

## AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**ABSTRACT** attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. This reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained, Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.

**Introduction** With the models and datasets growing, there is still no sign of saturating performance.

The latter models, while theoretically efficient, have not yet been scaled effectively on modern hardware accelerators due to the use of specialized attention patterns. Therefore, in large-scale image recognition, classic ResNet-like architectures are still state of the art.

We split an image into patches and provide the sequence of linear embeddings of these patches as an input to a Transformer. Image patches are treated the same way as tokens (words) in an NLP application.

Transformers lack some of the inductive biases inherent to CNNs:  
translation equivariance and locality  
do not generalize well when trained on insufficient amounts of data

**RELATED WORK** local multi-head dot-product self attention blocks can completely replace convolutions.

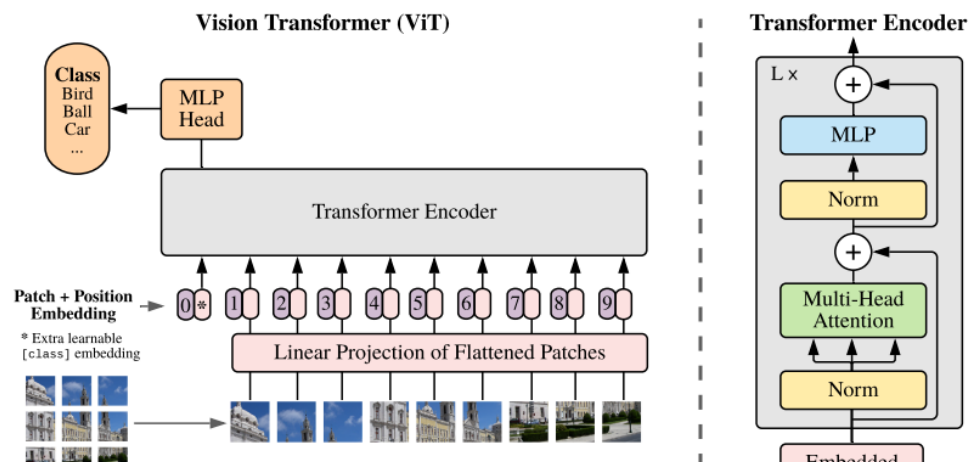
Many of these specialized attention architectures demonstrate promising results on computer vision tasks, but require complex engineering to be implemented efficiently on hardware accelerators.

## METHOD

original Transformer

An advantage of this intentionally simple setup is that scalable NLP Transformer architectures – and their efficient implementations – can be used almost out of the box.

## VISION TRANSFORMER (ViT)





↓  
↓

Embedded  
Patches

we prepend a learnable embedding to the sequence of embedded patches

Position embeddings are added to the patch embeddings to retain positional information

we have not observed significant performance gains from using more advanced 2D-aware position embeddings

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D} \quad (1)$$

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

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

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \quad (4)$$

#### Inductive bias.

much less image-specific inductive bias than CNNs

Other than that, the position embeddings at initialization time carry no information about the 2D positions of the patches and all spatial relations between the patches have to be learned from scratch.

#### Hybrid Architecture.

patch embedding projection E (Eq. 1) is applied to patches extracted from a CNN feature map

patches can have spatial size 1x1, which means that the input sequence is obtained by simply flattening the spatial dimensions of the feature map and projecting to the Transformer dimension.

The classification input embedding and position embeddings are added as described above

### FINE-TUNING AND HIGHER RESOLUTION

we remove the pre-trained prediction head and attach a zero-initialized  $D \times K$  feedforward layer

The Vision Transformer can handle arbitrary sequence lengths (up to memory constraints), however, the pre-trained position embeddings may no longer be meaningful

2D interpolation of the pre-trained position embeddings, according to their location in the original image

this resolution adjustment and patch extraction are the only points at which an inductive bias about the 2D structure of the images is manually injected into the Vision Transformer

## EXPERIMENTS

### SETUP

#### Datasets

use the ILSVRC-2012 ImageNet dataset with 1k classes and 1.3M images (we refer to it as ImageNet in what follows), its superset ImageNet-21k with 21k classes and 14M images (Deng et al., 2009), and JFT (Sun

#### Model Variants.

Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants

Table 1: Details of Vision Transformer model variants.

Transformer’s sequence length is inversely proportional to the square of the patch size  
smaller patch size are computationally more expensive.

