

Make A Long Image Short: Adaptive Token Length for Vision Transformers

Qiqi Zhou^{1,*} and Yichen Zhu^{2,✉}

¹ Shanghai University of Electric Power, Shanghai, China

² Midea Group, Shanghai, China
{zhouqq31, zhuyc25}@midea.com

Abstract. The vision transformer is a model that breaks down each image into a sequence of tokens with a fixed length and processes them similarly to words in natural language processing. Although increasing the number of tokens typically results in better performance, it also leads to a considerable increase in computational cost. Motivated by the saying "A picture is worth a thousand words," we propose an innovative approach to accelerate the ViT model by shortening long images. Specifically, we introduce a method for adaptively assigning token length for each image at test time to accelerate inference speed. First, we train a Resizable-ViT (ReViT) model capable of processing input with diverse token lengths. Next, we extract token-length labels from ReViT that indicate the minimum number of tokens required to achieve accurate predictions. We then use these labels to train a lightweight Token-Length Assigner (TLA) that allocates the optimal token length for each image during inference. The TLA enables ReViT to process images with the minimum sufficient number of tokens, reducing token numbers in the ViT model and improving inference speed. Our approach is general and compatible with modern vision transformer architectures, significantly reducing computational costs. We verified the effectiveness of our methods on multiple representative ViT models on image classification and action recognition.

Keywords: vision transformer · token compression.

1 Introduction

The transformer has achieved remarkable success in computer vision since the introduction of ViT [12]. It has demonstrated impressive performance compared to convolutional neural networks (CNNs) on various visual domains, including image classification [42,10], object detection [7,57], semantic segmentation [25], and action recognition [13,4], using both supervised and self-supervised [19,2] training configurations. Despite the development of ViT models, their deployment remains a challenge due to the high computational cost associated with them.

* Work done during internships at Midea Group.

✉ Corresponding authors.

Accelerating ViT is a crucial yet understudied area. While many techniques like pruning, distillation, and neural architecture search have been applied to accelerate CNNs, these cannot be directly applied to ViT due to significant differences between the models [31,29,35]. As the attention module in the transformer computes the fully-connected relations among all input patches [43], the computational cost becomes quadratic with respect to the length of the input sequence [9,3]. Consequently, the transformer can be computationally expensive, particularly for longer input sequences. In the ViT model, images are divided into a fixed number of tokens; following conventional practice [12], an image is represented by 16×16 tokens. We aim to reduce the computational complexity of ViT by reducing the number of tokens used to split the images. Our motivation is depicted in Figure 1, which shows three examples predicted by individually trained DeiT-S models [42] with different token lengths. The checkmark denotes correct prediction, and the cross denotes the wrong prediction. We observe that some "easy-to-classify" images only require a few tokens to determine their category accurately, while some images require more tokens to make the right prediction. These observations motivate us to reduce the computational complexity of the existing ViT model by accurately classifying the input using the minimum possible number of tokens.

In an ideal scenario, we would know the minimum number of tokens required to accurately predict an image, and we could train a model to assign the optimal token length to the ViT model. However, training multiple ViT models, each with a fixed token length, would be computationally infeasible. To address this, we propose a modification to the transformer architecture, changing it from "static" to "dynamic," enabling the ViT model to adaptively process images with varying token lengths. This dynamic transformer, called Resizable-ViT (ReViT), identifies the minimum token length required to achieve correct predictions for each image. We then train a lightweight Token-Length Assigner (TLA) to predict the appropriate token length for a given image, with the label obtained from the ReViT. Consequently, the ReViT can process images with lower computational costs based on the assigned token length.

The primary challenge of our approach is training the ReViT to enable the ViT model to process images of any size provided by the TLA. To tackle this

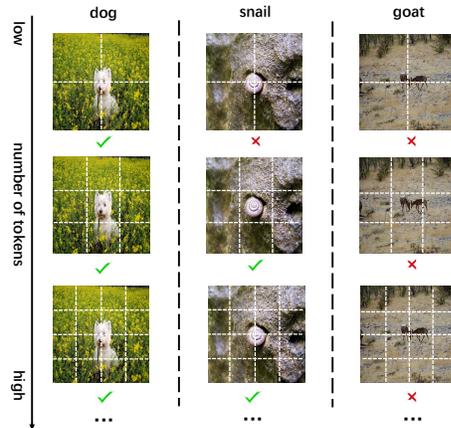


Fig. 1. The motivation for our approach. While some images (right) may need many tokens to predict their category, some images are easy to recognize. Thus, only a small number of tokens is sufficient to classify them correctly.

challenge, we introduce a token length-aware layer normalization that switches the normalization statistics for each type of token length, and a self-distillation module that enhances the model’s performance when using short token lengths in ReViT. Additionally, the ViT model needs to see the images with the corresponding token lengths beforehand to handle various token lengths effectively. However, as the number of predefined token-length choices increases, the training cost linearly increases. To overcome this, we introduce a parallel computing strategy for efficient training that makes the ReViT training almost as inexpensive as a vanilla ViT model’s training.

We showcase the efficacy of our approach on several prominent ViT models, such as DeiT [42] and LV-ViT [22] for image classification, and TimesFormer [4] for video recognition. Our experiments demonstrate that our method can significantly reduce computational costs while maintaining performance levels. For instance, we achieve a 50% acceleration in DeiT-S [42] model with an accuracy reduction of only 0.1%. On action recognition, the computational cost of TimesFormer [4] can be reduced up to 33% on Kinetic 400 with only a 0.5% loss in recognition accuracy.

