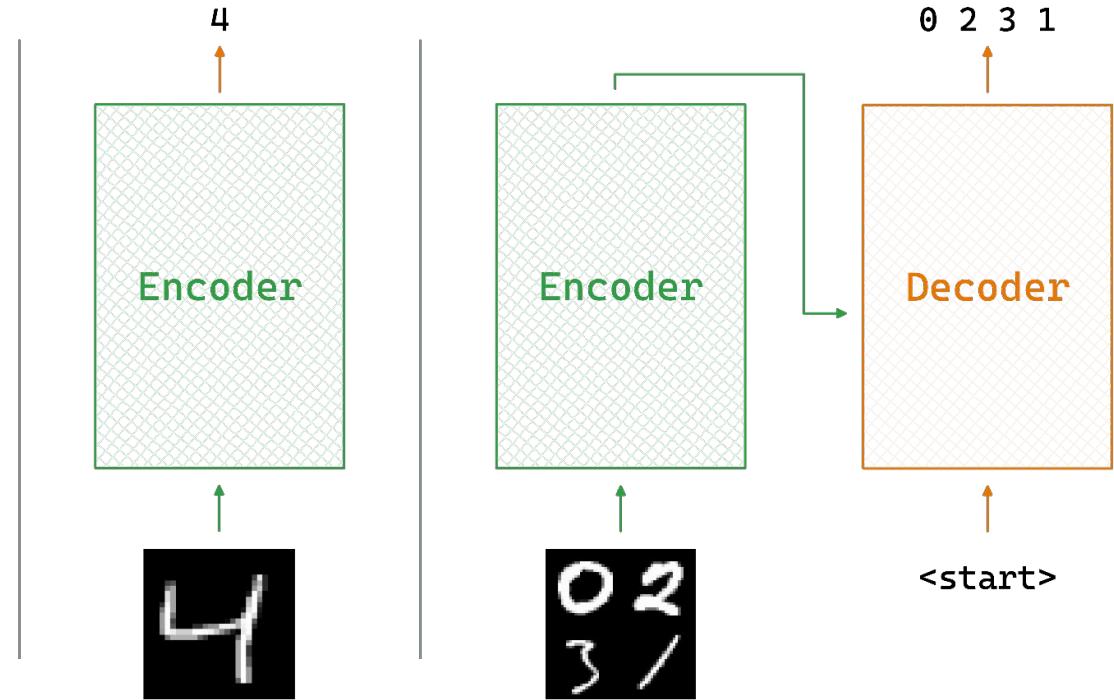


# A picture is worth a thousand words

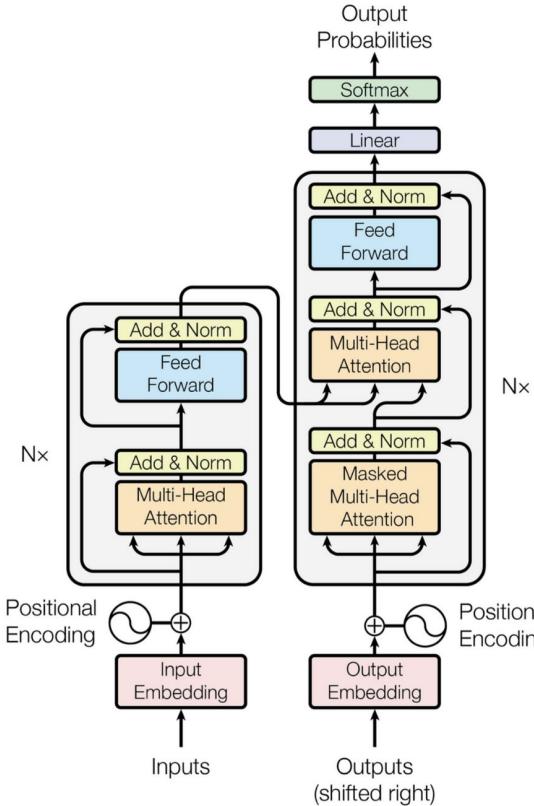
Can you write them?

