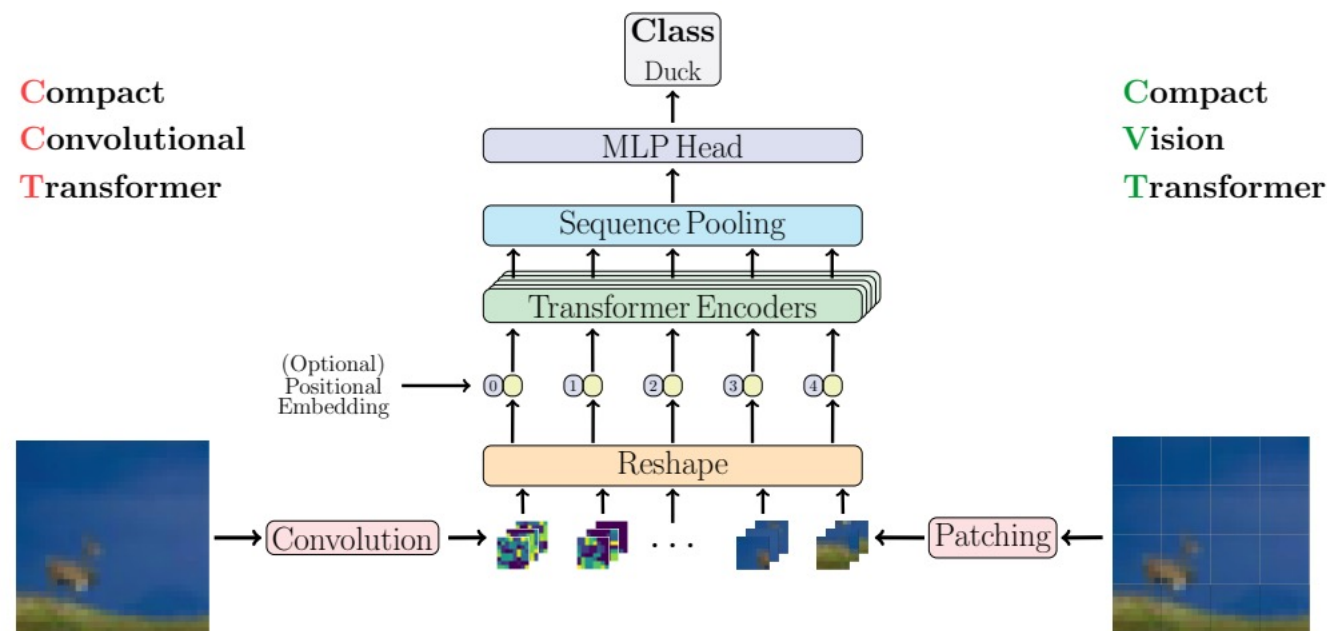


Escaping the Big Data Paradigm with Compact Transformers

- Transformers are not suitable for small sets of data.
- Compact Transformers with the right size, convolutional tokenization, transformers can avoid overfitting and outperform state-of-the-art CNNs on small datasets.

