of downstream Transformers.

206 Given the input representation z_l , the self-attention is performed as follows. First, the query, key and
 207 value matrices Q_l , K_l and V_l are computed via linear projections:

$$Q_l = z_l W_l^Q, \quad K_l = z_l W_l^K, \quad V_l = z_l W_l^V, \quad (6)$$

208 where \mathbf{W}_l^Q , \mathbf{W}_l^K and \mathbf{W}_l^V are weight matrices. Then the attention map is calculated by a scaled
 209 dot-product operation with softmax to aggregate the values of all tokens, namely

$$\text{Attention}(\mathbf{z}_l) = \text{Softmax}(\mathbf{A}_l)\mathbf{V}_l, \quad \mathbf{A}_l = \mathbf{Q}_l\mathbf{K}_l^\top / \sqrt{d}. \quad (7)$$

210 Here d is the hidden dimension of \mathbf{Q} or \mathbf{K} , and $\mathbf{A}_l \in \mathbb{R}^{N \times N}$ denotes the logits of the attention
 211 map. Note that we omit the details on the multi-head attention mechanism for clarity, where \mathbf{A}_l may
 212 include multiple attention maps. Such a simplification does not affect the description of our method.

213 For relationship reuse, we first concatenate the attention logits produced by all layers of the upstream
 214 model (i.e., \mathbf{A}_l^{up} , $l \in \{1, \dots, L\}$):

$$\mathbf{A}^{\text{up}} = \text{Concat}(\mathbf{A}_1^{\text{up}}, \mathbf{A}_2^{\text{up}}, \dots, \mathbf{A}_L^{\text{up}}) \in \mathbb{R}^{N_{\text{up}} \times N_{\text{up}} \times N_{\text{up}}^{\text{Att}}}, \quad (8)$$

215 where N_{up} and $N_{\text{up}}^{\text{Att}}$ denote the number of tokens and all attention maps in the upstream model,
 216 respectively. Typically, we have $N_{\text{up}}^{\text{Att}} = N^H L$, where N^H is the number of heads for the multi-head
 217 attention and L is the number of layers. Then the downstream Transformer learns attention maps by
 218 leveraging both its own tokens and \mathbf{A}^{up} simultaneously. Formally, we replace Eq. (7) by

$$\text{Attention}(\mathbf{z}_l) = \text{Softmax}(\mathbf{A}_l + r_l(\mathbf{A}^{\text{up}}))\mathbf{V}_l, \quad \mathbf{A}_l = \mathbf{Q}_l\mathbf{K}_l^\top / \sqrt{d}, \quad (9)$$

219 where $r_l(\cdot)$ is a transformation network that integrates the information provided by \mathbf{A}^{up} to refine the
 220 downstream attention logits \mathbf{A}_l . The architecture of $r_l(\cdot)$ is presented in Figure 4, which includes a
 221 MLP for nonlinearity followed by an upsample operation to match the size of attention maps. For
 222 multi-head attention, the output dimension of the MLP will be set to the number of heads.

223 Notably, Eq. (9) is a simple but flexible formulation. For one
 224 thing, each self-attention block in the downstream model has
 225 access to all the upstream attention heads in both shallow and
 226 deep layers, and hence can be trained to leverage multi-level
 227 relationship information on its own basis. For another, as the
 228 newly generated attention maps and the reused relationships
 229 are combined in logits, their relative importance can be auto-
 230 matically learned by adjusting the magnitude of logits. It is
 231 also worth noting that the regular upsample operation cannot
 232 be directly applied in $r_l(\cdot)$. To illustrate this issue, we take
 233 upsampling a $HW \times HW$ ($H=W=2$) attention map to $H'W' \times H'W'$ ($H'=W'=3$) for example
 234 in Figure 5. Since each of its rows and columns corresponds to $H \times W$ image tokens, we reshape the
 235 rows or columns back to $H \times W$, scale them to $H' \times W'$, and then flatten them to $H'W'$ vectors.

