

Vision Transformer in Healthcare: Harnessing the Power and Unraveling the Trade-offs

Presented by Sadaf Farkhani



Danish Research
Center for
Magnetic
Resonance



Aim

- Introduce the fundamental concepts of Vision Transformer (ViT)
- Explore the application of ViT in the context of 3D medical images
- Conducting a comparison between CNNs and ViT
- ViT or CNN? That is the question :)

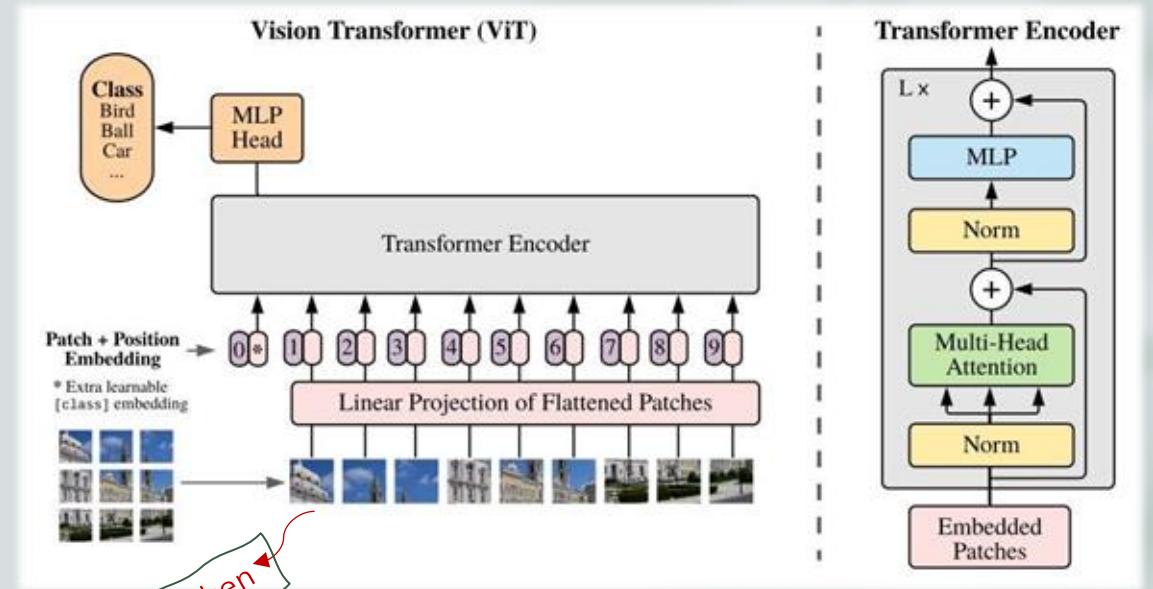
Vision Transformer (ViT)

- Position embedding
- Multi-head attention

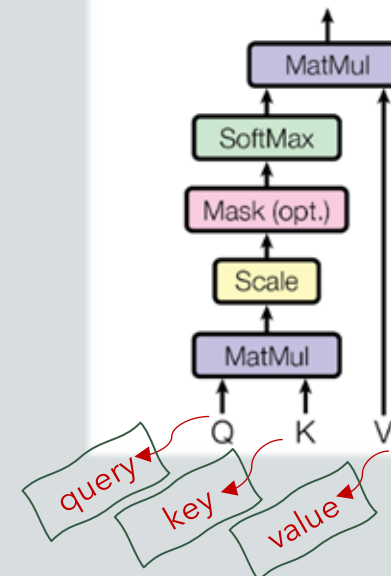
$$Attention = \text{Softmax}\left(\frac{QK^T}{d_k}\right)$$

- Feed-forward layers (MLP)