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



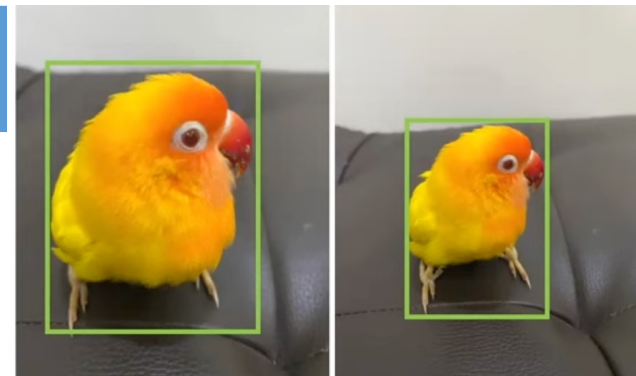
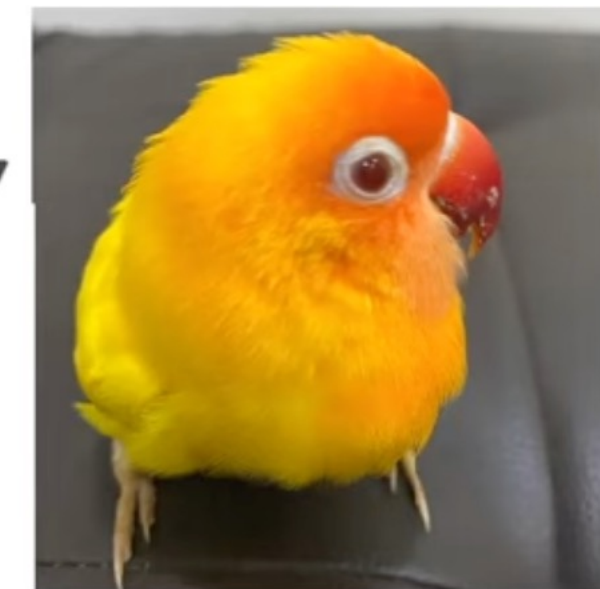
Introduction

- Convolution is translation and scale equivariant
- Pooling is translation and scale invariant
- Equivariance (Conv) + Invariance (Pool): Important for recognition, object detection and segmentation.
- ViT is Invariant to Position of Patches