### Training & Fine-tuning

#### Metrics

few-shot or fine-tuning accuracy

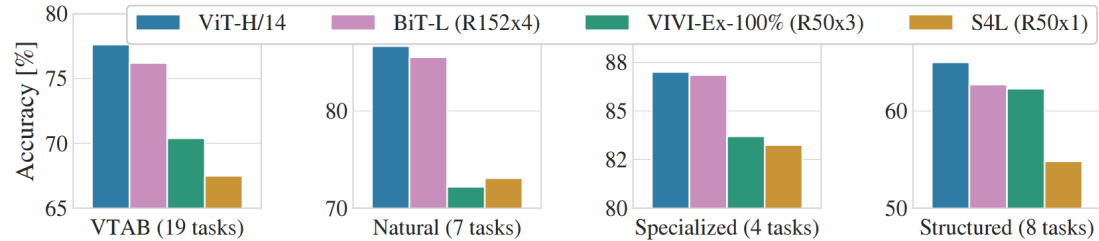
Few-shot accuracies are obtained by solving a regularized least-squares regression problem that maps the (frozen) representation of a subset of training images to  $\{-1, 1\}^K$  target vectors  
we mainly focus on fine-tuning performance, we sometimes use linear few-shot accuracies for fast on-the-fly evaluation where fine-tuning would be too costly

### COMPARISON TO STATE OF THE ART

Noisy Student

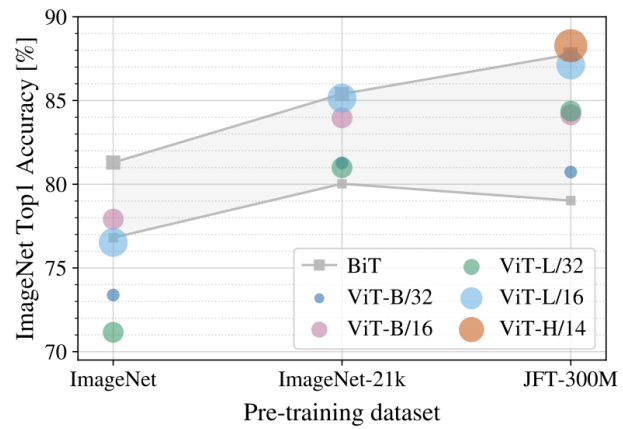
	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> $\pm 0.04$	87.76 $\pm 0.03$	85.30 $\pm 0.02$	87.54 $\pm 0.02$	88.4/88.5*
ImageNet ReaL	<b>90.72</b> $\pm 0.05$	90.54 $\pm 0.03$	88.62 $\pm 0.05$	90.54	90.55
CIFAR-10	<b>99.50</b> $\pm 0.06$	99.42 $\pm 0.03$	99.15 $\pm 0.03$	99.37 $\pm 0.06$	—
CIFAR-100	<b>94.55</b> $\pm 0.04$	93.90 $\pm 0.05$	93.25 $\pm 0.05$	93.51 $\pm 0.08$	—
Oxford-IIIT Pets	<b>97.56</b> $\pm 0.03$	97.32 $\pm 0.11$	94.67 $\pm 0.15$	96.62 $\pm 0.23$	—
Oxford Flowers-102	99.68 $\pm 0.02$	<b>99.74</b> $\pm 0.00$	99.61 $\pm 0.02$	99.63 $\pm 0.03$	—
VTAB (19 tasks)	<b>77.63</b> $\pm 0.23$	76.28 $\pm 0.46$	72.72 $\pm 0.21$	76.29 $\pm 1.70$	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

pre-training efficiency may be affected not only by the architecture choice, but also other parameters, such as training schedule, optimizer, weight decay, etc



### PRE-TRAINING DATA REQUIREMENTS

pre-train ViT models on datasets of increasing size:

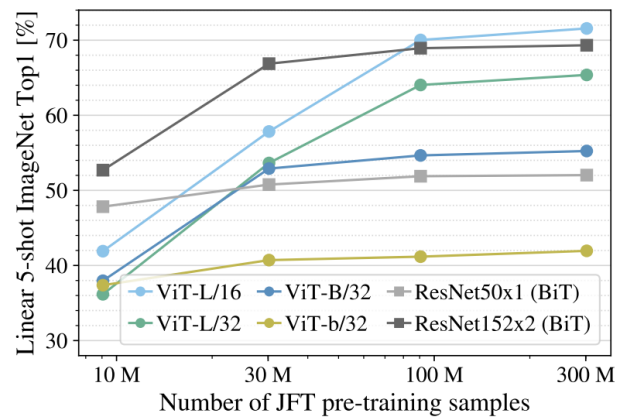


name	Epochs	ImageNet	ImageNet Real	CIFAR-10	CIFAR-100	Pets	Flowers	exaFLOPs
ViT-B/32	7	80.73	86.27	98.61	90.49	93.40	99.27	55
ViT-B/16	7	84.15	88.85	99.00	91.87	95.80	99.56	224
ViT-L/32	7	84.37	88.28	99.19	92.52	95.83	99.45	196
ViT-L/16	7	86.30	89.43	99.38	93.46	96.81	99.66	783
ViT-L/16	14	87.12	89.99	99.38	94.04	97.11	99.56	1567
ViT-H/14	14	88.08	90.36	99.50	94.71	97.11	99.71	4262
<hr/>								
ResNet50x1	7	77.54	84.56	97.67	86.07	91.11	94.26	50
ResNet50x2	7	82.12	87.94	98.29	89.20	93.43	97.02	199
ResNet101x1	7	80.67	87.07	98.48	89.17	94.08	95.95	96
ResNet152x1	7	81.88	87.96	98.82	90.22	94.17	96.94	141
ResNet152x2	7	84.97	89.69	99.06	92.05	95.37	98.62	563
ResNet152x2	14	85.56	89.89	99.24	91.92	95.75	98.75	1126
ResNet200x3	14	87.22	90.15	99.34	93.53	96.32	99.04	3306
<hr/>								
R50x1+ViT-B/32	7	84.90	89.15	99.01	92.24	95.75	99.46	106
R50x1+ViT-B/16	7	85.58	89.65	99.14	92.63	96.65	99.40	274
R50x1+ViT-L/32	7	85.68	89.04	99.24	92.93	96.97	99.43	246
R50x1+ViT-L/16	7	86.60	89.72	99.18	93.64	97.03	99.40	859
R50x1+ViT-L/16	14	87.12	89.76	99.31	93.89	97.36	99.11	1668

boost the performance on the smaller datasets, we optimize three basic regularization parameters – weight decay, dropout, and label smoothing  
JFT-300M, do we see the full benefit of larger models

BiT CNNs outperform ViT on ImageNet, but with the larger datasets, ViT overtakes.

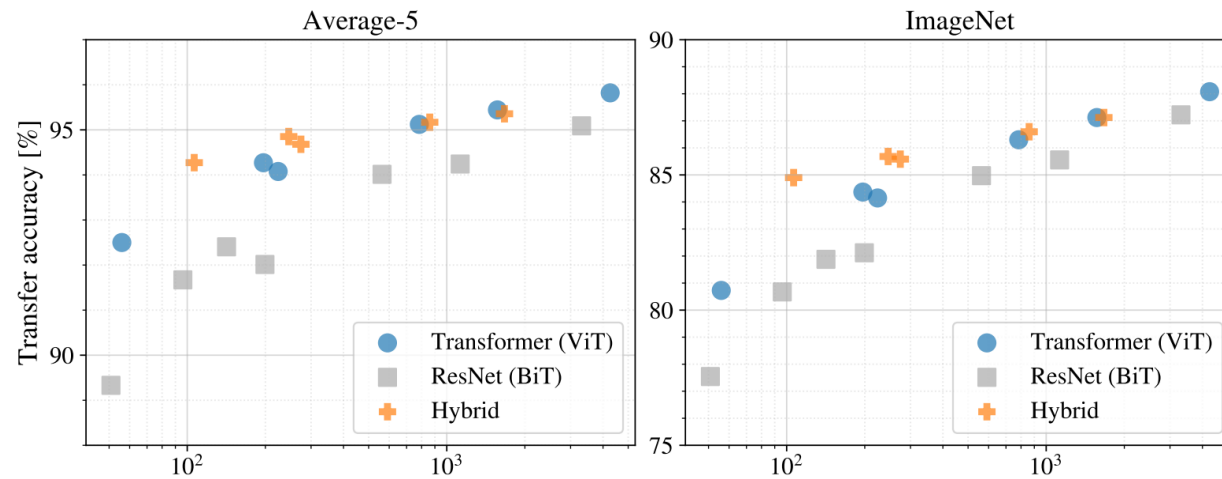
perform additional regularization on the smaller subsets and use the same hyper-parameters for all settings



Vision Transformers overfit more than ResNets with comparable computational cost on smaller datasets

This result reinforces the intuition that the convolutional inductive bias is useful for smaller datasets, but for larger ones, learning the relevant patterns directly from data is sufficient, even beneficial

## SCALING STUDY



name	Epochs	ImageNet	ImageNet Real	CIFAR-10	CIFAR-100	Pets	Flowers	exaFLOPs
------	--------	----------	---------------	----------	-----------	------	---------	----------