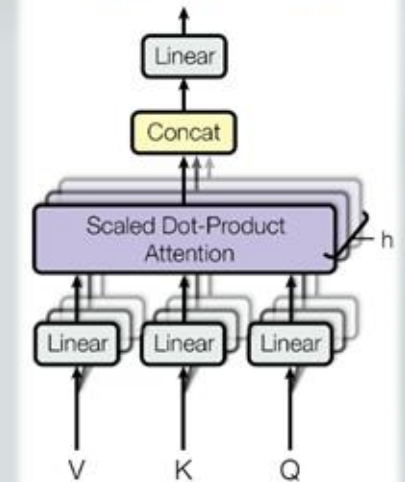
Dosovitskiy, Alexey, et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." International Conference on Learning Representations. 2020



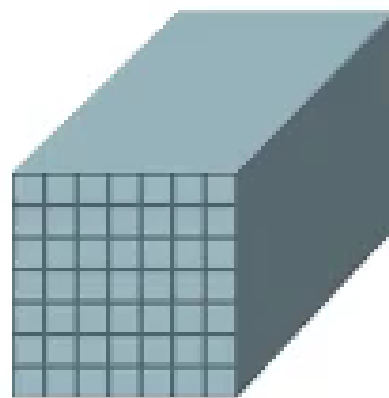
Scaled Dot-Product Attention



Multi-Head Attention

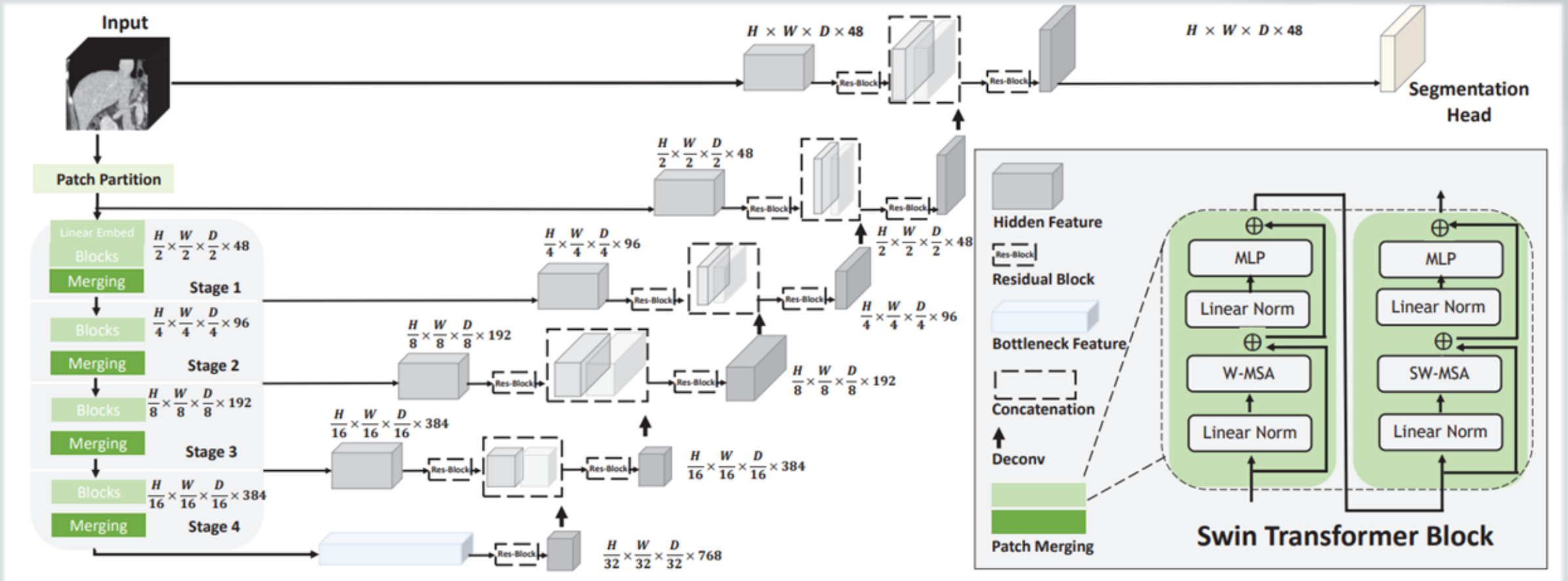


Multi-head Attention



Input tensor

3D Medical images (Swin-UNETR)