2 Related Works

Vision transformer. ViT have recently gained much attention in computer vision due to their strong capability to model long-range relations. Many attempts have been made to integrate long-range modeling into CNNs, such as non-local networks [45,51], relation networks [21], among others. Vision Transformer (ViT)[12] introduced a set of pure Transformer backbones for image classification, and its follow-ups have soon modified the vision transformer to dominate many downstream tasks for computer vision, such as object detection[7,57], semantic segmentation [25], action recognition [4,13], 2D/3D human pose estimation [50,56], 3D object detection [33], and even self-supervision [19]. ViT has shown great potential to be an alternative backbone for convolutional neural networks.

Dynamic vision transformer. The over-parameterized model is known to have many attractive merits and can achieve better performance than smaller models. However, in real-world scenarios, computational efficiency is critical as executed computation is translated into power consumption or carbon emission. To address this issue, many works have attempted to reduce the computational cost of Convolutional Neural Networks (CNNs) through methods such as neural architecture search [24,6,61], knowledge distillation [59,55,20,60,58], and pruning [18,15].

Recent work has shift its attention to reduce the number of tokens used for inference, as the number of tokens can be a computational bottleneck to the vision transformer. There are two major approaches: unstructured token sparsification and structured token division. The majority of works, including PatchSlim [41], TokenSparse [36], GlobalEncoder [39], IA-RED [32], and Tokenlearner [37], focus

on the former. TokenLearner [37] uses an MLP to reduce the number of tokens. TokenPooling [30] merges tokens via a k-mean based algorithm. TokenMerge [5] calculates the token similarity and merges tokens via bipartite soft matching. .

They aim to remove uninformative tokens, such as those that learn features from the background of the image, thereby boosting inference speed by reserving only informative tokens. These approaches typically need to progressively reduce the number of tokens based on the inputs and can be performed either jointly with ViT training or afterward. However, pruning tokens sparsely can bring unstable training issues, especially when the model is huge [23].

The latter, which is known as unstructured token sparsification, is the most relevant work to our research. Wang et al.[46] proposed Dynamic Vision Transformer (DVT) to dynamically determine the number of patches required to divide an image. They employed a cascade of ViT models, with each ViT responsible for a specific token length. The cascade ViT model makes a sequential decision and stops inference for an input image if it has sufficient confidence in the prediction at the current token length. In contrast to DVT[46], our method is more practical and accessible, as it only requires a *single* ViT model. Additionally, we focus on how to *accurately* determine the minimum number of token lengths required in the transformer to provide correct predictions for each image.

3 Methodology

The vision transformers treat an image as a sentence by dividing the 2D image into 1D tokens and modeling the long-range dependencies between them using the multi-head self-attention mechanism. However, the self-attention is considered the computational bottleneck in the transformer model, as its computational cost increases quadratically with the number of incoming tokens. As mentioned earlier, our approach is motivated by the observation that many “easy-to-recognize” images do not require 16×16 tokens [12] to be correctly classified. Therefore, computational costs can be reduced by processing fewer tokens on “easy” images while using more tokens on “hard” images. It is worth noting that the key to a successful input-dependent token-adaptive ViT model is to determine precisely the minimum number of tokens required to accurately classify the image.

To achieve our goal, we propose a two-stage model training approach. In the first stage, we train a ViT model that can handle images with any predefined token lengths. Usually, a single ViT model can only handle one token length. We describe the model design and training strategy of this ViT model in detail in Section 3.2. In the second stage, we train a model to determine the appropriate token length for each image. We first obtain the token-length label, which represents the minimum number of tokens required for accurate classification, from the previously trained ViT model. Then, we train a Token-Length Assigner (TLA) using the training data, where the input is an image and the label is the corresponding token length. This decoupled procedure allows the TLA to make a better decision regarding the number of tokens required for each image. Dur-

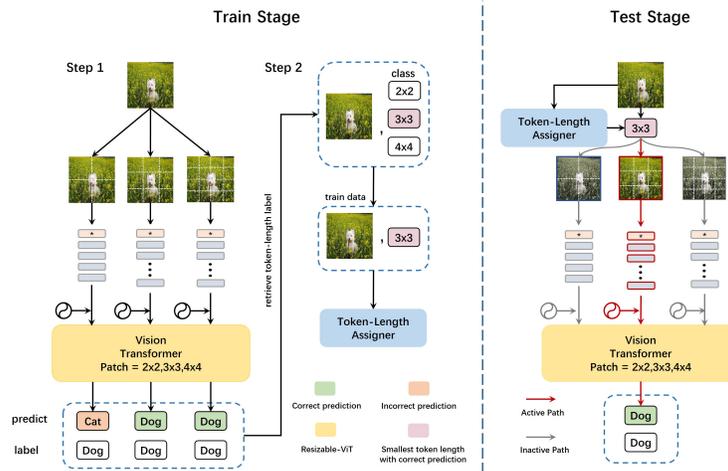


Fig. 2. Left: There are two steps in the training procedure. First, we train the Resizable-ViT that can split an image into any predefined token length. Secondly, we train a Token-Length Assigner based on the token-length label that is retrieved from ReViT. It is the smallest number of tokens that can correctly predicate the class of the image. **Right:** In inference, the TLA first assigns a token-length for the image, then ReViT uses this setting to make predication.

ing inference, the TLA guides the ViT model on the optimal number of tokens required for accurate classification based on the input. The complete training and testing process is illustrated in Figure 2.

In the following, we first introduce the Token-Label Assigner, then present the training method on the Resizable-ViT model and improved techniques.

3.1 Token-Length Assigner

The purpose of the Token-Length Assigner (TLA) is to make accurate predictions based on the feedback from ReViT. TLA training is performed after ReViT. We first define a list of token lengths $L = [l_1, l_2, \dots, l_n]$ in descending order. For simplicity, we use a single number to represent the token length, such as $L = [14 \times 14, 10 \times 10, 7 \times 7]$. The model with a token length of 7×7 has the lowest computational cost among the three token lengths.

