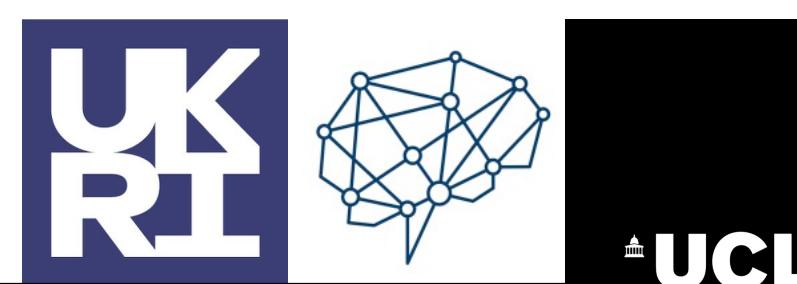


Setting the Record Straight on Transformer Oversmoothing

Gbetondji J-S Dovonon¹, Michael M Bronstein², Matt J Kusner¹

¹University College London

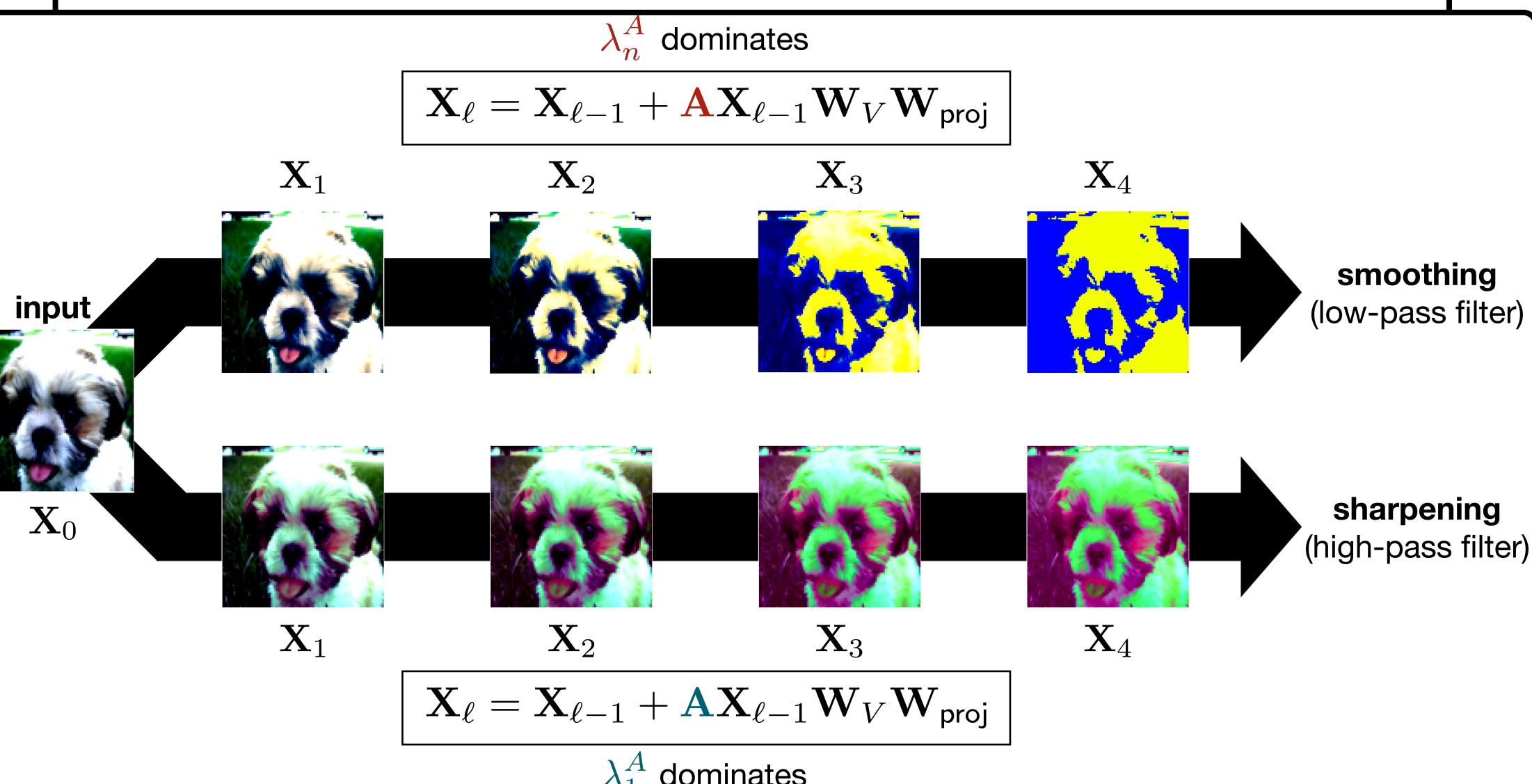
²University of Oxford



Recent work has argued that transformers are inherently low-pass filters. That is not always the case and we show why that matters for robustness and data efficiency.