236 3.3 Adaptive Inference

237 As aforementioned, the proposed DVT framework progressively increases the number of tokens for
 238 each test sample and performs early-termination, such that “easy” and “hard” images can be processed
 239 using varying tokens with uneven computation cost, improving the overall efficiency. Specifically, at
 240 the i^{th} exit that produces the softmax prediction \mathbf{p}_i , the largest entry of \mathbf{p}_i , i.e., $\max_j p_{ij}$ (defined as
 241 confidence [21, 46, 41]), is compared with a threshold η_i . If $\max_j p_{ij} \geq \eta_i$, the inference will stop
 242 by adopting \mathbf{p}_i as the output. Otherwise, the image will be represented using more tokens to activate
 243 the downstream Transformer. We always adopt a zero-threshold for the final Transformer.

244 The values of $\{\eta_1, \eta_2, \dots\}$ are solved on the validation set. We assume a *budgeted batch classification*
 245 [21] setting, where DVT needs to recognize a set of samples \mathcal{D}_{val} within a given computational
 246 budget $B > 0$. Let $\text{Acc}(\mathcal{D}_{\text{val}}, \{\eta_1, \eta_2, \dots\})$ and $\text{FLOPs}(\mathcal{D}_{\text{val}}, \{\eta_1, \eta_2, \dots\})$ denote the accuracy and
 247 computational cost on \mathcal{D}_{val} when using the thresholds $\{\eta_1, \eta_2, \dots\}$. The optimal thresholds can be
 248 obtained by solving the following optimization problem:

$$\underset{\eta_1, \eta_2, \dots}{\text{maximize}} \quad \text{Acc}(\mathcal{D}_{\text{val}}, \{\eta_1, \eta_2, \dots\}), \quad \text{s.t. } \text{FLOPs}(\mathcal{D}_{\text{val}}, \{\eta_1, \eta_2, \dots\}) \leq B. \quad (10)$$

249 Due to the non-differentiability, we solve this problem with the genetic algorithm [42] in this paper.

250 4 Experiments

251 In this section, we empirically validate the proposed DVT on ImageNet [8] and CIFAR-10/100 [23].
 252 Ablation studies and visualization are presented on ImageNet to give a deeper understanding of our
 253 method. All the code and pre-trained models will be released upon the acceptance of this paper.

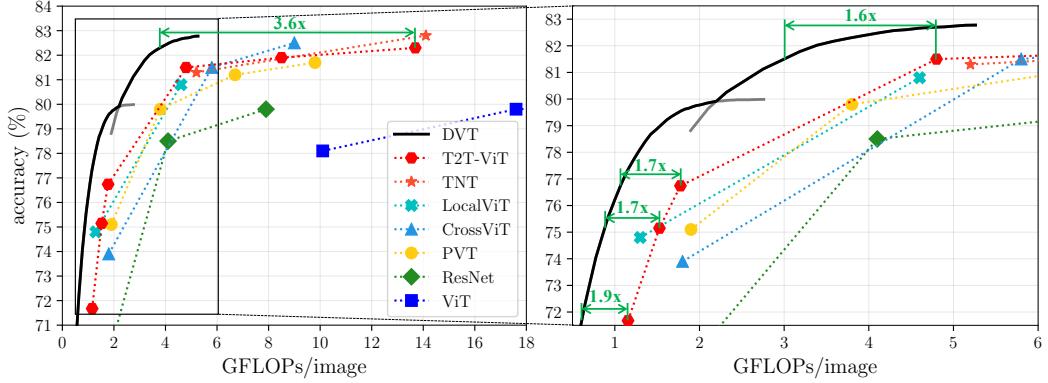


Figure 6: Top-1 accuracy v.s. GFLOPs on ImageNet. DVT is implemented on top of T2T-ViT-12/14.

Table 2: The practical speed of DVT.

Models	ImageNet (NVIDIA 2080Ti, bs=128)	
	Top-1 acc.	Throughput
T2T-ViT-7 DVT	71.68% 78.48% ($\uparrow 6.80\%$)	1574 img/s 1574 img/s
T2T-ViT-10 DVT	75.15% 79.74% ($\uparrow 4.59\%$)	1286 img/s 1286 img/s
T2T-ViT-12 DVT	76.74% 80.43% ($\uparrow 3.69\%$)	1121 img/s 1128 img/s
T2T-ViT-14 DVT	81.50% 81.50%	619 img/s 877 img/s ($\uparrow 1.42x$)
T2T-ViT-19 DVT	81.93% 81.93%	382 img/s 666 img/s ($\uparrow 1.74x$)

Table 3: Performance of DVT on CIFAR-10/100.