Patch									
Position	1	2	3	4	5	6	7	8	9

Bird ✓
→

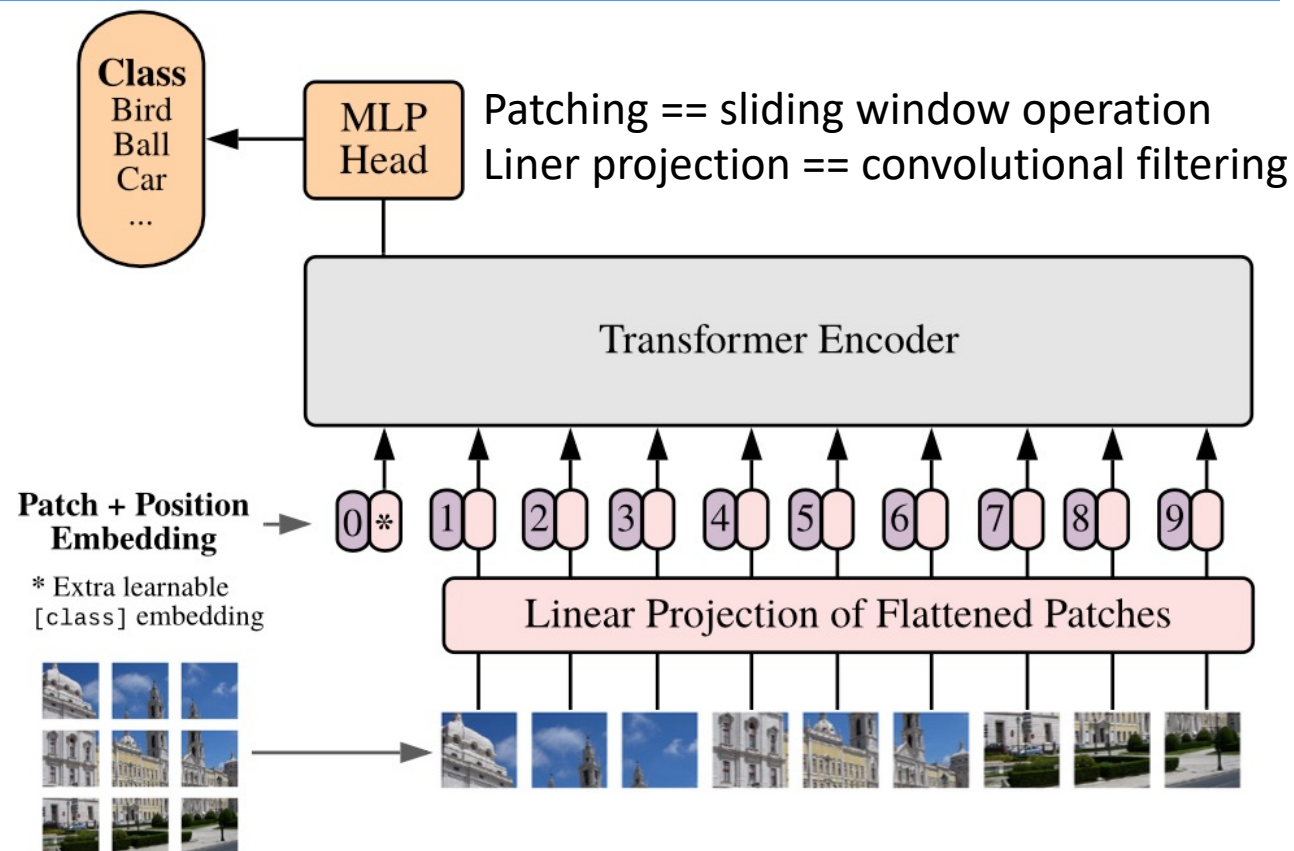


“Transformers lack some of the inductive biases inherent to CNNs, such as **translation equivariance and locality**, and therefore do not generalize well when trained on insufficient amounts of data.”

Introduction

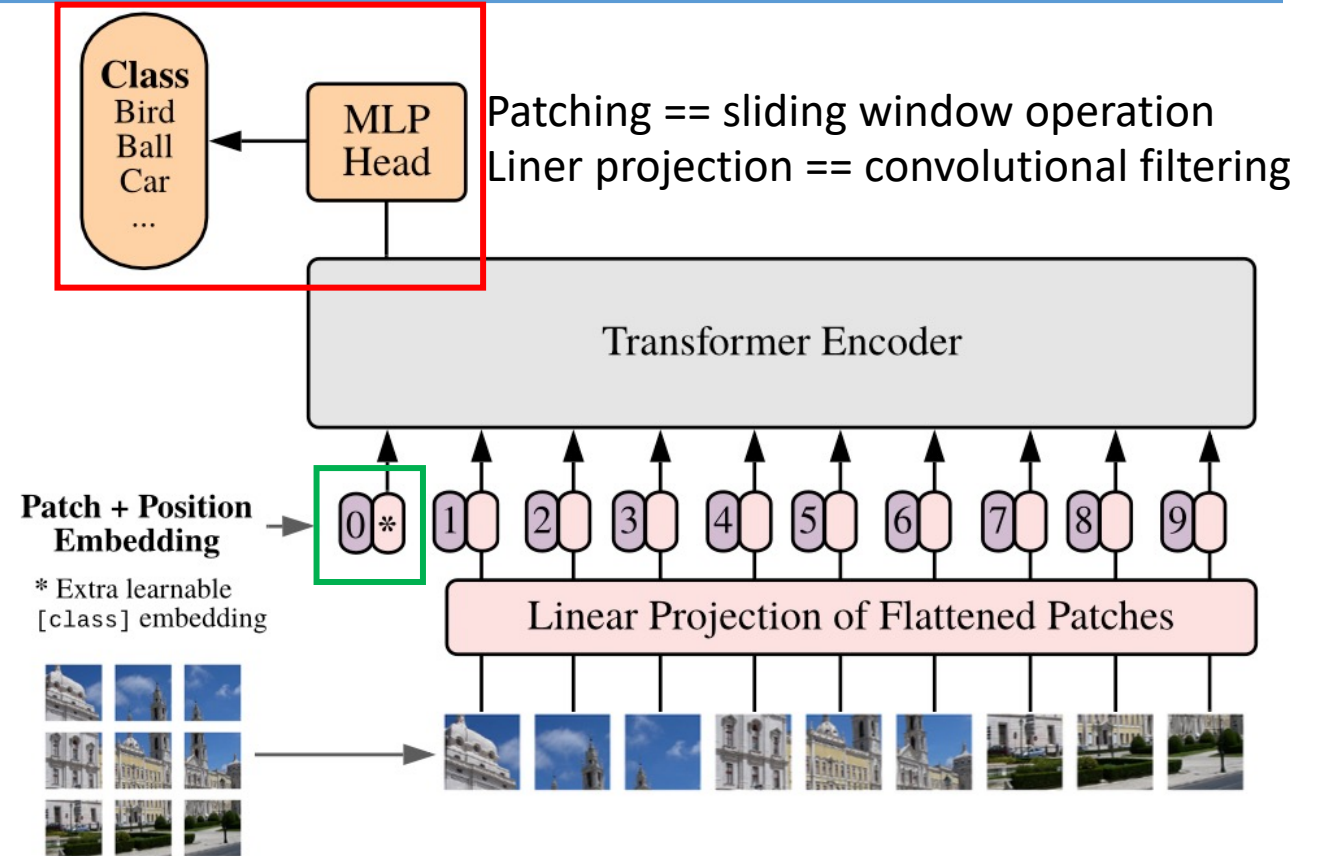
- **CNNs** are still the go-to **models for smaller datasets** because they are more efficient, both computationally and in terms of memory, when compared to transformers. Additionally, **local inductive bias** shows to be more important in smaller images.
- **CNNs** require less time and data to train while also requiring a lower number of parameters to accurately fit data. **However**, they do not enjoy the **long range interdependence** that **attention mechanisms** in **transformers** provide.