WHAT IS OVERSMOOTHING?



Low Frequency component $LFC[\mathbf{X}] := (1/n)\mathbf{1}\mathbf{1}^\top \mathbf{X}$.
High Frequency component $HFC[\mathbf{X}] := (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top)\mathbf{X}$.

Oversmoothing happens when the LFC dominates the HFC

$$\lim_{L \rightarrow \infty} \frac{\|HFC[f^L(\mathbf{X})]\|_2}{\|LFC[f^L(\mathbf{X})]\|_2} = 0.$$

Issue: Smoothing reduces detail, making classification difficult.
How do vision transformers work if they oversmooth?
They learn not to.

We analyze a single self attention layer with a skip connection

$$\mathbf{A} := \text{Softmax}\left(\frac{1}{\sqrt{k}}\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^\top\mathbf{X}^\top\right)$$

$$\mathbf{X}_\ell := \mathbf{X}_{\ell-1} + \mathbf{A}_\ell \mathbf{X}_{\ell-1} \mathbf{W}_{V,\ell} \mathbf{W}_{\text{proj},\ell}$$

We can rewrite it in vector form using the vec operator and Kronecker product

$$\text{vec}(\mathbf{X}_\ell) = (\mathbf{I} + \underbrace{\mathbf{W}_{\text{proj}}^\top \mathbf{W}_V^\top \otimes \mathbf{A}}_{:= \mathbf{H}}) \text{vec}(\mathbf{X}_{\ell-1})$$

MAIN THEORETICAL RESULTS

$\lambda_1^A \leq \dots \leq \lambda_n^A$ and $\lambda_1^H \leq \dots \leq \lambda_d^H$ the eigenvalues of A and H

$(1 + \lambda_j^H \lambda_n^A)$ dominates \rightarrow over smoothing

$(1 + \lambda_j^H \lambda_1^A)$ dominates \rightarrow no over smoothing

$\lambda_j^H \in [-1, 0], \forall j$ guarantees we avoid over smoothing

TRANSFORMERS LEARN TO AVOID OVERSMOOTHING

Model	Random Init.		ImageNet		CIFAR100		
	ViT-B/16	ViT-B/16	DeiT-B/16	ViT-Ti	ViT-Ti ⁺	ViT-Ti ⁻	
$(1 + \lambda_j^H \lambda_n^A)$	100%	70.27%	57.82%	98.92%	100%	0%	
$(1 + \lambda_j^H \lambda_1^A)$	0%	29.73%	42.18%	1.08%	0%	100%	

Table 1: Distribution of dominating eigenvalues. We compare different models trained (or not) on ImageNet and CIFAR100 and count the percentage of cases where the dominating eigenvalue is $(1 + \lambda_j^H \lambda_n^A)$ or $(1 + \lambda_j^H \lambda_1^A)$.

A NEW PARAMETERIZATION

We parameterize H as a symmetric eigendecomposition

$$\mathbf{H} = \mathbf{V}_H \Lambda_H \mathbf{V}_H^\top$$

$$[\mathbf{V}_H, \mathbf{R}] := \text{QR}(\Theta) \text{diag}(\Lambda_H) := \pm(\psi^2)$$

↑
Learnable parameters

We refer to the model with negative eigenvalues Λ_H^- as ViT⁻ and the one with positive eigenvalues as ViT⁺