# Task

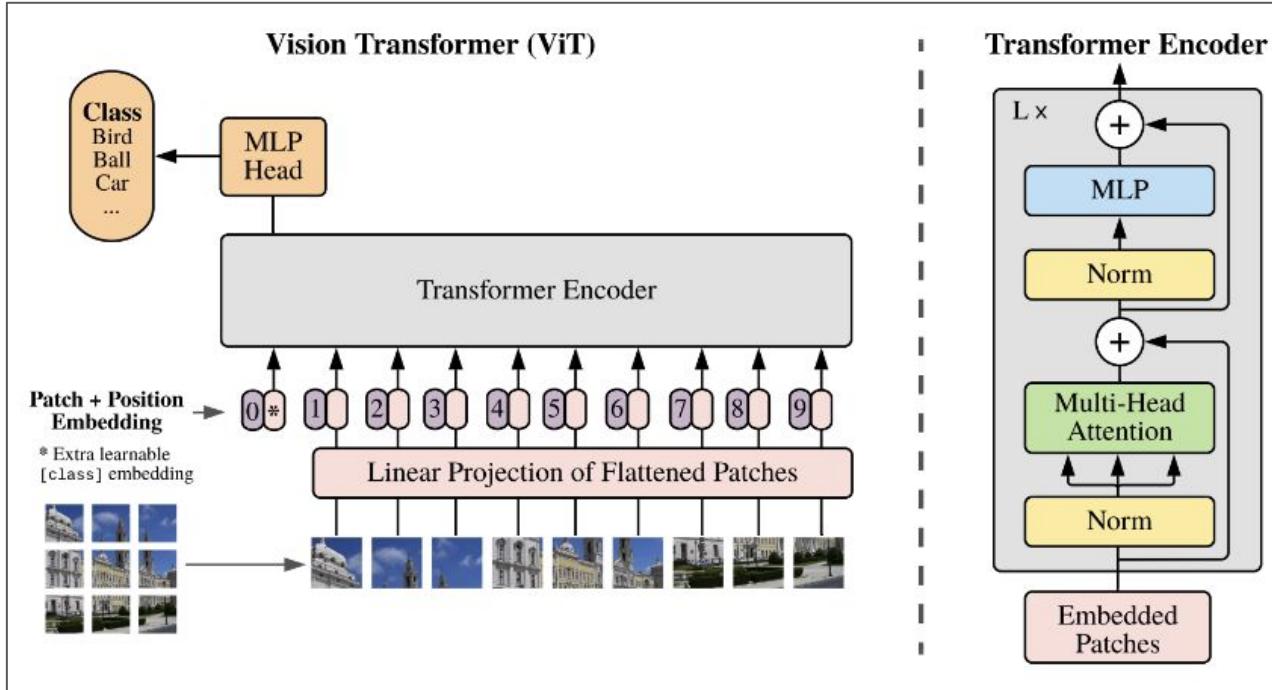
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