- **Image Tokenization:** ViT subdivides an image into non-overlapping square patches in raster-scan order. The sequence of patches, $x_p \in \mathbb{R}^{H \times (P^2C)}$ with patch size P , are flattened into 1D vectors and transformed into latent vectors of dimension d .



- **Positional Embedding:** Positional embedding adds spatial information into the sequence. Since the model does not actually know anything about the spatial relationship between tokens.

- A learnable [cls] token is prepended to every sequence of patches. ($D \times 1$ dimensional vector)
- After applying PE to the sequences, they are fed to the Transformer encoder, which outputs $(N + 1) \times D$ dimensional vector (N is the sequence length).
- At the classification stage only the first $1 \times D$ row (corresponding to the learnable [cls] token) is taken and fed to the final classifier.



- ViT trained from scratch on ImageNet, the results were slightly worse. A pilot experiment: the ViT base model achieved only 69.82% Top-1 on CiFAR-10 if directly trained from scratch.
- ViT has achieve a really great performance on a large dataset as ImageNet.

"Transformers lack some of the inductive biases inherent to CNNs, such as translation equivariance and locality, and therefore do not generalize well when trained on insufficient amounts of data."

~ Dosovitskiy et al. 2020

Accuracy:
0.98222232553694,
Precision:
0.98222232553694,
Recall:
0.98222232553694,
F1 score:
0.98222232553694

Accuracy:
0.6305799550480313,
Precision:
0.6352797018157111,
Recall:
0.6305799550480313,
F1 score:
0.6224562724431356