Models	CIFAR-10		CIFAR-100	
	Top-1 acc.	GFLOPs	Top-1 acc.	GFLOPs
T2T-ViT-10 DVT	97.21% 97.21%	1.53 0.50 ($\downarrow 3.1x$)	85.44% 85.45%	1.53 0.54 ($\downarrow 2.8x$)
T2T-ViT-12 DVT	97.45% 97.46%	1.78 0.52 ($\downarrow 3.4x$)	86.23% 86.26%	1.78 0.61 ($\downarrow 2.9x$)
T2T-ViT-14 DVT	98.19% 98.19%	4.80 0.77 ($\downarrow 6.2x$)	89.10% 89.11%	4.80 1.62 ($\downarrow 3.0x$)
T2T-ViT-19 DVT	98.43% 98.43%	8.50 1.44 ($\downarrow 5.9x$)	89.37% 89.38%	8.50 1.74 ($\downarrow 4.9x$)
T2T-ViT-24 DVT	98.53% 98.53%	13.69 1.49 ($\downarrow 9.2x$)	89.62% 89.63%	13.69 1.86 ($\downarrow 7.4x$)

254 **Datasets.** (1) ImageNet is a 1,000-class dataset from ILSVRC2012 [8], containing 1.2 million images
255 for training and 50,000 images for validation. (2) CIFAR-10/100 datasets [23] contain 32x32 colored
256 images in 10/100 classes. Both of them consist of 50,000 images for training and 10,000 images
257 for testing. For all the three datasets, we adopt the same data pre-processing and data augmentation
258 policy as [17, 22, 21]. In addition, we solve the confidence thresholds stated in Section 3.3 on the
259 training set, which we find achieves similar performance to adopting cross-validation.

260 **Backbones.** Our experiments are based on several state-of-the-art vision Transformers, namely
261 T2T-ViT-12 [49], T2T-ViT-14 [49], and DeiT-small (w/o distillation) [32]. Unless otherwise specified,
262 we deploy DVT with three exits, corresponding to representing the images as 7x7, 10x10 and 14x14
263 tokens¹. For fair comparisons, our implementation exploits the official code of the backbones, and
264 adopts exactly the same training hyper-parameters. More training details can be found in Appendix
265 A. The number of FLOPs is calculated using the *fvcore* toolkit provided by Facebook AI Research,
266 which is also used in Detectron2 [44], PySlowFast [12], and ClassyVision [1].

267 **Implementation details.** For feature reuse, the hidden size and output size of the MLP in $f_l(\cdot)$ are
268 set to 128 and 48. In relationship reuse, for implementation efficiency, we share the same hidden
269 state across the MLPs of all $r_l(\cdot)$, such that $r_l(\mathbf{A}^{\text{up}})$, $l \in \{1, \dots, L\}$ can be obtained at one time in
270 concatenation by implementing a single large MLP, whose hidden size and output size are $3N^H L$ and
271 $N^H L$. Note that N^H is the head number of multi-head attention and L is the layer number.

272 4.1 Main Results

273 **Results on ImageNet** are shown in Figures 6 and 7, where T2T-ViT [49] and DeiT [32] are im-
274 plemented as backbones respectively. As stated in Section 3.3, we vary the average computational
275 budget, solve the confidence thresholds, and evaluate the corresponding validation accuracy. The
276 performance of DVT is plotted in gray curves, with the best accuracy under each budget plotted in
277 black curves. We also compare our method with several highly competitive baselines, i.e., TNT [15],
278 LocalViT [25], CrossViT [5], PVT [38], ViT [10] and ResNet [17]. It can be observed that DVT
279 consistently reduces the computational cost of the backbones. For example, DVT achieves the 82.3%
280 accuracy with 3.6x less FLOPs compared with the vanilla T2T-ViT. When the budget ranges among

¹Although 4x4 tokens are also used as an example in Section 1, we find starting with 7x7 is more efficient.