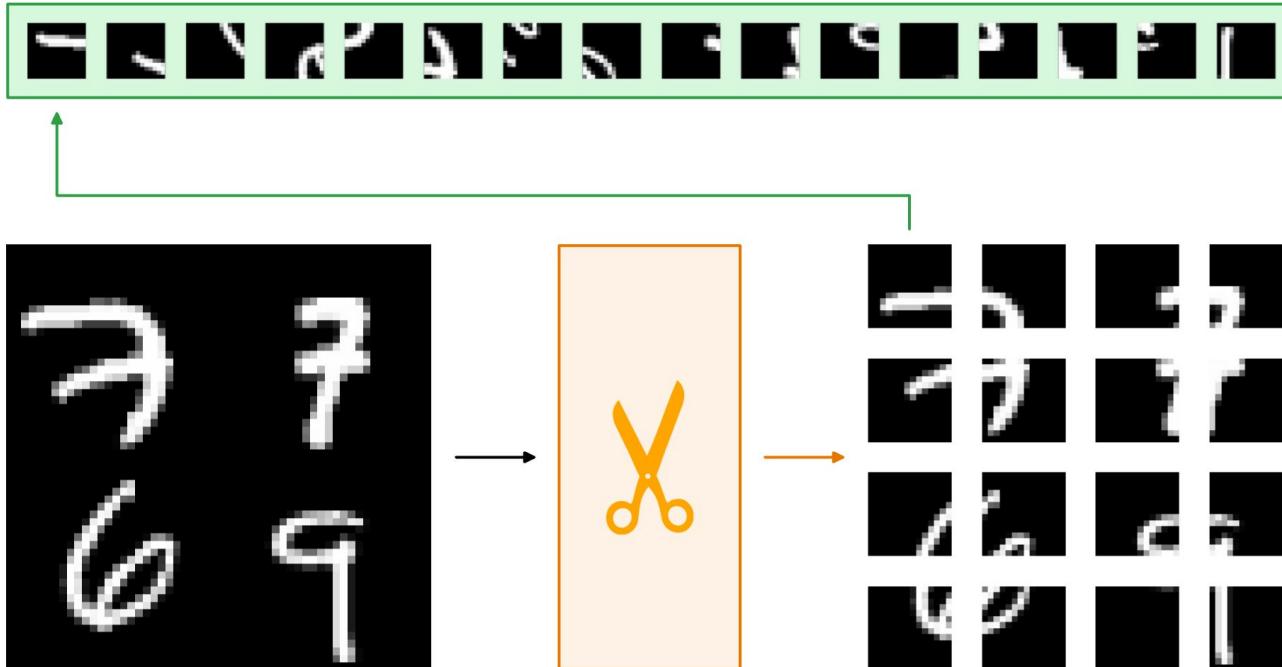
# Attention Is All You Need



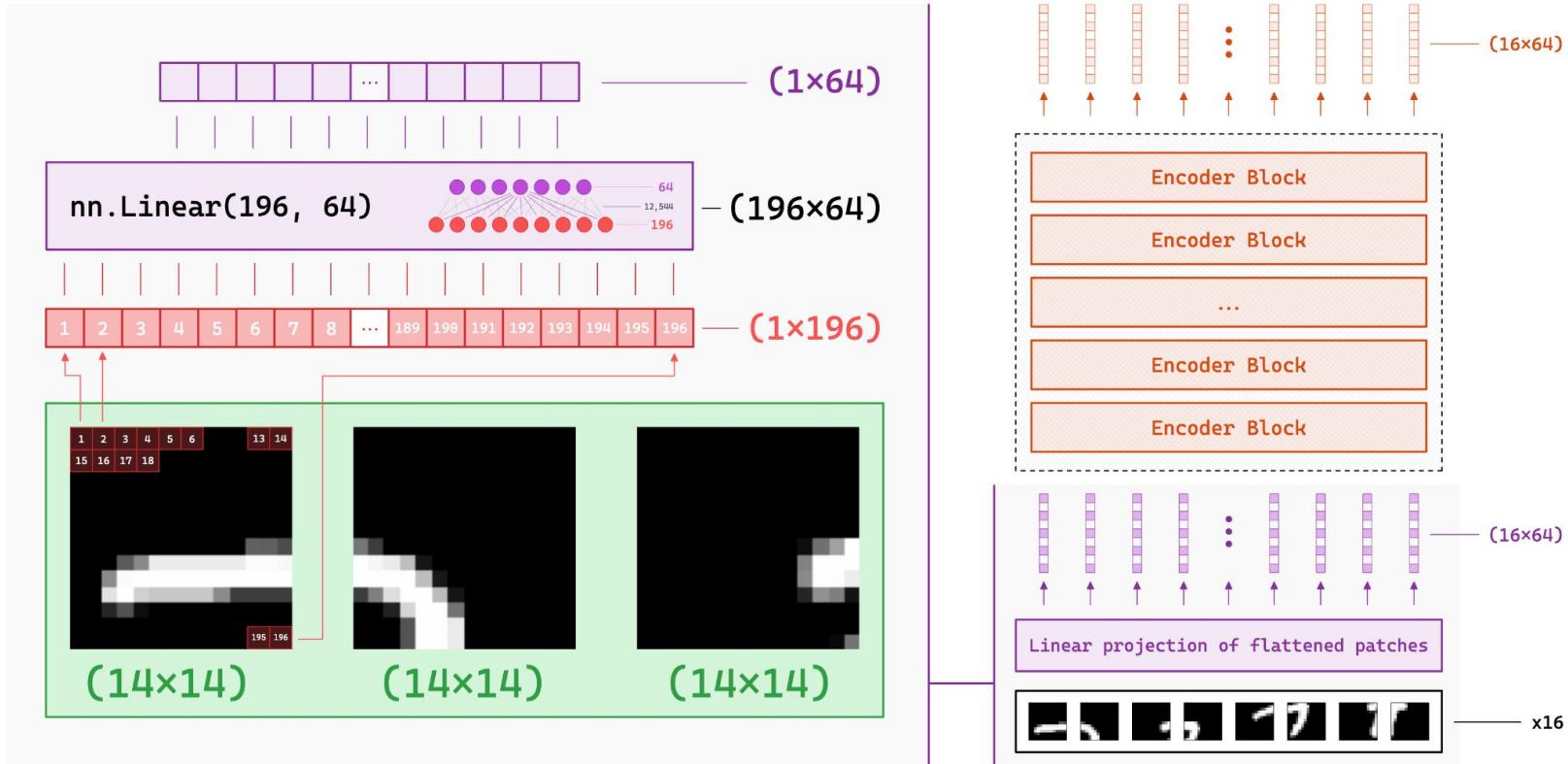
# Vision Transformer (ViT)