Swin Transformer training results on my dataset:
fine-tuning** vs **scratch

Compact Vision Transformer

2.4. Comparison

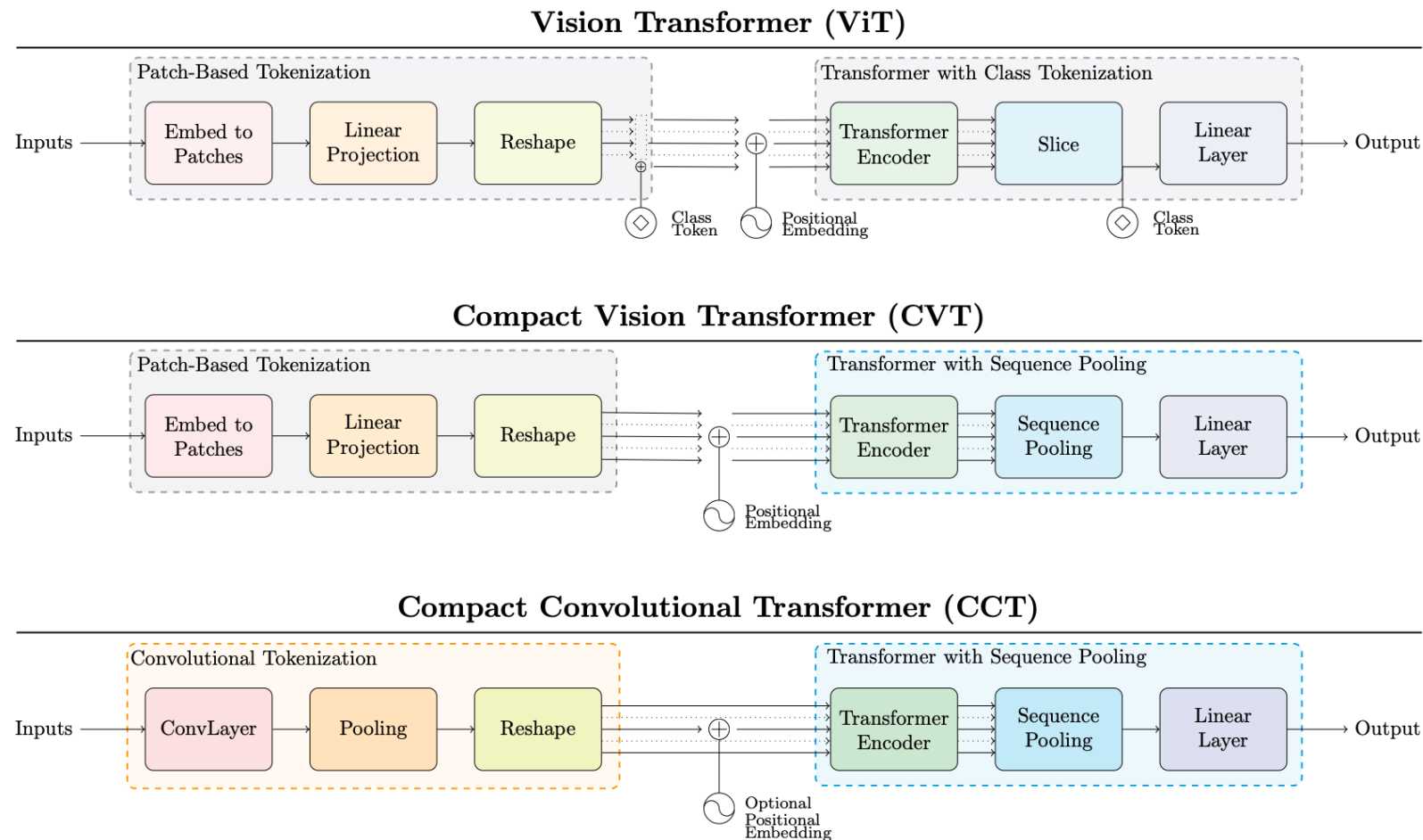
Our work differs from the aforementioned in several ways, in that it focuses on answering the following question: **Can vision transformers be trained from scratch on small datasets?** Focusing on a small datasets, we seek to create a model that can be trained, from scratch, on datasets that are orders of magnitude smaller than ImageNet. Having a model that is compact, small in size, and efficient allows greater accessibility, as training on ImageNet is still a difficult and data intensive task for many researchers. Thus our focus is on an accessible model, with few parameters, that can quickly and efficiently be trained on smaller platforms while still maintaining SOTA results.