CIFAR100 RESULTS

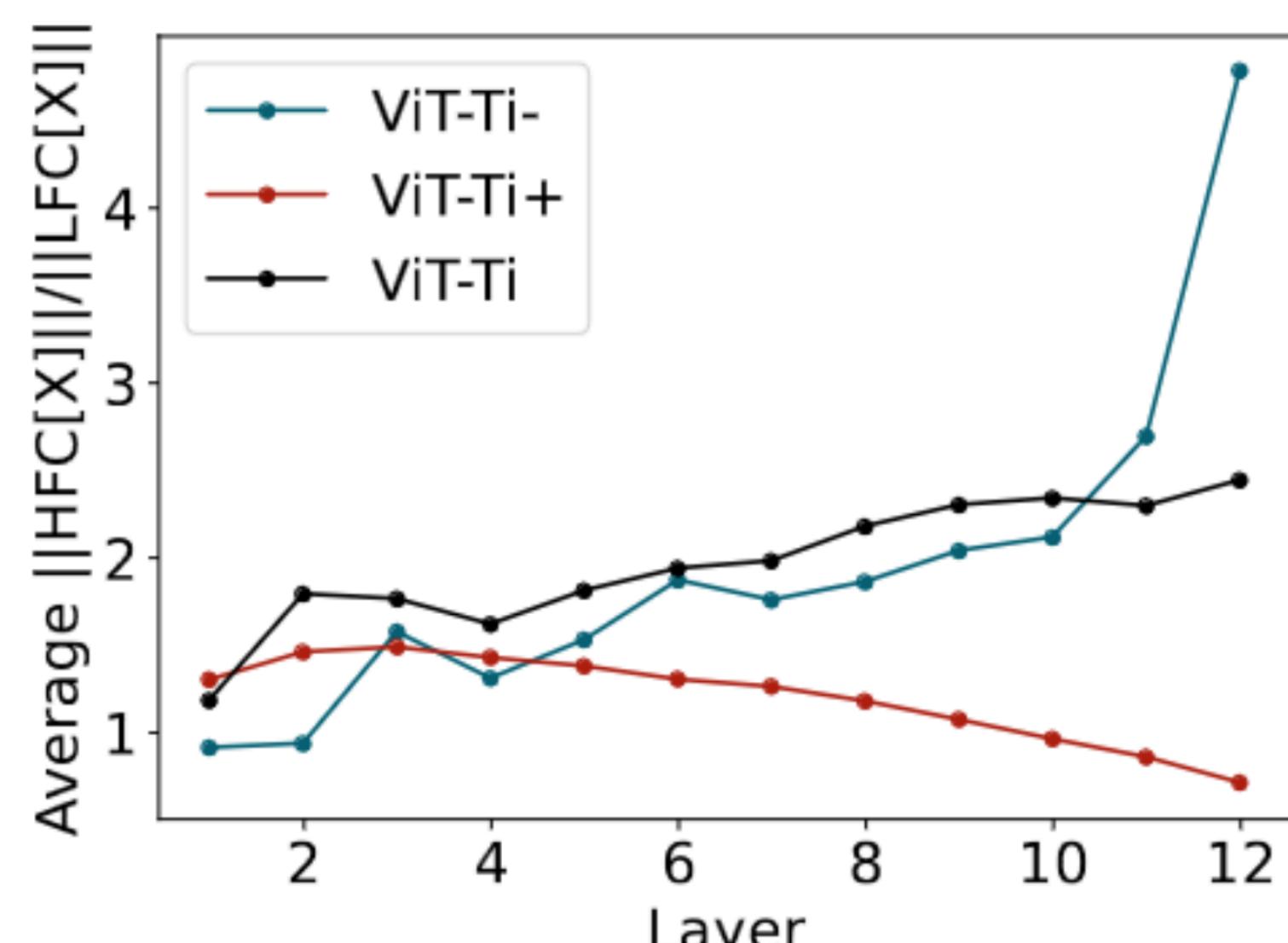
Data kept	100%	50%	10%
ViT-Ti	66.78 ± 0.1	56.38 ± 0.5	32.47 ± 0.9
ViT-Ti ⁺	64.34 ± 0.3	55.21 ± 0.5	28.67 ± 0.5
ViT-Ti ⁻	66.62 ± 0.3	56.58 ± 0.3	33.99 ± 0.4

Table 2: Data efficiency. Results for different percentages of the CIFAR100 used in training.

