Foundation Models on a Budget: Approximating Blocks in Large Vision Models

Large vision foundation models (e.g., ViTs, DiNO, DEiT) achieve state-of-the-art performance but are extremely resource-hungry, limiting accessibility and sustainability. While compression methods like pruning, distillation, or early exiting exist, they typically require retraining or fine-tuning. This paper introduces **Transformer Block Approximation (TBA)**, a novel, training-free approach to reduce model size by exploiting **intra-network similarities**. Prior work mainly studied inter-model representation similarities (e.g., for model stitching), but this work shows that within a single model, multiple transformer blocks often produce **redundant representations**.

Our method works in two-stages by first identifying detect block pairs whose outputs are highly similar using a simple metric yet effective metric (i.e., MSE), then it replaces intermediate blocks with a simple closed-form **linear transformation estimated on a small subset** that maps the output of an earlier block directly to the later one. This removes parameters and computations without retraining or fine-tuning. See Figure 1 for a visualization of the method.

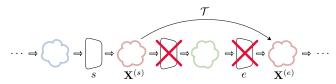


Figure 1: **Framework Description**. Given two latent spaces $\mathbf{X}^{(s)}$ and $\mathbf{X}^{(e)}$ representing the output of two blocks s and e for a random subset of n data points from the training set, we define a transformation matrix \mathcal{T} such that: $\mathbf{X}^{(e)} \approx \mathcal{T}(\mathbf{X}^{(s)})$.

Results in Table 1 demonstrate that TBA reduces parameter count while largely preserving representation fidelity. On image classification benchmarks (CIFAR-10/100, ImageNet-1k), TBA compresses different models with **minimal accuracy drop**, sometimes even **improving performance**. Additionally, learned transformations generalize across datasets, showing potential for privacy-sensitive domains where retraining is impossible.

Table 1: **Image Classification Performance Across Architectures.** Classification accuracy scores for DiNO-B and DEiT-S using multiple datasets, and 3 seeds. The "Approx." column specifies the blocks used for approximation, where the first value represents the block whose output is used to approximate the second block's output. The "Params." column shows the number of parameters removed by the approximation compared to the original model.

