Topic: Escaping the big data paradigm with compact transformers

Abstract:

With the rise of transformers as the standard for language processing and their advancements in computer vision, there has been a corresponding growth in parameter size and amount of training data. Many have come to believe that because of this, transformers are not suitable for small sets of data. This trend leads to concerns such as: limited availability of data in certain scientific domains and the exclusion of those with limited resource from research in the field. In this paper, we aim to present an approach for small-scale learning by introducing compact Transformers. We show for the first time that with the right size, convolutional tokenization, transformers can avoid overfitting and outperform state of the art CNNs on small datasets. Our models are flexible in terms of model size, and can have as little as 0.28M parameters while achieving competitive results. Our best model can reach 98% accuracy when training on CIFAR-10 with only 3.7M parameters from scratch, which is a significant improve

in data efficiency over previous Transformer based models being over 10x smaller than other transformers and is 15% the size of ResNet 50 while achieving similar performance, CCT also outperforms many modern CNN based approaches and even some recent NAS based approaches. Additionally, we obtain a new SOTA result on flowers-102 with 99.76% top-1 accuracy, and improve upon the existing baseline on Imagenet as well as NLP tasks. Our simple and compact design for transformers makes them more frasible to study for those with limited computing resources dealing with small datasets while extending research efforts in data efficient transformers"

## Conclusion:

Transformers have commonly been percieved to be only applicable to large scale or medium scale training. While their scalability is undeniable, we have shown within this paper that with proper configuration, a transformer can be success-fully used in small data regimes as well and outperform convulational models of equivalent and even larger sizes. Our method is simple, flexible in size and the smallest of our variants can be easily loaded on even a minimal GPU or even a CPU. While part of which has been focused on large

scale models and data sets, we focus on smaller scales in which there is still much research to be done in data efficiency. we show that CCT can outperform other transformer based models on small datasets while also having a significant reduction in computational costs & memory constraints. This work demonstrates that transformers dont require vast computational resources and can allow for their applications in even the most modest of settings. This type of research is important to many scientstific domains where data is far more limited that the conventional machine learning datasets which are used in general research. continuing research in this direction will help open research up to more people and domains, extending machine learning research.