In order to train a token-length adapter (TLA), it is necessary to obtain a token-length label from the ReViT model at convergence. For an image, the token-length label is defined as the minimum token length required by the ViT model to accurately classify that image. The inference speed of the ReViT model, denoted by M , can be ranked as $Speed(M_{l_1}) < Speed(M_{l_2}) < \dots < Speed(M_{l_k})$, where $k = len(L)$ represents the total number of options for token length. For each input x , we can obtain the prediction $y_{l_i} = M_{l_i}(X)$ for all $i \in n$. The label of the input x is determined by the smallest token size l_j for which any smaller token length would result in an incorrect prediction, i.e., $y_{l_{j-1}} \neq y^{gt}$,

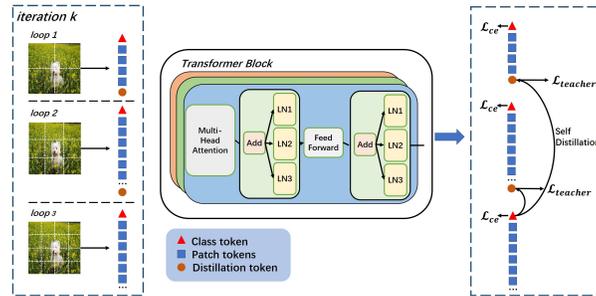


Fig. 3. Example of self-distillation and token-length aware layer normalization in ReViT. Each token length corresponds to a LayerNorm (LN in this figure) and pass-through this LayerNorm during both training and inference. The self-distillation is only conducted in training, where smaller token lengths have an extra distillation token to learn from the teacher’s knowledge.

where gt is the ground truth label. Therefore, a set of input-output pairs (x, l_j) can be obtained and used to train the TLA. Since token-label assignment is straightforward, the TLA is a lightweight module, with minimal computational overhead introduced. Moreover, since unnecessary tokens are reduced in the ViT model, the additional computational overhead is relatively small.

3.2 Resizable-ViT

In this section, we present the Resizable-ViT (ReViT), a dynamic ViT model capable of accurately classifying images with various token lengths. We introduce two techniques that enhance the performance of ReViT and subsequently present the training strategy. Additionally, we offer an efficient training implementation that accelerates the training process of ReViT.

Token-Aware Layer Normalization. The Layer Normalization (LN/LayerNorm) layer is a widely used normalization technique that accelerates training and improves the generalization of the Transformer architecture. In both natural language processing and computer vision, it is common to adopt an LN layer after addition in the transformer block. However, as the feature maps of the self-attention matrices and feed-forward networks constantly change, the number of token sizes changes as well. Consequently, inaccurate normalization statistics across different token lengths are shared in the same layer, which impairs test accuracy. Additionally, we found empirically that LN cannot be shared in ReViT.

To address this issue, we propose a Token-Length-Aware LayerNorm (TAL-LN), which uses an independent LayerNorm for each choice of token length in the predefined token length list. In other words, we use $Add \& \{LN_1, \dots, LN_k\}$ as a building block, where k represents the number of predefined token lengths. Each LayerNorm layer calculates layer-wise statistics specifically and learns the parameters of the corresponding feature map. Furthermore, the number of extra

Algorithm 1: Training Resizable-ViT M .

Require: Define Token-Length Assigner T , token-length list \mathbf{R} , for example, $\{16, 24, 32\}$. The iterations N_M for training M . The $CE(\cdot)$ denotes cross-entropy loss, and $DisT(\cdot)$ denotes distillation loss.

```

for  $t = 1, \dots, N_M$  do
    Get data  $x$  and class label  $y_c$  of current mini-batch.
    Clear gradients for all parameters,  $optimizer.zero\_grad()$ 
    for  $i = 1, \dots, len(\mathbf{R}) - 1$  do
        Convert ReViT to selected token-length  $M_i$ ,
        Execute current scaling configuration.  $\hat{y}_i = M_i(x)$ .
        if  $\mathbf{R}[i] == 16$  then
            set teacher label.  $\hat{y}_i^{teacher} = \hat{y}_i$ 
            Compute loss  $loss_i = CE(\hat{y}_i, y)$ 
        end
        else
            Compute loss  $loss_i = DisT(\hat{y}_i^{teacher}, \hat{y}_i, y)$ 
        end
        Compute gradients,  $loss_i.backward()$ 
    end
    Update weights,  $optimizer.step()$ .
end
Obtain token-length label for all train data  $(x, y_t)$ .
Train  $T$  with  $(x, y_t)$ .

```

parameters in TAL-LN is negligible since the number of parameters in normalization layers typically takes less than one percent of the total model size [54]. A brief summary is illustrated in Figure 3.

Self-Distillation It is aware that the performance of ViT is strongly correlated to the number of patches, and experiments have shown that reducing the token size significantly hampers the accuracy of small token ViT. Directly optimizing via supervision from the ground truth poses a challenge for the small token length sub-model. Motivated by self-attention, a variant of knowledge distillation techniques, where the teacher can be insufficiently trained, or even the student model itself [54,52,53], we propose a token length-aware self-distillation (TLSD). In the next section, we will show that the model with the largest token length M_1 is always trained first. For M_{l_1} , the training objective is to minimize the cross-entropy loss \mathcal{L}_{CE} . When it comes to the model with other token lengths $M_{l_i}, i \leq k, i \neq 1$, we use a distillation objective to train the target model:

$$\mathcal{L}_{teacher} = (1 - \lambda)\mathcal{L}_{CE}(\phi(Z_s), y) + \lambda\tau^2 KL(\phi(Z_s/\tau), \phi(Z_t/\tau)) \quad (1)$$

where Z_s and Z_t is the logits of the student model M_{l_i} and teacher model M_{l_1} , respectively. τ is the temperature for the distillation, λ is the coefficient balancing the KL loss (Kullack-Leibler divergence) and the CE loss (cross-entropy) on ground truth label y , and ϕ is the softmax function. Similar to DeiT, we add

a distillation token for student models. Figure 3 gives an overview. Notably, this distillation scheme is computational-free: we can directly use the predicted label of the model with the largest token length as the training label for other sub-model, while for the largest token length model, we use ground truth.