Method

3 different models:

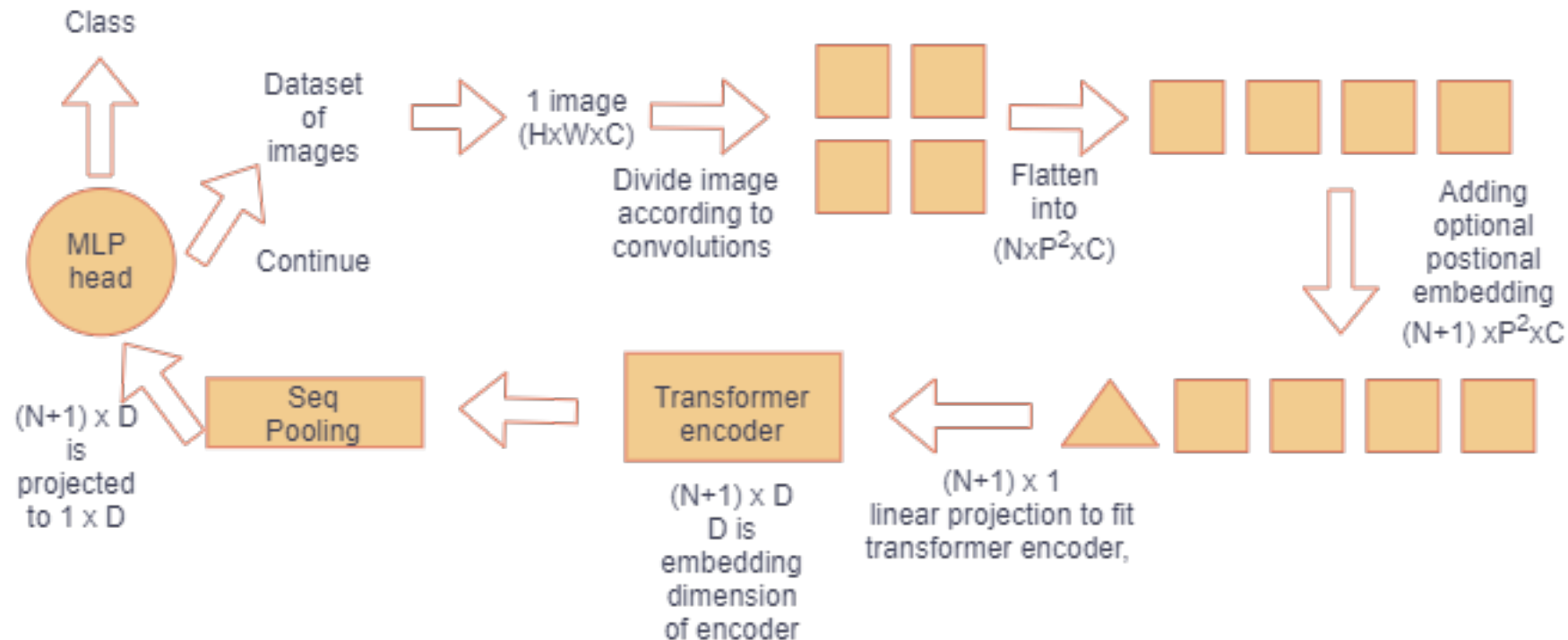
1. ViT-Lite is nearly identical to the original ViT: suitable size and patch size for small-scale learning.
2. CVT builds on ViT by using our Sequence Pooling method (SeqPool), that pools the entire sequence of tokens produced by the transformer encoder. SeqPool replaces the conventional [class] token.
3. CCT is builds on CVT and utilizes a convolutional tokenizer, generating richer tokens and preserving local information.