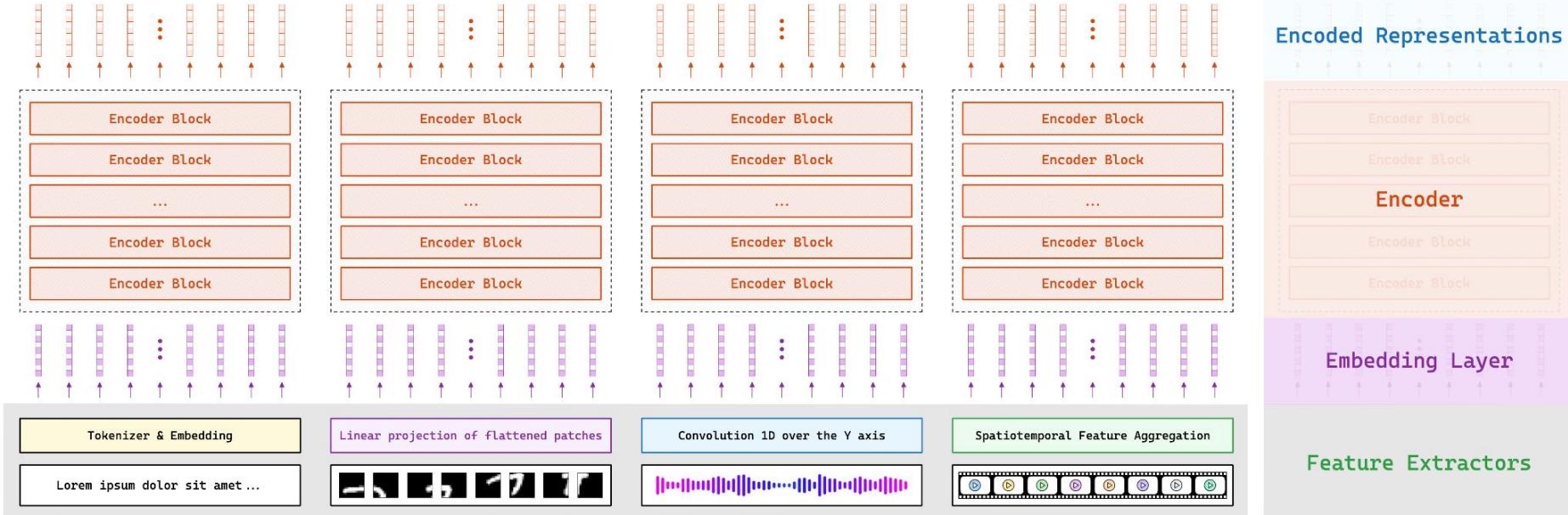
# Prep



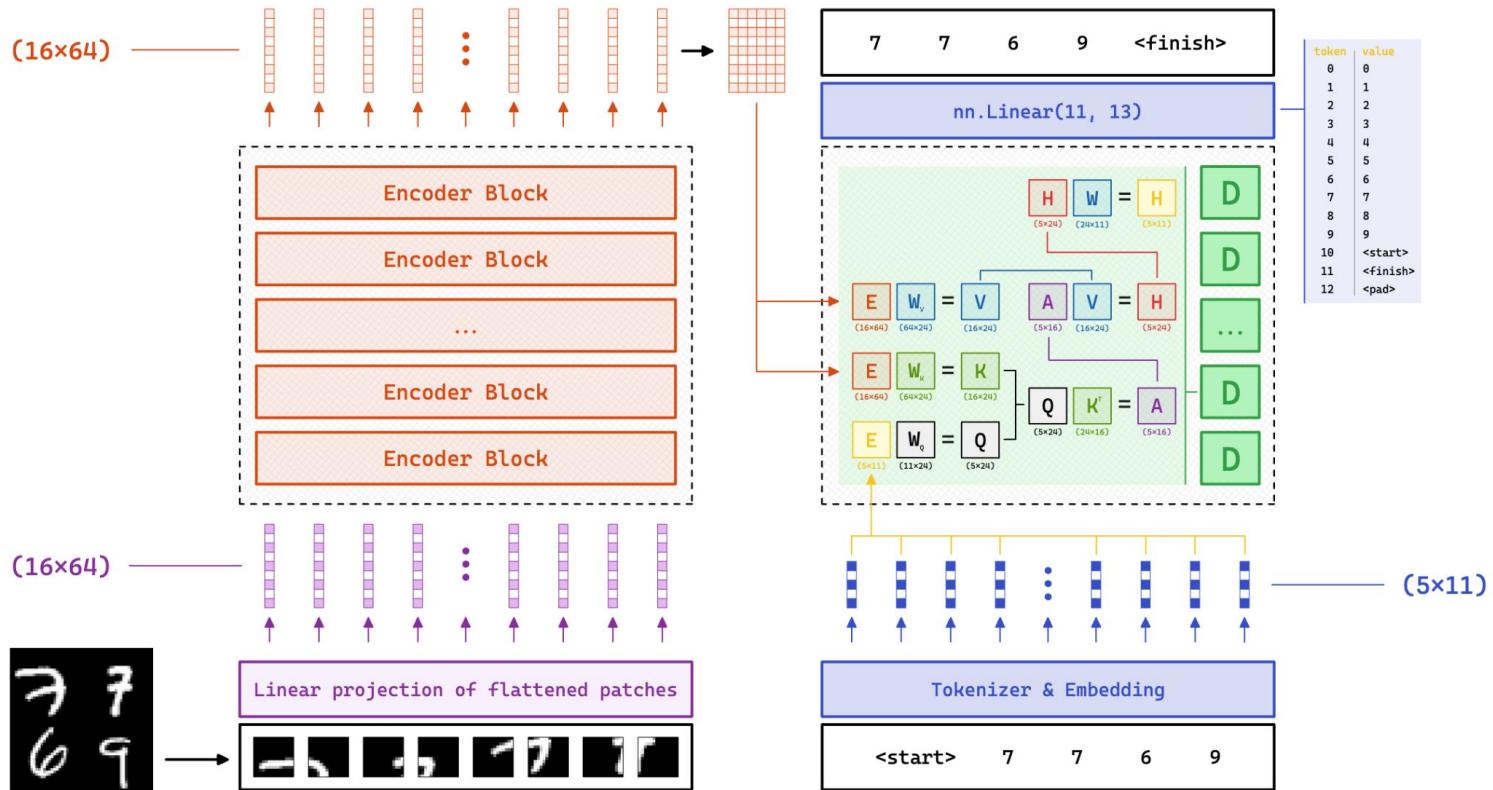
# Patch Projection



# Encoders



# Transformer



# Suggestions

- do not just copy-paste from ChatGPT
- understand and practice PyTorch ops
- use google colab for little snippets
- pair programming, swap pairs within the team
- watch tutorials but make sure to talk
- no need for GPUs yet, do everything local