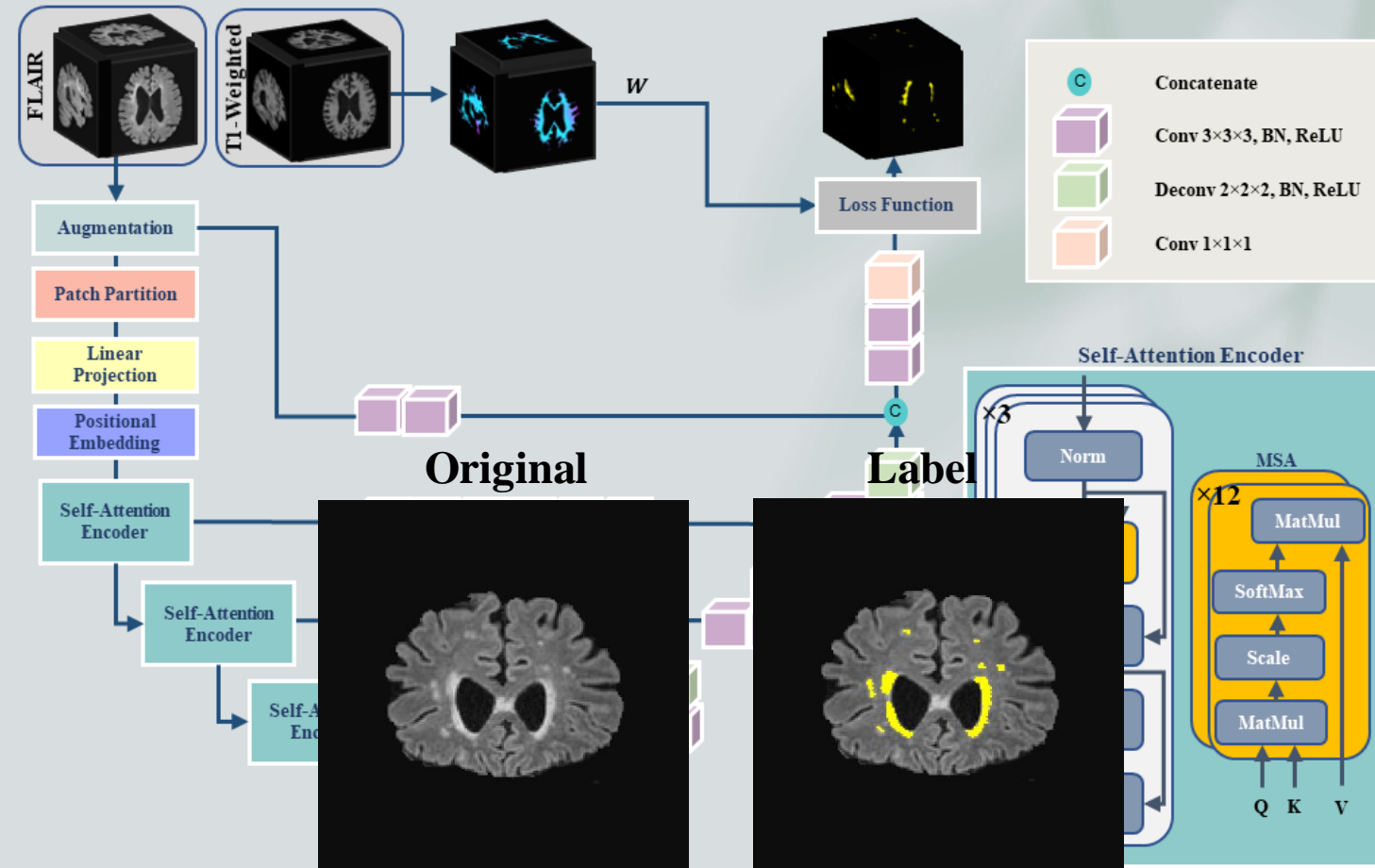


Tang, Yucheng, et al. "Self-supervised pre-training of swin transformers for 3d medical image analysis." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