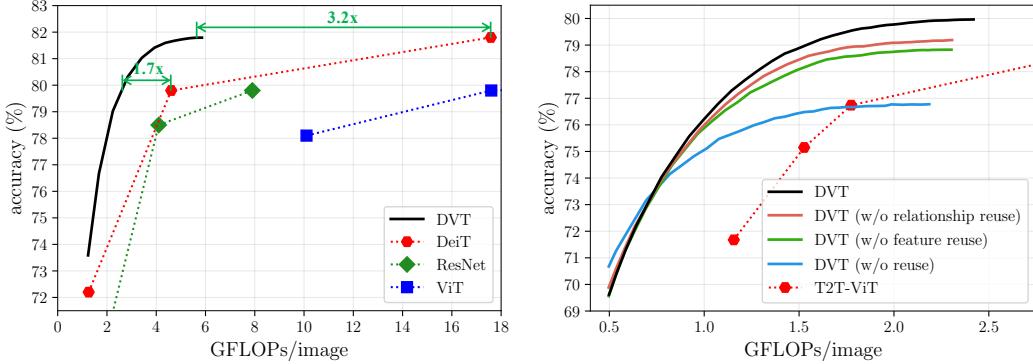


Figure 7: Performance of DeiT-based DVT on ImageNet. DeiT-small is used as the backbone. Figure 8: Performance of the DVT based on T2T-ViT-12 with and without the reuse mechanisms.

281 0.5-2 GFLOPs, DVT has \sim 1.7-1.9x less computation than T2T-ViT with the same performance.
282 Notably, our method can flexibly attain all the points on each curve by simply adjusting the values of
283 confidence thresholds with a single DVT.

284 **Practical efficiency of DVT.** We test the actual speed of DVT on a NVIDIA 2080Ti GPU under a
285 batch inference setting, where a mini-batch of data is fed into the model at a time. After inferring each
286 Transformer, the samples that meet the early-termination criterion will exit, with the remaining images
287 fed into the downstream Transformer. The results are presented in Table 2. Here we adopt a two-exit
288 DVT based on T2T-ViT-12 using 7x7 and 14x14 tokens, which we find more efficient in practice. All
289 other implementation details remain unchanged. One can observe that DVT improves the accuracy of
290 small models (T2T-ViT-7/10/12) by 3.7-6.8% with the same inference speed, while accelerates the
291 inference of the large T2T-ViT-14/19 models by 1.4-1.7x without sacrificing performance.

292 **Results on CIFAR** are presented in Table 3. Following the common practice [10, 49, 15, 32], we
293 resize the CIFAR images to 224x224, and fine-tune the T2T-ViT and DVT models in Figure 6. The
294 official code and training configurations provided by [49] are utilized. We report the computation
295 cost of DVT when it achieves the competitive performance with baselines. Our proposed method is
296 shown to consume \sim 3-9x less computation compared with T2T-ViT.

297 4.2 Ablation Study

298 **Effectiveness of feature and relationship reuse.** We conduct experiments by ablating one or both of
299 the reuse mechanisms. For a clear
300 comparison, we first deactivate the
301 early-termination, and report the
302 accuracy and GFLOPs corresponding
303 to each exit in Table 4. The
304 three-exit DVT based on T2T-ViT-
305 12 is considered. One can observe

Table 4: Effects of feature (F) and relationship (R) reuse. The percentages in brackets denote the additional computation compared to baselines involved by the reuse mechanisms.