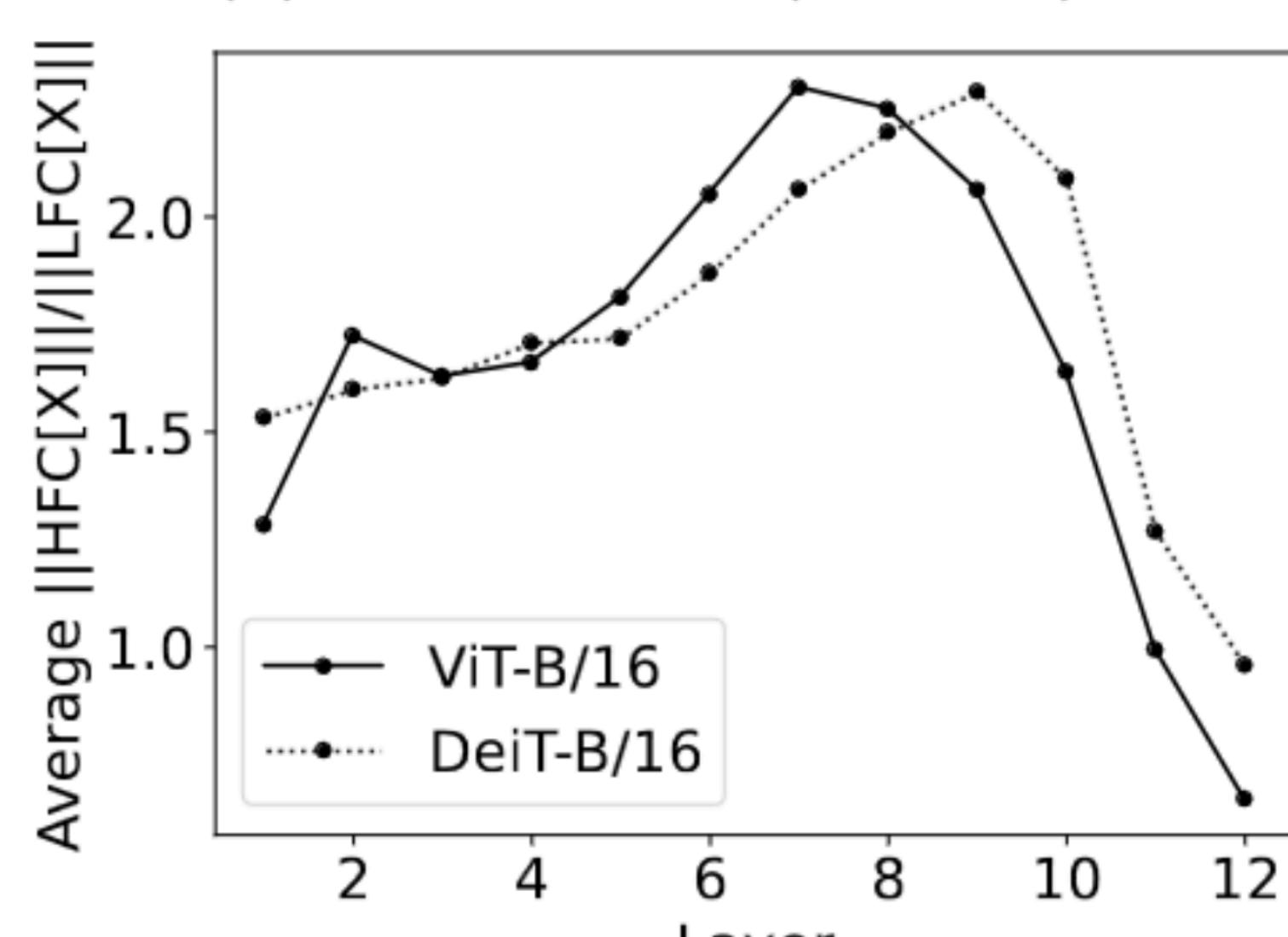
Corruption Intensity	0	1	2	3	4	5
ViT-Ti	66.78 ± 0.1	61.59 ± 0.2	57.75 ± 0.2	53.12 ± 0.1	49.32 ± 0.1	39.17 ± 0.2
ViT-Ti ⁺	64.34 ± 0.3	58.36 ± 0.4	53.83 ± 0.5	49.03 ± 0.6	45.15 ± 0.6	35.86 ± 0.6
ViT-Ti ⁻	66.62 ± 0.3	61.62 ± 0.2	58.00 ± 0.3	53.56 ± 0.2	49.70 ± 0.3	40.07 ± 0.2

Table 3: Corruption Robustness. Results for different corruption intensities on CIFAR100.

MEASURING OVERSMOOTHING



(b) CIFAR100 (trained)



(c) ImageNet (trained)