3.3 Training Strategy

To enable the ViT model to adaptively process various token lengths in the predefined choice list, it is necessary to expose it to images with different token lengths. Inspired by batch gradient accumulation, a technique used to overcome the problem of small batch size by accumulating gradient and batch statistics in a single iteration, we propose a mixing token length training. As shown in Algorithm 1, a batch of images is processed with different token lengths to compute the loss through feed-forward, and individual gradients are obtained. After looping through all token length choices, the gradients of all parameters calculated by feeding different token lengths are accumulated to update the parameters.

Efficient Training Implementation. An issue with the aforementioned training strategy is that the training time increases linearly with the number of predefined token length

choices. To address this issue, we propose an efficient implementation strategy that trades memory cost for training time. As shown in Figure 4, we replicate the model, with each model corresponding to a specific token length. At the end of each iteration, the gradients of the different replicas are synchronized and accumulated. Notably, we always send the gradient of replicas in which the token length is small to the one with a larger token length, as they are the training bottleneck. Thus, the communication cost in the gradient synchronization step is negligible. Then, the model parameters are updated through back-propagation. After the parameter updating is complete, the main process distributes the learned parameters to the rest of the replicas. These steps are repeated until the end of training, after which all replicas except the model in the main process can be removed. As such, the training time of the Resizable Transformer reduces from $O(k)$ to $O(1)$, where k is the number of predefined token lengths. Though the number of k is small, i.e., $k = 3$, in practice, the

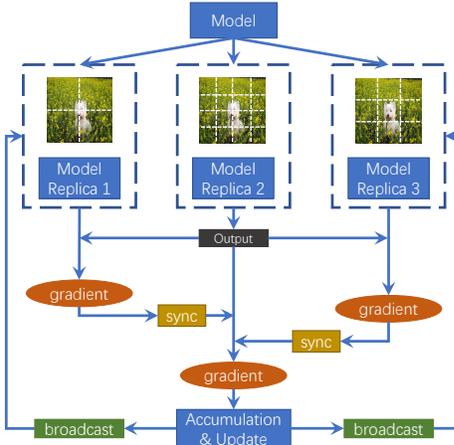


Fig. 4. Efficient training implement for Resizable Transformer through parallel computing. All gradient from the replicate nodes are synchronize on the node that have the largest token length to save the cost of communication.

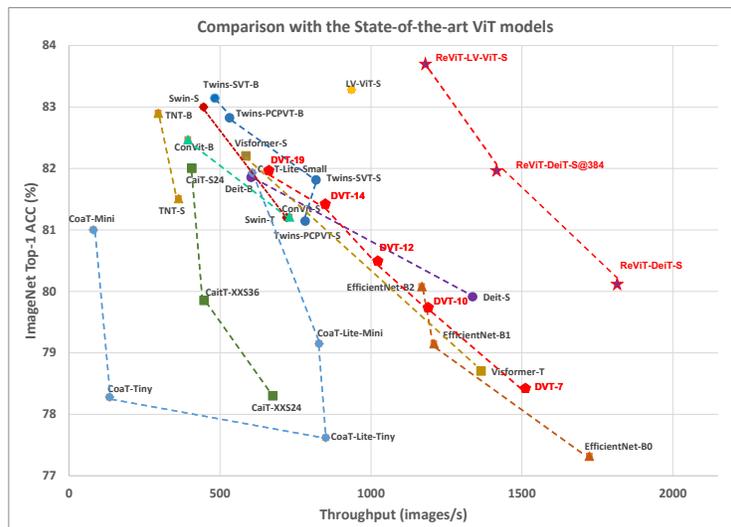


Fig. 5. Comparison of different models with various accuracy-throughput trade-off. The throughput is measured on an NVIDIA RTX 3090 GPU with batch size fixed to 32. The input image size is 224×224 unless indicate otherwise. The ReViT (red stars in the figure) achieves better trade-off than other methods, including DVT [46].

computational cost of training k ViT is high. Through our designed parallel computing, the training cost for ReViT is almost the same as that of naive ViT, where the cost of communication between replicas is negligible compared to the model training cost. In exchange for fast training, extra computational power is required for parallel computing.

4 Experiments

Implementation details. For image classification, we trained all models on the ImageNet [11] training set consisting of around 1.2 million images and reported their accuracy on the 50k test images. The predefined token lengths were set to 14×14 , 10×10 , and 7×7 by default, with the token length of 4×4 excluded due to a significant accuracy drop. We conducted experiments on DeiT-S [42] and LV-ViT-S [22] using an image resolution of 224×224 , unless otherwise specified. We followed the training settings and optimization methods described in the original papers of DeiT [42] and LV-ViT [22]. For LV-ViT, we obtained token labels for smaller token lengths using their proposed method. We also trained the ReViT on resized images with higher resolutions, such as 384 on DeiT-S. To avoid optimization difficulties caused by large kernel and stride convolutional layers required for patch embedding, we replaced them with consecutive convolutions followed by the method in Xiao et al. [48]. After training the ReViT, we obtained token-length labels for all training data and trained

Table 1. Video recognition on Something-Something V2. Our ReViT outperforms state-of-the-art CNN-based and ViT-based methods. IN-21 and K400 are abbreviations for ImageNet-21K and Kinetic-400 datasets.