Reuse	1 th Exit (7x7)		2 th Exit (10x10)		3 th Exit (14x14)			
	F	R	Top-1 acc.	GFLOPs	Top-1 acc.	GFLOPs	Top-1 acc.	GFLOPs
✓			70.33%	0.47	73.54%	1.37	76.74%	3.15
✓		✓	69.42%	0.47	75.31%	1.43 _(4.4%)	79.21%	3.31 _(5.1%)
✓	✓		69.03%	0.47	75.34%	1.41 _(2.9%)	78.86%	3.34 _(6.0%)
✓	✓	✓	69.04%	0.47	75.65%	1.46 _(6.6%)	80.00%	3.50 _(11.1%)

306 that both the two reuse mechanisms are able to significantly boost the accuracy of DVT at the 2th
307 and 3th exits with at most 6% additional computation, while they are compatible with each other to
308 further improve the performance. We also find that involving computation reusing slightly hurts the
309 accuracy at the 1th exit, which may be attributed to the compromise made by the first Transformer
310 for downstream models. However, once the early-termination is adopted, this difference only results
311 in trivial disadvantage when the computational budget is very small, as shown in Figure 8. DVT
312 outperforms the baseline significantly in most cases.

315 **Design choices for the reuse mechanisms.** Here we study the design of the feature and relationship
316 reuse mechanisms. For experimental efficiency, we consider a two-exit DVT based on T2T-ViT-12
317 using 7x7 and 10x10 tokens, while enlarge the batch size and the initial learning rate by 4 times. Such
318 a training setting slightly degrades the accuracy of DVT, but it is still reliable to reflect the difference
319 between different design variants. We deactivate early-termination, and report the performance of
320 each exit. Notably, as the FLOPs of 1th exit remain unchanged (i.e., 0.47G), we do not present it.

Table 5: Ablation studies for feature reuse.

Ablation	1 th Exit (7x7) Top-1 acc.	2 th Exit (10x10) Top-1 acc.	GFLOPs
w/o reuse	70.08%	73.61%	1.37
Layer-wise feature reuse	69.84%	74.31%	1.43
Reuse classification token	69.79%	74.70%	1.43
Remove $f_l(\cdot)$, $l \geq 2$	69.33%	74.73%	1.38
Remove LN in $f_l(\cdot)$	69.63%	75.05%	1.42
Ours	69.44%	75.23%	1.43

“Easy”

“Hard”

“Easy”

“Hard”

“Easy”

“Hard”

“Easy”

“Hard”

Figure 9: Visualization of the “easy” and “hard” samples in DVT.

321 We consider four variants of feature reuse in Table 5: (1) reusing features from the corresponding
 322 upstream layer instead of the final layer; (2) reusing classification token; (3) only performing feature
 323 reuse in the first layer of the downstream model; (4) removing the LN in $f_l(\cdot)$. One can see that taking
 324 final tokens of the upstream model and reusing them in each downstream layer are both important.

325 Ablation results for relationship reuse are presented in Table 6. We consider four variants as well:
 326 (1) only reusing the attention logits from the corresponding upstream layer; (2) only reusing the
 327 attention logits from the final upstream layer; (3) replacing the MLP by a linear layer in $r_l(\cdot)$; (4)
 328 adopting naive upsample operation instead of what is shown in Figure 5; The results indicate that it is
 329 beneficial to enable each downstream layer to flexibly reuse all upstream attention logits. Besides,
 330 naive upsampling significantly hurts the performance.

331 **Early-termination.** We vary the criterion for
 332 adaptive inference and report the accuracy under
 333 several computational budgets in Table 7. Two
 334 variants are considered: (1) adopting the entropy
 335 of softmax prediction to determine whether to
 336 exit [31]; (2) performing random exiting with
 337 the same exit proportion as DVT. The simple
 338 but effective confidence-based criterion achieves better performance than both of them.

339 4.3 Visualization

340 Figure 9 shows the images that are first correctly classified at the 1th and 3th exits of the DVT (T2T-
 341 ViT-12). The former are recognized as “easy” samples, while the later are considered to be “hard”.
 342 One can observe that “easy” samples usually depict the recognition objectives in clear and canonical
 343 poses and sufficiently large resolution. On the contrary, “hard” samples may contain complex scenes
 344 and non-typical poses or only include a small part of the objects, and require a finer representation
 345 using more tokens. Figure 10 presents the numbers of images that exit at different exits when the
 346 computational budget increases. The plot shows that the accuracy of DVT is significantly improved
 347 with more images exiting later, which is achieved by changing the confidence thresholds online.