## Multimodal Transformers and Image Captioning



2 Lessons | 40 hours  
with **Ardavan Afshar**





# Papers



arXiv:1706.03762v7 [cs.CL] 2 Aug 2023

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

## Attention Is All You Need

Ashish Vaswani*	Noam Shazeer*	Niki Parmar*	Jakob Uszkoreit*
Google Brain	Google Brain	Google Research	Google Research
avaswani@google.com	noam@google.com	nikip@google.com	usz@google.com

Llion Jones*	Aidan N. Gomez*	Lukasz Kaiser*
Google Research	University of Toronto	Google Brain
llion@google.com	aidan@cs.toronto.edu	lukasz.kaiser@google.com

Ilia Polosukhin <sup>†</sup>
ilia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional networks that stack multiple layers of encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior to their recurrent and seq2seq counterparts at a fraction of the cost in terms of less time to train. Our model achieves 26.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model achieves a new state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

\*Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Ilia, designed and implemented the first Transformer models and helped to evaluate them. Aidan helped to implement the first version of the codebase and helped to evaluate attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and encouraged us to publish this work. Ilia and Aidan spent countless long days designing various parts of and implementing tensor2tensor, replacing our original codebase, greatly improving results and massively accelerating our research.

<sup>†</sup>Work performed while at Google Brain.

<sup>\*</sup>Work performed while at Google Research.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

Published as a conference paper at ICLR 2021

## AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy\*, Lucas Beyer\*, Alexander Kolesnikov\*, Dirk Weissenborn\*, Xiaohua Zhai\*, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby\*<sup>†</sup>  
 \*equal contribution, <sup>†</sup>equal advising  
 Google Research, Brain Team  
 {adosovitskiy, neilhoulsby}@google.com

### ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision are limited. In vision, the common approach is to apply a convolution with a convolutional network to replace certain components of convolutional networks while keeping their overall structure intact. We show that 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. We pre-train on large datasets of image and text pairs and multiple mid-sized or small-scale recognition benchmarks (ImageNet, CIFAR-100, VTBAB, etc.). Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.

### 1 INTRODUCTION

Self-attention-based architectures, in particular Transformers (Vaswani et al., 2017), have become the dominant choice for many natural language processing (NLP). The Transformer approach is to pre-train on large text corpora and then fine-tune on a smaller specific task (Devlin et al., 2019). Thanks to Transformers’ computational efficiency and scalability, it is now possible to train models of unprecedented size, with over 100B parameters (Brown et al., 2020; Lepikhin et al., 2020). With the models and datasets growing, there is still no sign of saturating performance.

In computer vision, however, convolutional architectures remain dominant (LeCun et al., 1989; Krizhevsky et al., 2012; He et al., 2016). Inspired by NLP success, multiple works try combining Conv-like blocks with attention mechanisms (Wen et al. & Guo, 2020), sometimes adding the convolutions entirely (Ramanathan et al., 2019; Wang et al., 2020). These latter models, while theoretically efficient, have not yet 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 (Mahajan et al., 2018; Xie et al., 2020; Kolesnikov et al., 2020).

Inspired by the Transformer scaling successes in NLP, we experiment with applying a standard Transformer directly to images, with the fewest possible modifications. To do so, we split an image into patches and project 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. We train the model on image classification in supervised fashion.

When trained on mid-sized datasets such as ImageNet without strong regularization, these models yield modest accuracies of a few percentage points below ResNets of comparable size. This seemingly discouraging outcome may be expected: Transformers lack some of the inductive biases

<sup>\*</sup>Fine-tuning code and pre-trained models are available at [https://github.com/google-research/vision\\_transformer](https://github.com/google-research/vision_transformer)



# Good luck!



# Teams

## Recurrent Rebels

Dimitris  
Kenton  
Emil  
Liam

## Gradient Gigglers

Josh  
Loredana  
Filippo  
Kori

## Overfitting Overlords

Milo  
Stanley  
Artemis  
Pyry

## Perceptron Party

Daniel  
Ollie  
Maxi  
Andrea

## Backprop Bunch

Guillaume  
Amy  
Yurii  
Aygun

## Dropout Disco

Ardrit  
James  
Evelyn  
Gaurav

## Activation Aces

Nnamdi  
Neville  
Coline  
David  
Dimitar