Method



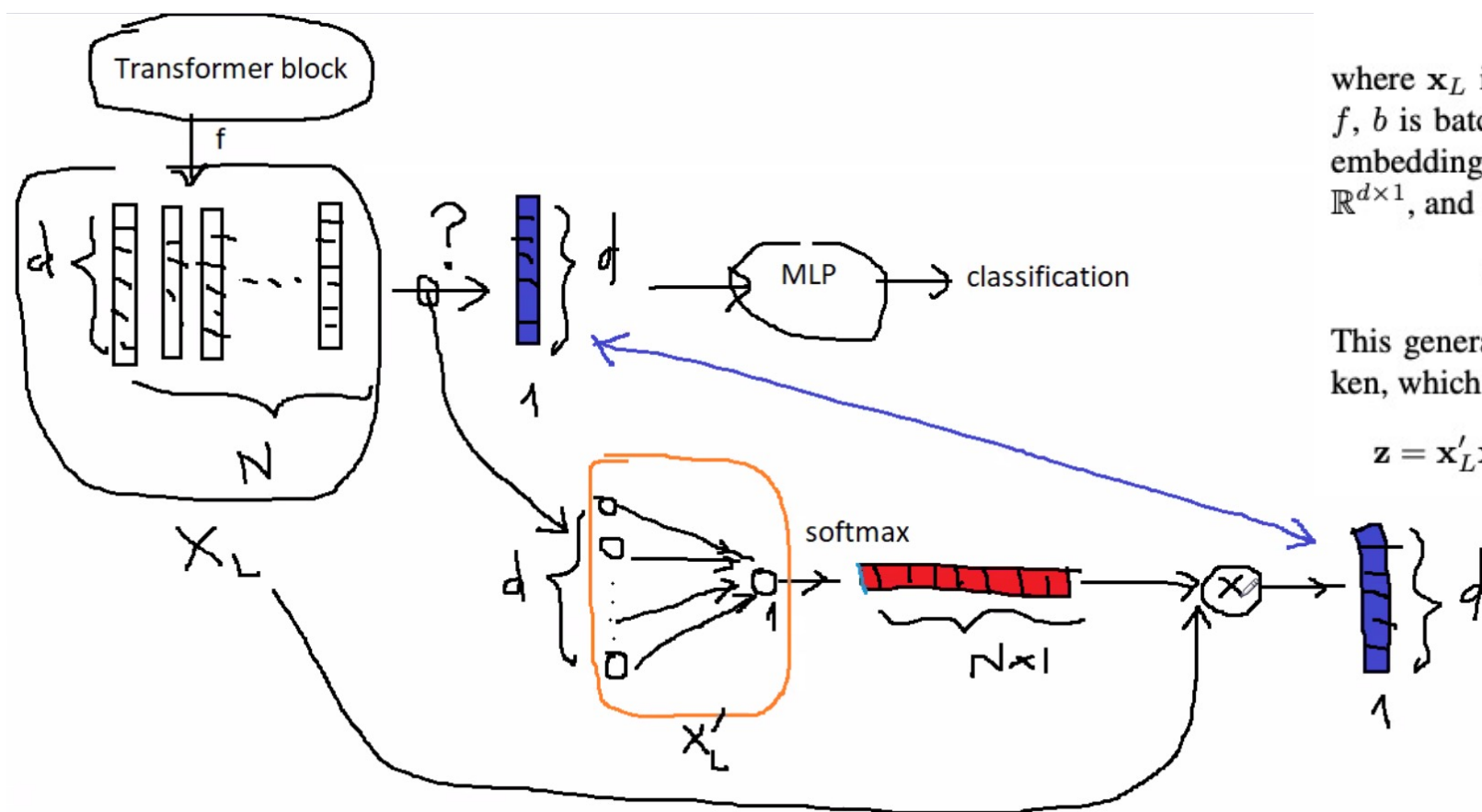
Transformer-based Backbone
Small and Compact Models
SeqPool

Method



[Shreyas-Bhat/CompactTransformers: Implementation of "Escaping the big data paradigm with Compact Transformers" by Ali Hassani et al. in Keras \(github.com\)](#)