Method	Backbone	FLOPs (G)	Top-1 (%)	Top-5 (%)	Frames	Extra Data
TEINet [27]	ResNet50	$99 \times 10 \times 3$	66.5	-		
TANet [28]	ResNet50	$99 \times 2 \times 3$	66.0	90.1	8+16	ImageNet-1K
TDN [44]	ResNet101	$198 \times 1 \times 3$	69.6	92.2		
SlowFast [14]	ResNet101	$106 \times 1 \times 3$	63.1	87.6	8+32	
MViTv1 [13]	MViTv1-B	$455 \times 1 \times 3$	67.7	90.9	64	Kinetics-400
TimeSformer [4]	ViT-B	$196 \times 1 \times 3$	59.5	-	8	
TimeSformer [4]	ViT-L	$5549 \times 1 \times 3$	62.4	-	64	ImageNet21K
ViViT [1]	ViT-L	$995 \times 4 \times 3$	65.9	89.9	32	
Video Swin [26]	Swin-B	$321 \times 1 \times 3$	69.6	92.7	32	
Motionformer [34]	ViT-B	$370 \times 1 \times 3$	66.5	90.1	16	IN-21K + K400
Motionformer [34]	ViT-L	$1185 \times 1 \times 3$	68.1	91.2	32	
ReViT _motionformer	ViT-B	$183 \times 1 \times 3$	66.6	89.9	16	
ReViT _motionformer	ViT-L	$570 \times 1 \times 3$	67.6	90.8	32	IN-21K + K400

the Token-Length Assigner (TLA), which was a small version of EfficientNet-B0 compared to the ViT model. We also included feature map transfer and attention transfer as part of self-distillation, which we found empirically useful. We use Something-Something V2 [16] to conduct experiments on action recognition. The Something-Something V2 is another large-scale video dataset, having around 169k videos for training and 20k videos for validation. We follow the training setting of MotionFormer [34]. Specifically, two versions of MotionFormer are tested. The default version operates on $16 \times 224 \times 224$ video clips, and a high spatial resolution variant operates on $32 \times 448 \times 448$ video clips.

4.1 Experimental Results

Main Results on ImageNet Classification. We present the main results of our ReViT based on DeiT-S and LV-ViT-S in Figure 5. Our approach is compared with several models, including DeiT [42], CaiT [38], LV-ViT [22], CoaT [49], Swin [25], Twins [10], Visformer [8], ConViT [47], TNT [17], and EfficientNet [40]. The results show that our method achieves a favorable accuracy-throughput trade-off. Specifically, ReViT reduces the computational cost of the baseline counterpart by decreasing the token number used for inference. By increasing the input resolution, we manage to outperform the baseline counterpart, given a similar computational cost. We also highlight the experimental results of DVT [46] in red. Our method achieves significantly better performance in terms of both accuracy and throughput. We hypothesize that despite the low FLOPs of DVT, the practical speed of DVT is high due to its multiple cascade ViT structure.

Main Results on Video Recognition. One of the core motivations behind ReViT is to address the issue of high computational costs in extremely long

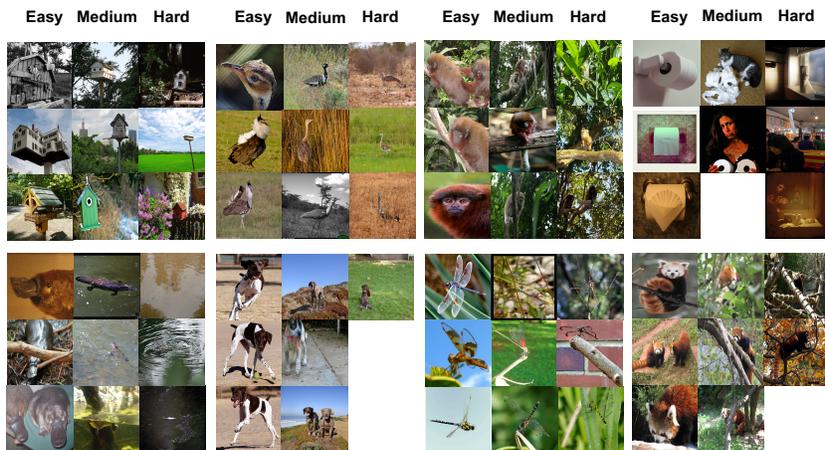


Fig. 6. Visualization of “hard”, “medium”, and “easy” samples that predicted by Token-Length Assigner and which the ReViT-DeiT-S got correction prediction. Most of the “easy” images have clear sight on the object, while size of objects is mostly small for “hard” samples.

token lengths during inference for image classification tasks. To further explore this idea, we investigate the applicability of our method to video recognition tasks, where the token length in transformers is typically much longer than that in image classifiers.

To this end, we train the ReViT-MotionFormer models with ViT-B and ViT-L, two different backbones, and compare them with the baseline models, respectively. The results are presented in Table 1. Our method demonstrates a significant speedup over the MotionFormer baseline, with a computational cost reduction of approximately 51% and a 0.1% accuracy increase. By training on larger image resolutions, we correspondingly reduce the model size by 48% with a 0.5% accuracy drop, which is slightly worse than the smaller resolution counterpart. Nonetheless, our experiments demonstrate that ReViT is effective for action recognition tasks.