3D Medical Images (VoSHT)

Task: brain lesion segmentation

- Sparse imbalanced lesions
- Multiple modalities are required
- More input modalities \rightarrow more parameters
- Weighted loss function



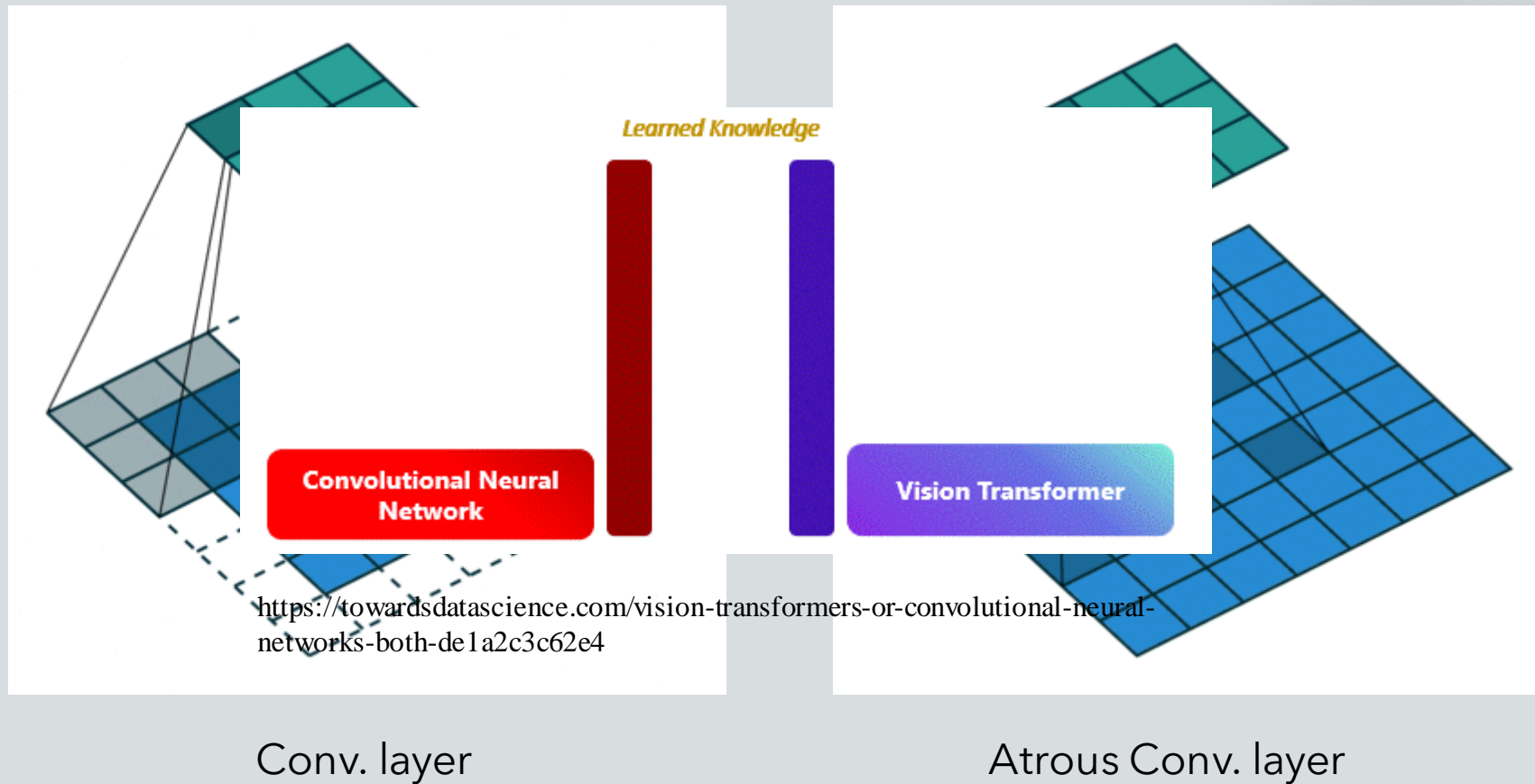
S. Farkhani, et. al, "End-to-end Volumetric Segmentation of White Matter Hyperintensities: Effect of Data, Model, and Loss Function", under review.