Sequence Pooling



$$\mathbf{x}_L = f(\mathbf{x}_0) \in \mathbb{R}^{b \times n \times d}$$

where \mathbf{x}_L is the output of an L layer transformer encoder f , b is batch size, n is sequence length, and d is the total embedding dimension. \mathbf{x}_L is fed to a linear layer $g(\mathbf{x}_L) \in \mathbb{R}^{d \times 1}$, and softmax activation is applied to the output:

$$\mathbf{x}'_L = \text{softmax}(g(\mathbf{x}_L)^T) \in \mathbb{R}^{b \times 1 \times n}$$

This generates an importance weighting for each input token, which is applied as follows:

$$\mathbf{z} = \mathbf{x}'_L \mathbf{x}_L = \text{softmax}(g(\mathbf{x}_L)^T) \times \mathbf{x}_L \in \mathbb{R}^{b \times 1 \times d} \quad (1)$$

Method

- Instead of prepending a D dimensional [cls] token vector to the sequence, we place that vector as a $D \times 1$ linear layer after the Transformer encoder. It therefore maps the $N \times D$ output matrix to an $N \times 1$ vector.
- We apply softmax to the $N \times 1$ output, transpose it and multiply it by the original $N \times D$ output, which will result in a $1 \times D$ output. We pool over the entire sequence of data since it contains relevant information across different parts of the input image, making our model compact.
- The $1 \times D$ is then taken and fed to the final classifier.

Method

Original ViT-Base

12 layers

Non-overlapping Convolutions _{16x16}

Class Token

69.82% Top-1

CVT

6 layers

Non-overlapping Convolutions _{4x4}

SeqPool

92.58% Top-1

CCT

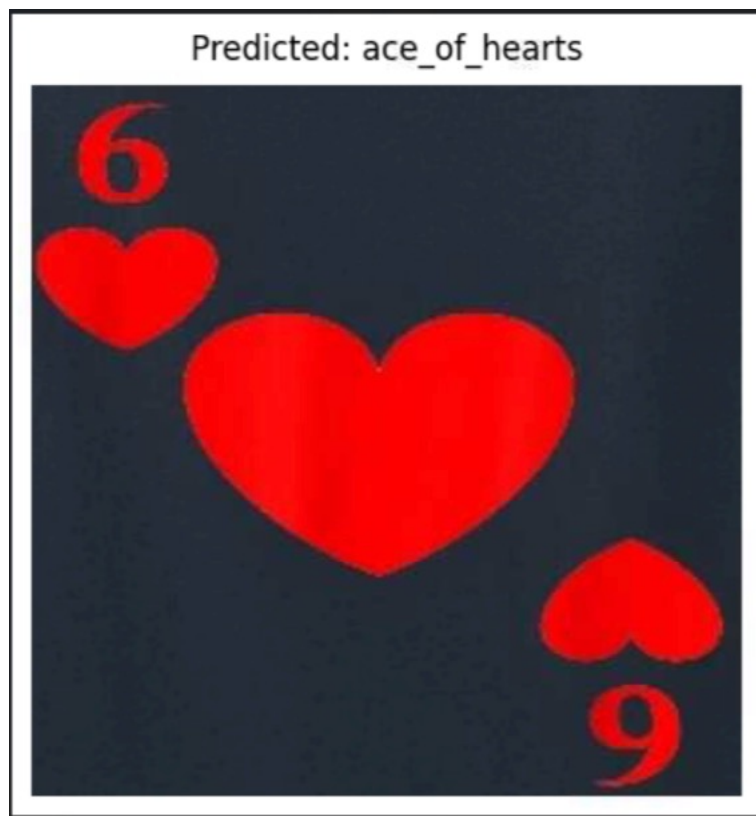
6 layers

Overlapping Convolutions _{3x3}

Seqpool

94.81 % Top-1

Implementation



Accuracy: 0.7445975277158949,
Precision: 0.7451575464672513,
Recall: 0.7445975277158949,
F1 score: 0.7384613355000814