348 5 Conclusion

349 In this paper, we sought to optimally configure a proper number of tokens for each individual image in
 350 vision Transformers, and hence proposed the *Dynamic Vision Transformer* (DVT) framework. DVT
 351 processes each test input by sequentially activating multiple Transformers using increasing tokens,
 352 until an appropriate token number is reached (measured by the corresponding prediction confidence).
 353 We further introduce the feature and relationship reuse mechanisms to facilitate efficient computation
 354 reuse. Extensive experiments indicate that DVT significantly improves the computational efficiency
 355 on top of the state-of-the-art vision Transformers, both theoretically and empirically.

Table 6: Ablation studies for relationship reuse.

Ablation	1 th Exit (7x7) Top-1 acc.	2 th Exit (10x10) Top-1 acc.	GFLOPs
w/o reuse	70.08%	73.61%	1.37
Layer-wise relationship reuse	69.63%	73.89%	1.38
Reuse final-layer relationships	69.25%	74.31%	1.39
MLP→Linear	69.20%	73.84%	1.38
Naive upsample	69.60%	73.34%	1.41
Ours	69.50%	74.91%	1.41

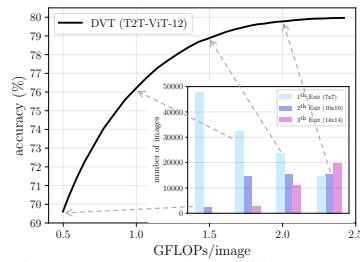


Figure 10: Numbers of images exiting at different exits with varying computational budgets.

Table 7: Comparisons of early-termination criterions. The accuracy under each budget is reported.

Ablation	Top-1 acc.			
	0.75G	1.00G	1.25G	1.50G
Randomly Exit	70.19%	71.66%	72.61%	73.59%
Entropy-based	73.41%	75.21%	77.08%	78.40%
Confidence-based (ours)	73.70%	76.22%	77.89%	78.89%

356 **References**

- 357 [1] A. Adcock, V. Reis, M. Singh, Z. Yan, L. van der Maaten, K. Zhang, S. Motwani, J. Guerin,
358 N. Goyal, I. Misra, L. Gustafson, C. Changhan, and P. Goyal. Classy vision. <https://github.com/facebookresearch/ClassyVision>, 2019.
- 360 [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- 362 [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
363 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
364 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
365 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
366 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
367 Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle,
368 M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *NeurIPS*, volume 33, pages 1877–
369 1901. Curran Associates, Inc., 2020.
- 370 [4] Shijie Cao, Lingxiao Ma, Wencong Xiao, Chen Zhang, Yunxin Liu, Lintao Zhang, Lanshun Nie,
371 and Zhi Yang. Seernet: Predicting convolutional neural network feature-map sparsity through
372 low-bit quantization. In *CVPR*, pages 11216–11225, 2019.
- 373 [5] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision
374 transformer for image classification. *arXiv preprint arXiv:2103.14899*, 2021.
- 375 [6] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and
376 Chunhua Shen. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021.
- 378 [7] Stéphane d’Ascoli, Hugo Touvron, Matthew Leavitt, Ari Morcos, Giulio Biroli, and Levent
379 Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. *arXiv
380 preprint arXiv:2103.10697*, 2021.
- 381 [8] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
382 image database. In *ICML*, pages 248–255, 2009.
- 383 [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training
384 of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages
385 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- 386 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
387 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
388 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
389 recognition at scale. In *ICLR*, 2021.
- 390 [11] Alaaeldin El-Nouby, Natalia Neverova, Ivan Laptev, and Hervé Jégou. Training vision trans-
391 formers for image retrieval. *arXiv preprint arXiv:2102.05644*, 2021.
- 392 [12] Haoqi Fan, Yanghao Li, Bo Xiong, Wan-Yen Lo, and Christoph Feichtenhofer. Pyslowfast.
393 <https://github.com/facebookresearch/slowfast>, 2020.
- 394 [13] Michael Figurnov, Maxwell D Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry Vetrov,
395 and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. In *CVPR*,
396 pages 1039–1048, 2017.
- 397 [14] Ben Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou,
398 and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. *arXiv
399 preprint arXiv:2104.01136*, 2021.
- 400 [15] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in
401 transformer. *arXiv preprint arXiv:2103.00112*, 2021.
- 402 [16] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic
403 neural networks: A survey. *arXiv preprint arXiv:2102.04906*, 2021.