				Accuracy ↑	
	Approx.	Params.	CIFAR-10	CIFAR-100F	ImageNet1k
DiNO-B	1 →5	-26.5M	86.53 ± 0.21 (-11.28%)	62.11 ± 0.86 (-28.62%)	30.74 ± 0.45 (-58.57%)
	$\begin{array}{c} 2 \rightarrow 5 \\ 7 \rightarrow 10 \end{array}$	-19.5M -19.5M	$92.03 \pm 0.10 (-5.87\%) 93.41 \pm 0.34 (-4.46\%)$	$70.04 \pm 0.72 \text{ (-19.50\%)} $ $74.81 \pm 1.32 \text{ (-14.02\%)}$	$53.23 \pm 0.09 \text{ (-28.26\%)} $ $47.58 \pm 0.77 \text{ (-35.88\%)}$
	$ \begin{array}{c} 2 \to 4 \\ 9 \to 11 \\ 1 \to 2, 4 \to 5 \end{array} $	-13M -13M -13M	$\begin{array}{c} 96.31 \pm 0.18 \ (\text{-}1.49\%) \\ 87.01 \pm 0.30 \ (\text{-}10.76\%) \\ 96.33 \pm 0.18 \ (\text{-}1.47\%) \end{array}$	$81.25 \pm 0.57 \text{ (-6.62\%)}$ $72.65 \pm 1.86 \text{ (-14.36\%)}$ $81.69 \pm 0.07 \text{ (-6.11\%)}$	$68.49 \pm 0.12 (-7.70\%) 46.46 \pm 0.23 (-27.74\%) 68.42 \pm 0.20 (-7.79\%)$
	$\begin{array}{c} 0 \to 1 \\ 3 \to 4 \\ 9 \to 10 \\ 10 \to 11 \end{array}$	-6.5M -6.5M -6.5M -6.5M	$\begin{array}{c} \textbf{97.72} \pm 0.05 \ (\textbf{-0.05\%}) \\ 97.40 \pm 0.04 \ (\textbf{-0.38\%}) \\ 97.01 \pm 0.16 \ (\textbf{-0.78\%}) \\ 96.78 \pm 0.22 \ (\textbf{-1.01\%}) \end{array}$	$\begin{array}{c} \textbf{86.13} \pm 0.14 \ (\textbf{-1.01\%}) \\ 84.59 \pm 0.22 \ (\textbf{-2.78\%}) \\ 84.17 \pm 0.28 \ (\textbf{-3.26\%}) \\ 83.33 \pm 0.52 \ (\textbf{-4.23\%}) \end{array}$	
	original	86.58M	97.77 ± 0.24	87.01 ± 0.30	74.20 ± 0.68
DE1T-S	$1 \rightarrow 5$	-6.51M	78.35 ± 0.17 (-13.55%)	50.57 ± 0.29 (-28.90%)	$43.70 \pm 0.27 (\text{-}40.88\%)$
	$\begin{array}{c} 2 \rightarrow 5 \\ 7 \rightarrow 10 \end{array}$	-4.88M -4.88M	$85.73 \pm 0.31 \ (-5.41\%) 89.17 \pm 0.04 \ (-1.61\%)$	$60.55 \pm 0.16 \text{ (-14.87\%)} $ $69.15 \pm 0.33 \text{ (-2.78\%)}$	$62.04 \pm 0.21 \text{ (-16.05\%)} $ $57.48 \pm 0.06 \text{ (-22.24\%)}$
	$ \begin{array}{c} 2 \to 4 \\ 9 \to 11 \\ 1 \to 2, 4 \to 5 \end{array} $	-3.26M -3.26M -3.26M	$88.95 \pm 0.05 \text{ (-1.85\%)}$ $90.90 \pm 0.12 \text{ (+0.30\%)}$ $85.43 \pm 0.25 \text{ (-5.74\%)}$	$66.60 \pm 0.50 \text{ (-6.37\%)} $ $71.92 \pm 0.17 \text{ (+1.11\%)} $ $61.66 \pm 0.13 \text{ (-13.31\%)} $	$70.00 \pm 0.32 \text{ (-5.32\%)} $ $69.95 \pm 0.24 \text{ (-5.39\%)} $ $66.04 \pm 0.13 \text{ (-10.68\%)} $
	$\begin{array}{c} 0 \to 1 \\ 3 \to 4 \\ 9 \to 10 \\ 10 \to 11 \end{array}$	-1.63M -1.63M -1.63M -1.63M	$85.00 \pm 0.27 \text{ (-6.21\%)}$ $90.50 \pm 0.10 \text{ (-0.14\%)}$ $90.90 \pm 0.20 \text{ (+0.30\%)}$ $91.07 \pm 0.18 \text{ (+0.49\%)}$	$\begin{array}{c} 61.95 \pm 0.39 \ (\hbox{-}12.91\%) \\ 70.25 \pm 0.20 \ (\hbox{-}1.24\%) \\ 71.74 \pm 0.09 \ (\hbox{+}0.86\%) \\ \textbf{71.95} \pm 0.17 \ (\hbox{+}1.15\%) \end{array}$	$\begin{array}{c} 62.55 \pm 0.12 \ (\hbox{-}15.38\%) \\ 73.03 \pm 0.10 \ (\hbox{-}1.22\%) \\ 72.34 \pm 0.16 \ (\hbox{-}2.15\%) \\ \textbf{73.97} \pm 0.17 \ (\textbf{+}\textbf{0.05}\%) \end{array}$
	original	21.82M	90.63 ± 0.22	71.13 ± 0.26	73.93 ± 0.14

Conclusion

TBA is a practical, training-free method for simplifying foundation vision models by leveraging blockwise similarities. It improves computational accessibility and sustainability, offering a promising direction for efficient foundation models.