Results

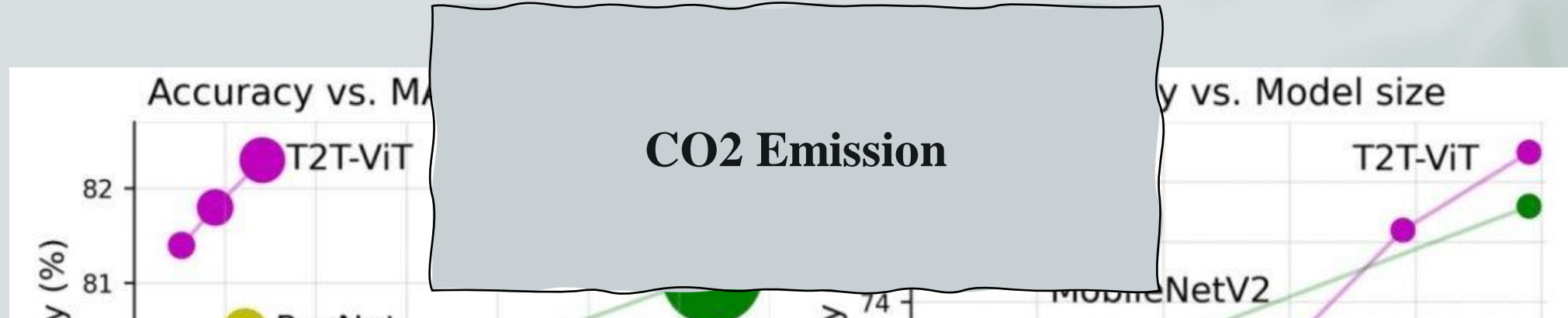
Out-of-distribution Testset

	METHOD		DSC	MSD	HD	RECALL	AUC-PR	
Metric	LST	DSC	71.71	6.30	22.84	76.61	46.92	AUC-PR
UNET	BIANCA	82.115	74.62	5.73	20.53	87.08	57.41	67.255
UNETR	UNET	78.811	84.25	2.28	12.59	86.62	70.19	62.374
TrUENet	UNETR	83.453	81.88	2.76	15.61	80.83	65.76	64.720
UNET+	TrUENet	85.533	80.94	2.61	13.79	76.63	66.50	72.376
VoSHT	UNET+	84.132	85.43	1.82	11.12	87.48	72.41	70.199
	VoSHT	85.59	85.59	1.76	10.30	87.48	72.67	