- 404 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 405 recognition. In *CVPR*, pages 770–778, 2016.
- 406 [18] Shuteng He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. Transreid: Transformer-
 407 based object re-identification. *arXiv preprint arXiv:2102.04378*, 2021.
- 408 [19] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan,
 409 Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3.
 410 In *ICCV*, pages 1314–1324, 2019.
- 411 [20] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias
 412 Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural
 413 networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- 414 [21] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q
 415 Weinberger. Multi-scale dense networks for resource efficient image classification. In *ICLR*,
 416 2018.
- 417 [22] Gao Huang, Zhuang Liu, Geoff Pleiss, Laurens Van Der Maaten, and Kilian Weinberger.
 418 Convolutional networks with dense connectivity. *IEEE Transactions on Pattern Analysis and
 419 Machine Intelligence*, pages 1–1, 2019.
- 420 [23] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
 421 Technical report, Citeseer, 2009.
- 422 [24] Hao Li, Hong Zhang, Xiaojuan Qi, Ruigang Yang, and Gao Huang. Improved techniques for
 423 training adaptive deep networks. In *ICCV*, pages 1891–1900, 2019.
- 424 [25] Yawei Li, Kai Zhang, Jiezhang Cao, Radu Timofte, and Luc Van Gool. Localvit: Bringing
 425 locality to vision transformers. *arXiv preprint arXiv:2104.05707*, 2021.
- 426 [26] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In *NeurIPS*, pages
 427 2181–2191, 2017.
- 428 [27] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
 429 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint
 430 arXiv:2103.14030*, 2021.
- 431 [28] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines
 432 for efficient cnn architecture design. In *ECCV*, pages 116–131, 2018.
- 433 [29] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen.
 434 Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, pages 4510–4520, 2018.
- 435 [30] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural
 436 networks. In *ICML*, volume 97, pages 6105–6114, 2019.
- 437 [31] Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. Branchynet: Fast inference
 438 via early exiting from deep neural networks. In *ICPR*, pages 2464–2469. IEEE, 2016.
- 439 [32] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
 440 Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv
 441 preprint arXiv:2012.12877*, 2020.
- 442 [33] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou.
 443 Going deeper with image transformers. *arXiv preprint arXiv:2103.17239*, 2021.
- 444 [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 445 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg,
 446 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *NeurIPS*, volume 30.
 447 Curran Associates, Inc., 2017.
- 448 [35] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In
 449 *ECCV*, pages 3–18, 2018.