Visualization of samples with different token-length. We selected eight classes from the ImageNet validation set and chose three samples from each category, classified as easy, medium, and hard, corresponding to tokens with dimensions of 14×14 , 10×10 , and 7×7 , respectively. The image samples were selected based on the token length assigned by the Token-Length Assigner. The resulting images are displayed in Figure 6. Notably, some classes do not have all images filled because less than three samples in the validation set belong to those categories. For example, only one image in the dog class requires the largest token length for classification. We observe that the number of tokens required to predict the category is highly correlated with the object’s size. For larger objects, only a few tokens are sufficient to predict their category.

Table 2. The ablation study of self-distillation in ReViT. The SD* denotes the self-distillation. We also evaluate the performance with different choices of τ . The self-distillation improves the performance notably, the small token length model outperforms the baseline when $\tau = 0.9$.

Method	SD*	τ	Top-1 Acc (%)		
			14×14	10×10	7×7
Deit-S	✗	-	79.85	74.68	72.41
ReViT	✗	-	80.12	74.24	70.15
	✓	0.5	79.92	76.16	71.33
	✓	0.9	79.83	76.86	74.21

Table 3. The ablation study of shared patch embedding and shared position encoding in ReViT. The Pos denotes the positional encoding module. We notice that sharing these two modules decrease the model accuracy.

Method	Shared		Top-1 Acc (%)		
	Patch	Pos	14×14	10×10	7×7
ReViT	✗	✓	65.14	61.30	58.35
	✓	✗	75.24	71.32	69.73
	✓	✓	79.83	76.85	74.21

4.2 Ablation Study

Shared patch embedding and position encoding. We conducted an experiment to evaluate the impact of using shared patch embedding and position encoding. As the token number changes during training, we applied some techniques to enable sharing of both operations. To handle position encoding, we followed the approach of ViT [12] and zero-padded the position encoding module whenever the token size changed. This technique was initially used to adjust the positional encoding in the pretrain-finetune paradigm. For shared patch embedding, we used a weight-sharing kernel [6]. A large kernel was constructed to process a large patch size, and when the patch size changed, a smaller kernel with shared weight on the center was adopted to flatten the image patch.

As shown in Table 3, both shared patch embedding and shared positional encoding decreased the model’s accuracy. In particular, the accuracy dropped by nearly 14% for the large token length model when using the shared patch strategy. The shared positional encoding module performed better than shared patch embedding but still significantly impacted the performance of ReViT.

The effect of self-distillation and choice of τ . We conducted experiments to verify the effectiveness of self-distillation in ReViT and investigated the impact of the hyper-parameter τ . We tested two different values of τ , 0.9 and 0.5, for all sub-networks and demonstrated the results in Table 2. Without self-distillation, the accuracy on small token lengths was comparable to tokens of size 10×10 , but

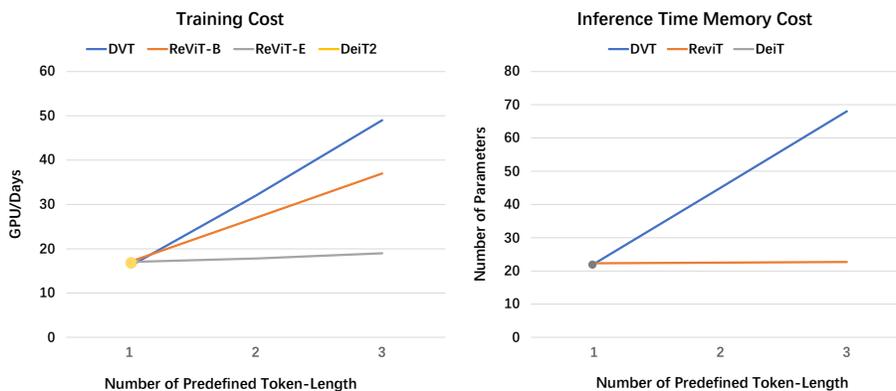


Fig. 7. Compare our approach with DeiT-S [42] and DVT [46] for training cost and memory cost at inference time in terms of the number of predefined token-length. Our proposed ReViT is almost as cheap as training the baseline DeiT-S, while DVT requires linearly increased budget on training and memory.

significantly worse on tokens of size 7×7 . When we applied self-distillation with $\tau = 5$, the accuracy of both models increased. To further evaluate the model, we used $\tau = 5$. The higher value of τ negatively impacted the accuracy of the largest token length, dropping the accuracy by around 0.3%, but significantly improving the performance of models with token size 7×7 . This highlights the necessity of using self-distillation in our scenario and demonstrates the importance of carefully selecting the hyper-parameter τ for optimal performance.

Training cost and Memory Consumption. We compared ReViT with DeiT-S and DVT [46] in terms of training cost and memory consumption, as shown in Figure 7. ReViT-B denotes the baseline approach of ReViT, while ReViT-E is the efficient implementation method. Both ReViT-B and DeiT-S show a linear increase in training cost as the number of choices in s increases. ReViT-B is cheaper because backpropagation of multiple token lengths is merged. However, the training time of ReViT-E slightly increases due to the communication cost between parallel models increasing.

As for memory consumption (number of parameters) during testing, since our method only has a single ViT where most computational heavy components are shared, the memory cost is slightly higher than the baseline. However, compared to DVT, the increase in the number of parameters with respect to the increasing number of token length choices is negligible. This indicates that our approach is more practical than DVT in terms of both training cost and memory cost. Furthermore, our method is easier to apply to existing ViT models than DVT.