1. Global vs. Local



2. Parameters

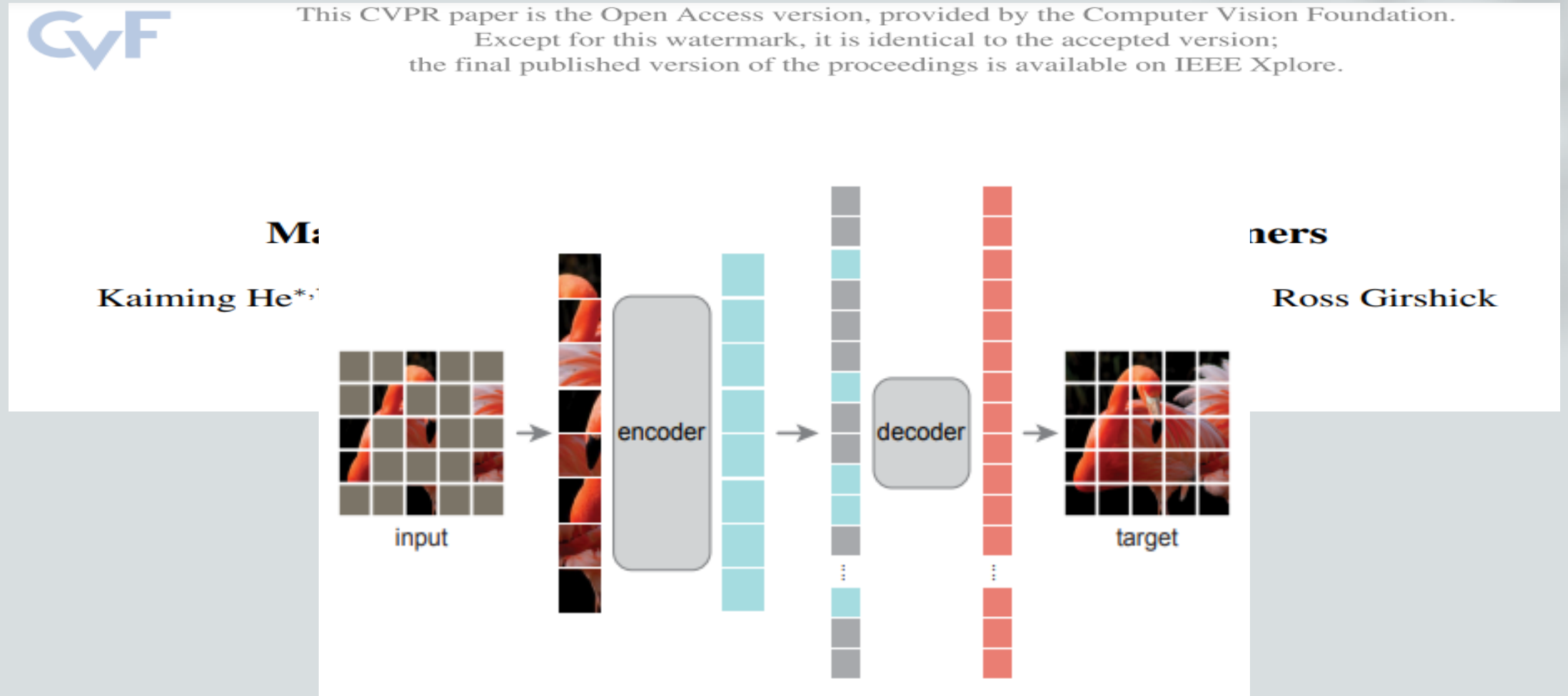


Carbontracker: Tracking and Predicting the Carbon Footprint of Training Deep Learning Models

Lasse F. Wolff Anthony^{*1} Benjamin Kanding^{*1} Raghavendra Selvan¹

Anthony, Lasse F. Wolff, Benjamin Kanding and Raghavendra Selvan "Carbontracker: Tracking and predicting the carbon footprint of training deep learning models" *arXiv preprint 2021* Xiv:2007.03051 (2020).

2. Parameters

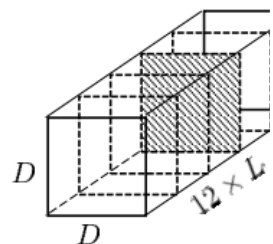


He, Kaiming, et al. "Masked autoencoders are scalable vision learners." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022.