- 450 [36] Thomas Verelst and Tinne Tuytelaars. Dynamic convolutions: Exploiting spatial sparsity for
 451 faster inference. In *CVPR*, pages 2320–2329, 2020.
- 452 [37] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
 453 GLUE: A multi-task benchmark and analysis platform for natural language understanding. In
 454 *ICLR*, 2019.
- 455 [38] Wenhui Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping
 456 Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction
 457 without convolutions. *arXiv preprint arXiv:2102.12122*, 2021.
- 458 [39] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E Gonzalez. Skipnet: Learning
 459 dynamic routing in convolutional networks. In *ECCV*, pages 409–424, 2018.
- 460 [40] Yulin Wang, Zhaoxi Chen, Haojun Jiang, Shiji Song, Yizeng Han, and Gao Huang. Adaptive
 461 focus for efficient video recognition. *arXiv preprint arXiv:2105.03245*, 2021.
- 462 [41] Yulin Wang, Kangchen Lv, Rui Huang, Shiji Song, Le Yang, and Gao Huang. Glance and focus:
 463 a dynamic approach to reducing spatial redundancy in image classification. In *NeurIPS*, 2020.
- 464 [42] Darrell Whitley. A genetic algorithm tutorial. *Statistics and computing*, 4(2):65–85, 1994.
- 465 [43] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang.
 466 CvT: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021.
- 467 [44] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2.
 468 <https://github.com/facebookresearch/detectron2>, 2019.
- 469 [45] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen
 470 Grauman, and Rogerio Feris. Blockdrop: Dynamic inference paths in residual networks. In
 471 *CVPR*, pages 8817–8826, 2018.
- 472 [46] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive
 473 networks for efficient inference. In *CVPR*, pages 2369–2378, 2020.
- 474 [47] Le Yang, Haojun Jiang, Ruojin Cai, Yulin Wang, Shiji Song, Gao Huang, and Qi Tian. CondenseNet v2: Sparse feature reactivation for deep networks. In *CVPR*, 2021.
- 476 [48] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating
 477 convolution designs into visual transformers. *arXiv preprint arXiv:2103.11816*, 2021.
- 478 [49] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and
 479 Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet.
 480 *arXiv preprint arXiv:2101.11986*, 2021.
- 481 [50] Fangao Zeng, Bin Dong, Tiancai Wang, Cheng Chen, Xiangyu Zhang, and Yichen Wei. Motr:
 482 End-to-end multiple-object tracking with transformer. *arXiv preprint arXiv:2105.03247*, 2021.
- 483 [51] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient
 484 convolutional neural network for mobile devices. In *CVPR*, pages 6848–6856, 2018.
- 485 [52] Moju Zhao, Kei Okada, and Masayuki Inaba. Trtr: Visual tracking with transformer. *arXiv
 486 preprint arXiv:2105.03817*, 2021.
- 487 [53] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou,
 488 and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*,
 489 2021.

490 **Checklist**

- 491 1. For all authors...
- 492 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
493 contributions and scope? **[Yes]**
- 494 (b) Did you describe the limitations of your work? **[Yes]** See Appendix C.
- 495 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See
496 Appendix C.
- 497 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
498 them? **[Yes]** We have read them and confirmed. This paper conforms to them.
- 499 2. If you are including theoretical results...
- 500 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 501 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 502 3. If you ran experiments...
- 503 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
504 perimental results (either in the supplemental material or as a URL)? **[Yes]** We have
505 introduced our implementation, the source of the training code, the hyper-parameters
506 and the datasets in details. See Section 4 and Appendix A. Our official code and
507 pre-trained models will be made public upon the acceptance of this paper.
- 508 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
509 were chosen)? **[Yes]** See Section 4 and Appendix A.
- 510 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
511 periments multiple times)? **[No]** Our main results are presented on the large scale
512 ImageNet dataset (with 1.2M training samples) using a long training schedule (300
513 epochs). This experimental protocol leads to satble performance (standard deviation \approx
514 0.15%), while our proposed method usually boosts the accuracy by $\geq 1.5\%$. Therefore,
515 we do not report the error bars, and we believe that this will not affects our contributions.
516 Importantly, this is actually a common practice on ImageNet.
- 517 (d) Did you include the total amount of compute and the type of resources used (e.g., type
518 of GPUs, internal cluster, or cloud provider)? **[Yes]** See Appendix A.
- 519 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 520 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We use publically
521 available datasets and the officially released code on Github. All citations have been
522 added.
- 523 (b) Did you mention the license of the assets? **[No]** We use the public ImageNet and
524 CIFAR datasets. The non-commercial use has been permitted.
- 525 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
- 526 (d) Did you discuss whether and how consent was obtained from people whose data you're
527 using/curating? **[No]**
- 528 (e) Did you discuss whether the data you are using/curating contains personally identifiable
529 information or offensive content? **[No]**
- 530 5. If you used crowdsourcing or conducted research with human subjects...
- 531 (a) Did you include the full text of instructions given to participants and screenshots, if
532 applicable? **[N/A]**
- 533 (b) Did you describe any potential participant risks, with links to Institutional Review
534 Board (IRB) approvals, if applicable? **[N/A]**
- 535 (c) Did you include the estimated hourly wage paid to participants and the total amount
536 spent on participant compensation? **[N/A]**