Comparison with DVT. We conducted a further investigation of our proposed method based on DeiT-S and compared it with DVT, which was also developed

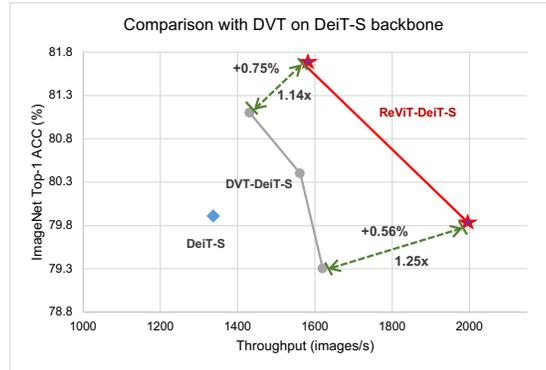


Fig. 8. Comparison with DVT [46] on DeiT-S backbone. Our method outperforms DVT by a large margin.

based on DeiT-S. Figure 8 shows that our proposed ReViT achieves superior performance compared to DVT. This could be due to our better selection of the number of patches that achieves the best accuracy-speed tradeoff.

5 Conclusions

This paper aims to reduce the token length to split the image in the ViT model to eliminate unnecessary computational costs. First, we propose the Resizable Transformer (ReViT), which adaptively processes any predefined token size for a given image. Then, we define a Token-Length Assigner to decide the minimum number of tokens that the transformer can use to classify the individual image correctly. Extensive experiments indicate that ReViT can significantly accelerate the state-of-the-art ViT model. Also, compared to the prior SOTA method, our approach achieves better training speed, inference cost, and model performance. Therefore, we believe our paper benefits practitioners who would like to adopt ViT in deployment.

References

1. Arnab, A., Deghani, M., Heigold, G., Sun, C., Lučić, M., Schmid, C.: Vivit: A video vision transformer. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6836–6846 (2021)
2. Bao, H., Dong, L., Wei, F.: Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254 (2021)
3. Beltagy, I., Peters, M.E., Cohan, A.: Longformer: The long-document transformer. arXiv preprint arXiv:2004.05150 (2020)
4. Bertasius, G., Wang, H., Torresani, L.: Is space-time attention all you need for video understanding? In: ICML. vol. 2, p. 4 (2021)
5. Bolya, D., Fu, C.Y., Dai, X., Zhang, P., Feichtenhofer, C., Hoffman, J.: Token merging: Your vit but faster. In: The Eleventh International Conference on Learning Representations (2023), <https://openreview.net/forum?id=JroZRarw7Eu>
6. Cai, H., Gan, C., Wang, T., Zhang, Z., Han, S.: Once-for-all: Train one network and specialize it for efficient deployment. arXiv preprint arXiv:1908.09791 (2019)
7. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European Conference on Computer Vision. pp. 213–229. Springer (2020)
8. Chen, Z., Xie, L., Niu, J., Liu, X., Wei, L., Tian, Q.: Visformer: The vision-friendly transformer. arXiv preprint arXiv:2104.12533 (2021)
9. Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin, A., Kaiser, L., et al.: Rethinking attention with performers. arXiv preprint arXiv:2009.14794 (2020)
10. Chu, X., Tian, Z., Wang, Y., Zhang, B., Ren, H., Wei, X., Xia, H., Shen, C.: Twins: Revisiting the design of spatial attention in vision transformers. arXiv preprint arXiv:2104.13840 1(2), 3 (2021)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Deghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
13. Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C.: Multiscale vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6824–6835 (2021)
14. Feichtenhofer, C., Fan, H., Malik, J., He, K.: Slowfast networks for video recognition. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6202–6211 (2019)
15. Frankle, J., Carbin, M.: The lottery ticket hypothesis: Finding sparse, trainable neural networks. arXiv preprint arXiv:1803.03635 (2018)
16. Goyal, R., Ebrahimi Kahou, S., Michalski, V., Materzynska, J., Westphal, S., Kim, H., Haenel, V., Freund, I., Yianilos, P., Mueller-Freitag, M., et al.: The "something something" video database for learning and evaluating visual common sense. In: Proceedings of the IEEE international conference on computer vision. pp. 5842–5850 (2017)
17. Han, K., Xiao, A., Wu, E., Guo, J., Xu, C., Wang, Y.: Transformer in transformer. arXiv preprint arXiv:2103.00112 (2021)

18. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. arXiv preprint arXiv:1510.00149 (2015)
19. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners (2021)
20. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
21. Hu, H., Gu, J., Zhang, Z., Dai, J., Wei, Y.: Relation networks for object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3588–3597 (2018)
22. Jiang, Z., Hou, Q., Yuan, L., Zhou, D., Jin, X., Wang, A., Feng, J.: Token labeling: Training a 85.5% top-1 accuracy vision transformer with 56m parameters on imagenet. arXiv preprint arXiv:2104.10858 (2021)
23. Li, Y., Fan, H., Hu, R., Feichtenhofer, C., He, K.: Scaling language-image pre-training via masking. arXiv preprint arXiv:2212.00794 (2022)
24. Liu, H., Simonyan, K., Yang, Y.: Darts: Differentiable architecture search. arXiv preprint arXiv:1806.09055 (2018)
25. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. arXiv preprint arXiv:2103.14030 (2021)
26. Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., Hu, H.: Video swin transformer. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3202–3211 (2022)
27. Liu, Z., Luo, D., Wang, Y., Wang, L., Tai, Y., Wang, C., Li, J., Huang, F., Lu, T.: Teinnet: Towards an efficient architecture for video recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 11669–11676 (2020)
28. Liu, Z., Wang, L., Wu, W., Qian, C., Lu, T.: Tam: Temporal adaptive module for video recognition. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 13708–13718 (2021)
29. Mahmood, K., Mahmood, R., Van Dijk, M.: On the robustness of vision transformers to adversarial examples. arXiv preprint arXiv:2104.02610 (2021)
30. Marin, D., Chang, J.H.R., Ranjan, A., Prabhu, A., Rastegari, M., Tuzel, O.: Token pooling in vision transformers. arXiv preprint arXiv:2110.03860 (2021)
31. Naseer, M., Ranasinghe, K., Khan, S., Hayat, M., Khan, F.S., Yang, M.H.: Intriguing properties of vision transformers. arXiv preprint arXiv:2105.10497 (2021)
32. Pan, B., Panda, R., Jiang, Y., Wang, Z., Feris, R., Oliva, A.: Ia-red²: Interpretability-aware redundancy reduction for vision transformers. arXiv preprint arXiv:2106.12620 (2021)
33. Pan, X., Xia, Z., Song, S., Li, L.E., Huang, G.: 3d object detection with pointformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7463–7472 (2021)
34. Patrick, M., Campbell, D., Asano, Y., Misra, I., Metze, F., Feichtenhofer, C., Vedaldi, A., Henriques, J.F.: Keeping your eye on the ball: Trajectory attention in video transformers. *Advances in neural information processing systems* **34**, 12493–12506 (2021)
35. Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., Dosovitskiy, A.: Do vision transformers see like convolutional neural networks? arXiv preprint arXiv:2108.08810 **4** (2021)
36. Rao, Y., Zhao, W., Liu, B., Lu, J., Zhou, J., Hsieh, C.J.: Dynamicvit: Efficient vision transformers with dynamic token sparsification. arXiv preprint arXiv:2106.02034 (2021)

37. Ryoo, M.S., Piergiovanni, A., Arnab, A., Dehghani, M., Angelova, A.: Token-learner: What can 8 learned tokens do for images and videos? arXiv preprint arXiv:2106.11297 (2021)
38. Sablayrolles, H.T.M.C.A., Jégou, G.S.H.: Going deeper with image transformers
39. Song, L., Zhang, S., Liu, S., Li, Z., He, X., Sun, H., Sun, J., Zheng, N.: Dynamic grained encoder for vision transformers. In: Thirty-Fifth Conference on Neural Information Processing Systems (2021), <https://openreview.net/forum?id=gnAIV-EKw2>
40. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: International Conference on Machine Learning. pp. 6105–6114. PMLR (2019)
41. Tang, Y., Han, K., Wang, Y., Xu, C., Guo, J., Xu, C., Tao, D.: Patch slimming for efficient vision transformers. arXiv preprint arXiv:2106.02852 (2021)
42. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. pp. 10347–10357. PMLR (2021)
43. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
44. Wang, L., Tong, Z., Ji, B., Wu, G.: Tdn: Temporal difference networks for efficient action recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1895–1904 (2021)
45. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7794–7803 (2018)
46. Wang, Y., Huang, R., Song, S., Huang, Z., Huang, G.: Not all images are worth 16x16 words: Dynamic vision transformers with adaptive sequence length. arXiv preprint arXiv:2105.15075 (2021)
47. Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., Zhang, L.: Cvt: Introducing convolutions to vision transformers. arXiv preprint arXiv:2103.15808 (2021)
48. Xiao, T., Singh, M., Mintun, E., Darrell, T., Dollár, P., Girshick, R.: Early convolutions help transformers see better. arXiv preprint arXiv:2106.14881 (2021)
49. Xu, W., Xu, Y., Chang, T., Tu, Z.: Co-scale conv-attentional image transformers. arXiv preprint arXiv:2104.06399 (2021)
50. Yang, S., Quan, Z., Nie, M., Yang, W.: Transpose: Towards explainable human pose estimation by transformer. arXiv preprint arXiv:2012.14214 (2020)
51. Yin, M., Yao, Z., Cao, Y., Li, X., Zhang, Z., Lin, S., Hu, H.: Disentangled non-local neural networks. In: European Conference on Computer Vision. pp. 191–207. Springer (2020)
52. Yu, J., Huang, T.S.: Universally slimmable networks and improved training techniques. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1803–1811 (2019)
53. Yu, J., Jin, P., Liu, H., Bender, G., Kindermans, P.J., Tan, M., Huang, T., Song, X., Pang, R., Le, Q.: Bignas: Scaling up neural architecture search with big single-stage models. In: European Conference on Computer Vision. pp. 702–717. Springer (2020)
54. Yu, J., Yang, L., Xu, N., Yang, J., Huang, T.: Slimmable neural networks. arXiv preprint arXiv:1812.08928 (2018)
55. Zhao, B., Cui, Q., Song, R., Qiu, Y., Liang, J.: Decoupled knowledge distillation. In: Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition. pp. 11953–11962 (2022)

56. Zheng, C., Zhu, S., Mendieta, M., Yang, T., Chen, C., Ding, Z.: 3d human pose estimation with spatial and temporal transformers. arXiv preprint arXiv:2103.10455 (2021)
57. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:2010.04159 (2020)
58. Zhu, Y., Liu, N., Xu, Z., Liu, X., Meng, W., Wang, L., Ou, Z., Tang, J.: Teach less, learn more: On the undistillable classes in knowledge distillation. In: Advances in Neural Information Processing Systems
59. Zhu, Y., Wang, Y.: Student customized knowledge distillation: Bridging the gap between student and teacher. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5057–5066 (2021)
60. Zhu, Y., Zhou, Q., Liu, N., Xu, Z., Ou, Z., Mou, X., Tang, J.: Scalekd: Distilling scale-aware knowledge in small object detector. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 19723–19733 (June 2023)
61. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. arXiv preprint arXiv:1611.01578 (2016)