2. Parameters

OpenReview

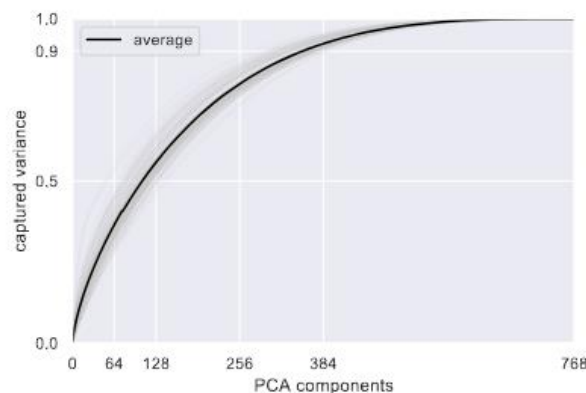
Login



I: raw weight matrices

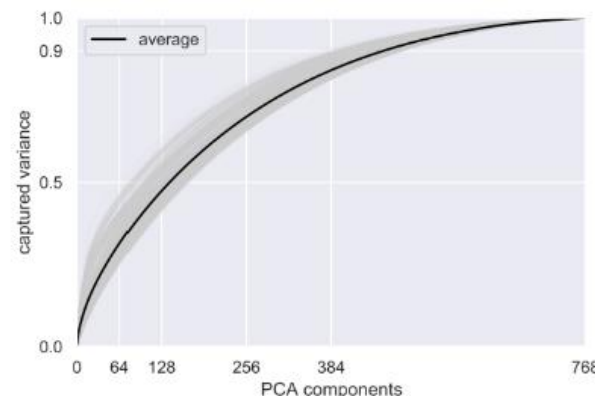
$$\mathbf{W}_i^I = \mathbf{W}_i$$

Intra-matrix

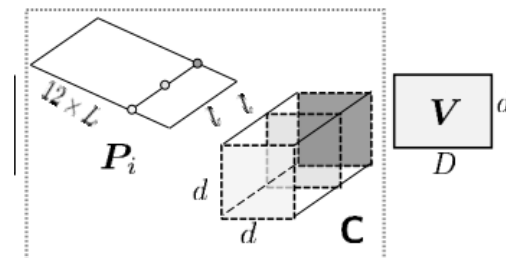


(a) PCA for each single weight matrix

Inter-matrix



(b) PCA for a pair of matrices along columns



V: tucker decomposition

$$\mathbf{W}_i^{IV} = \mathbf{U}(\mathbf{P}_i \mathbf{C}) \mathbf{V}$$

Published: 28 Jan 2022, Last modified: 28 Jan 2022

Figure 1: PCA for existing weight block matrices in BERT-base. We got nearly similar results in Fig. 5 for paired matrices along rows and columns, as shown in App. C.

Ren, Yuxin, et al. "Exploring extreme parameter compression for pre-trained language models." *arXiv preprint arXiv:2205.10036* (2022).

3. Generalization

In-distribution

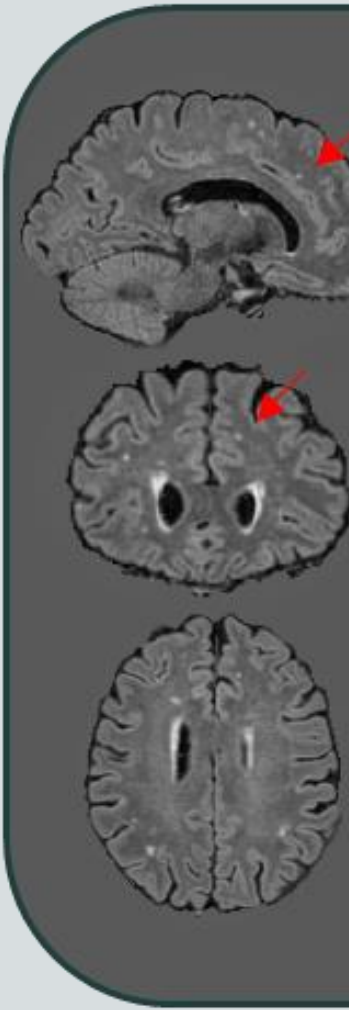
Are Transformers More Robust Than CNNs?

Published as a conference paper at ICLR 2023

Sagittal

Coronal

Axial



¹Johns Hopkins University
{ytongb

Transformers demonstrate superior generalization in recent work. In this paper, we compare the performance of Neural Networks and Vision Transformers on out-of-distribution samples. With our proposed modifications, Vision Transformers surprisingly perform better against adversarial attacks. While regular training on in-distribution samples suggests superior self-attention, we hope this work can help the community better understand the design of robust neural architectures. at <https://github.com/UCSC-VLAA/RobustCNN>.

n Test set

CAN CNNs BE MORE ROBUST THAN TRANSFORMERS?