ViT-B/32	7	80.73	86.27	98.61	90.49	93.40	99.27	55
ViT-B/16	7	84.15	88.85	99.00	91.87	95.80	99.56	224
ViT-L/32	7	84.37	88.28	99.19	92.52	95.83	99.45	196
ViT-L/16	7	86.30	89.43	99.38	93.46	96.81	99.66	783
ViT-L/16	14	87.12	89.99	99.38	94.04	97.11	99.56	1567
ViT-H/14	14	88.08	90.36	99.50	94.71	97.11	99.71	4262
ResNet50x1	7	77.54	84.56	97.67	86.07	91.11	94.26	50
ResNet50x2	7	82.12	87.94	98.29	89.20	93.43	97.02	199
ResNet101x1	7	80.67	87.07	98.48	89.17	94.08	95.95	96
ResNet152x1	7	81.88	87.96	98.82	90.22	94.17	96.94	141
ResNet152x2	7	84.97	89.69	99.06	92.05	95.37	98.62	563
ResNet152x2	14	85.56	89.89	99.24	91.92	95.75	98.75	1126
ResNet200x3	14	87.22	90.15	99.34	93.53	96.32	99.04	3306
R50x1+ViT-B/32	7	84.90	89.15	99.01	92.24	95.75	99.46	106
R50x1+ViT-B/16	7	85.58	89.65	99.14	92.63	96.65	99.40	274
R50x1+ViT-L/32	7	85.68	89.04	99.24	92.93	96.97	99.43	246
R50x1+ViT-L/16	7	86.60	89.72	99.18	93.64	97.03	99.40	859
R50x1+ViT-L/16	14	87.12	89.76	99.31	93.89	97.36	99.11	1668

Vision Transformers dominate ResNets on the performance/compute trade-off

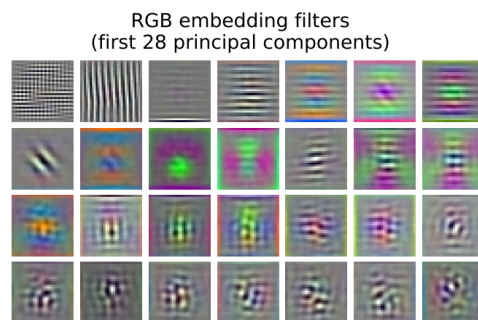
hybrids slightly outperform ViT at small computational budgets, but the difference vanishes for larger models

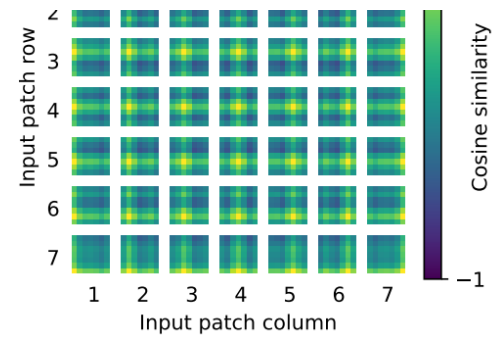
Vision Transformers appear not to saturate within the range tried, motivating future scaling efforts.

## INSPECTING VISION TRANSFORMER

internal representations

The first layer of the Vision Transformer linearly projects the flattened patches into a lower-dimensional space

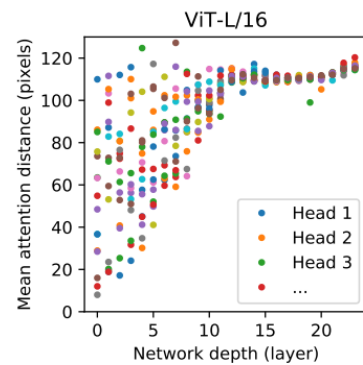




closer patches tend to have more similar position embeddings  
 patches in the same row/column have similar embeddings  
 a sinusoidal structure is sometimes apparent for larger grids

Self-attention allows ViT to integrate information across the entire image even in the lowest layers

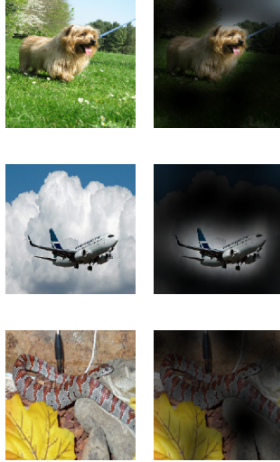
we compute the average distance in image space across which information is integrated, based on the attention weights



“attention distance” is analogous to receptive field size in CNNs  
 some heads attend to most of the image already in the lowest layers  
 This highly localized attention is less pronounced in hybrid models that apply a ResNet before the Transformer  
 it may serve a similar function as early convolutional layers in CNNs

attention distance increases with network depth  
 model attends to image regions that are semantically relevant for classification





## SELF-SUPERVISION

much of their success stems not only from their excellent scalability but also from large scale self-supervised pre-training  
 perform a preliminary exploration on masked patch prediction for self-supervision, mimicking the masked language modeling task used in BERT.

## CONCLUSION

Unlike prior works using self-attention in computer vision, we do not introduce image-specific inductive biases into the architecture apart from the initial patch extraction step  
 Instead, we interpret an image as a sequence of patches and process it by a standard Transformer encoder as used in NLP

many challenges remain.

One is to apply ViT to other computer vision tasks, such as detection and segmentation

continue exploring selfsupervised pre-training methods.

further scaling of ViT would likely lead to improved performance.