Zeyu Wang¹ Yutong Bai² Yuyin Zhou¹ Cihang Xie¹
¹UC Santa Cruz ²Johns Hopkins University

ABSTRACT

The recent success of Vision Transformers is shaking the long dominance of Convolutional Neural Networks (CNNs) in image recognition for a decade. Specifically, in terms of robustness on out-of-distribution samples, recent research finds that Transformers are inherently more robust than CNNs, regardless of different training setups. Moreover, it is believed that such superiority of Transformers should largely be credited to their *self-attention-like architectures per se*. In this paper, we question that belief by closely examining the design of Transformers. Our findings lead to three highly effective architecture designs for boosting robustness, yet simple enough to be implemented in several lines of code, namely a) patchifying input images, b) enlarging kernel size, and c) reducing activation layers and normalization layers. Bringing these components together, we are able to build pure CNN architectures without any attention-like operations that are as robust as, or even more robust than, Transformers. We hope this work can help the community better understand the design of robust neural architectures. The code is publicly available at <https://github.com/UCSC-VLAA/RobustCNN>.

4. Real-world Data

- Imperfect labels
- Imperfect metrics

Dice Score
(DSC)

2 x

Prediction

Ground truth

Prediction

+

Ground truth

Hausdorff Distance
(HD)

$$\sup_{x \in X} \inf_{y \in Y} d(x, y)$$

$$\sup_{y \in Y} \inf_{x \in X} d(x, y)$$

False Negative

CNN vs. ViT

ViT	CNN
<ul style="list-style-type: none">- Globally-capturing dependencies- High number of parameters- High redundancy (within attention head and between layers of Transformers)- Higher demand for GPU/data accessibility- Higher CO2 emission- Perform better on remarkably big objects	<ul style="list-style-type: none">- Locally-capturing dependencies- Lower number of parameters- High redundancy among kernels- Lower demand for GPU/data accessibility- Perform better on small objects- Better generalization performance

Summary

- Necessity of pre
- Achieving balan
- Environmental i
model
- Priority on generalization: significance in real-world applications

Achieving Synergy: Harmonizing Between
CNN and Transformer

modules

ects of your

Reference

- [1] Dosovitskiy, Alexey, et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." International Conference on Learning Representations. 2020.
- [2] Tang, Yucheng, et al. "Self-supervised pre-training of swin transformers for 3d medical image analysis. " *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.
- [3] <https://colab.research.google.com/drive/1hXIQ77A4TYS4y3UthWF-Ci7V7vVUoxmQ?usp=sharing#scrollTo=twSVFOM9SopW>
- [4] Yuan, Li, et al. "Tokens-to-token vit: Training vision transformers from scratch on imagenet." *Proceedings of the IEEE/CVF international conference on computer vision*. 2021.
- [5] He, Kaiming, et al. "Masked autoencoders are scalable vision learners." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022.
- [6] Ren, Yuxin, et al. "Exploring extreme parameter compression for pre-trained language models." *arXiv preprint arXiv:2205.10036* (2022).
- [7] S. Farkhani, N. Demnitz, C.J. Baroxbekk, H. Lundell, H. R. Siebner, E. T. Petersen, K. H. Madsen, "End-to-end Volumetric Segmentation of White Matter Hyperintensities: Effect of Data, Model, and Loss Function", submitted to *Computer Methods and Programming in Biomedicine*, 2023.
- [8] Bai, Yutong, et al. "Are transformers more robust than cnns?." *Advances in neural information processing systems* 34 (2021): 26831-26843.
- [9] Wang, Zeyu, et al. "Can CNNs Be More Robust Than Transformers?." *The Eleventh International Conference on Learning Representations*. 2022.