# SCHEME: SCALABLE CHANNEL MIXER FOR VISION TRANSFORMERS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031 032

034

040

041 042 043

044

045

046

047

048

Paper under double-blind review

# ABSTRACT

Vision Transformers have achieved impressive performance in many vision tasks. While the token mixer or attention block has been studied in great detail, much less research has been devoted to the channel mixer or feature mixing block (FFN or MLP), which accounts for a significant portion of the model parameters and computation. In this work, we show that the dense MLP connections can be replaced with a block diagonal MLP structure that supports larger expansion ratios by splitting MLP features into groups. To improve the feature clusters formed by this structure we propose the use of a lightweight, parameter-free, channel covariance attention (CCA) mechanism as a parallel branch during training. This enables gradual feature mixing across channel groups during training whose contribution decays to zero as the training progresses to convergence. As a result, the CCA block can be discarded during inference, enabling enhanced performance at no additional computational cost. The resulting Scalable CHannEl *MixEr* (SCHEME) can be plugged into any ViT architecture to obtain a gamut of models with different trade-offs between complexity and performance by controlling the block diagonal MLP structure. This is shown by the introduction of a new family of SCHEMEformer models. Experiments on image classification, object detection, and semantic segmentation, with different ViT backbones, consistently demonstrate substantial accuracy gains over existing designs, especially for lower complexity regimes. The SCHEMEformer family is shown to establish new Pareto frontiers for accuracy vs FLOPS, accuracy vs model size, and accuracy vs throughput, especially for fast transformers of small size.

## 1 INTRODUCTION



Figure 1: Comparison of the proposed SCHEMEformer family, derived from the Metaformer-PPAA-S12 model (52) with higher expansion ratios in the MLP blocks, and many SOTA transformers from the literature. The SCHEMEFormer family establishes a new Pareto frontier (optimal trade-off) for a) accuracy vs. FLOPs, b) accuracy vs model size, and c) accuracy vs, throughput. SCHEMEformer models are particularly effective for the design of fast transformers (throughput between 75 and 150 images/s) with small model size.See 5 for zoomed version.

Vision Transformers (ViTs) (11; 26; 42; 25) are now ubiquitous in computer vision. They decompose an image into a set of patches which are fed as tokens to a transformer model (41) of two

054 main components: a spatial attention module, which reweighs each token according to its similar-055 ity to the other tokens extracted from the image, enabling information fusion across large spatial 056 distances, and a *channel mixer* that combines the feature channels extracted from all patches using 057 a multi-layer perceptron (MLP or FFN). A bottleneck of this model is the quadratic complexity of 058 the attention mechanism on the number of patches. Numerous ViT variants have been proposed to address the problem, using improved attention mechanisms or hybrid architectures that replace attention or combine it with convolutions. Much less research has been devoted to the channel mixer. 060 Most models simply adopt the two-layer MLP block of (41), where channels are first expanded by a 061 specified *expansion ratio* and then compressed to the original dimension. This is somewhat surpris-062 ing since the mixer is critical for good transformer performance. For example, it is known that 1) 063 pure attention, without MLPs or residual connections, collapses doubly exponentially to a rank one 064 matrix (10), 2) training fails to converge without residual connections or MLP (52), and 3) replacing 065 MLPs with more attention blocks (both spatial and channel attention) of equivalent computational 066 complexity lowers the transformer accuracy (8). All these observations indicate that the channel 067 mixer is an indispensable ViT component. 068

In this work, quantify how much the channel mixer module contributes to ViT performance and 069 investigate how to improve the trade-off between complexity and accuracy of the ViT model. We show that enhanced design of the channel mixer can lead to significant improvements in transformer 071 performance by introducing a novel Scalable CHannEl MixEr (SCHEME) that enables the design 072 of models with larger expansion ratios. SCHEME is a generic channel mixer that can be plugged 073 into existing ViT variants to obtain effective scaled-down or scaled-up model versions. We replace 074 the channel mixer of a state of the art architecture (SOTA) for low-complexity transformers, the 075 MetaFormer-PPAA-S12 (52), with SCHEME to obtain a new family of SCHEMEformer models with improved accuracy/complexity trade-off. This is illustrated in Figure 1, where the SCHEME-076 former family is demonstrated to establish new Pareto frontiers for accuracy vs FLOPS, accuracy vs 077 model size, and accuracy vs throughput, showing that SCHEME allows fine control over all these variables, while guaranteeing SOTA performance. These properties are shown to hold for image 079 classification, object detection and semantic segmentation tasks, as well as for different architectures, such as T2T-ViT (54), CoAtNet (5), Swin Transformer (25), CSWin (9) and DaViT (8). 081

To develop SCHEME, we start by studying the impact of the mixer channel expansion. Transformer performance is shown to increase with expansion ration until it saturates for an expansion 083 ratio beyond 8. However, because mixer (MLP) complexity also increases with the dimension of the 084 intermediate representation, naive channel scaling with expansion ratios larger than 4 causes an ex-085 plosion of parameters and computation, leading to models of large complexity and prone to overfit. To achieve a better trade-off between dimensionality and computation, we leverage recent findings 087 about the increased hardware efficiency of mixers with block diagonal structure (1; 6). These group 880 the input and output feature vectors of a layer into disjoint subsets and perform matrix multipli-089 cations only within each group, as illustrated in Figure 2. We denote the resulting MLP block as Block Diagonal MLP (BD-MLP). Despite the lack of feature mixing across groups, we find that 091 a transformer model equipped with the BD-MLP and a larger expansion ratio (to match both the 092 parameter count and computation) achieves comparable or slightly higher accuracy than a baseline with dense MLP. This suggests that the lower accuracy of the block diagonal operations is offset by 093 the gains of larger expansion ratios. Further analyzing the features learned by the different groups 094 of the BD-MLP, we observe that they form feature clusters of similar ability to discriminate the tar-095 get classes. To learn better feature clusters, we seek a mechanism capable of restoring inter-group 096 feature communication during training without increasing parameters. For this, we propose a channel attention branch that reweighs the input features of the BD-MLP according to their covariance 098 matrix, as illustrated in Figure 2. This attention mechanism is denoted as the channel covariance attention (CCA) block. The re-weighted features are then fused with the BD-MLP output by means 100 of a weighted residual addition with learned weights  $(\alpha, 1 - \alpha)$ . 101

As shown in Figure 2, SCHEME combines the sparse block diagonal structure of the BD-MLP, and the parameter-free CCA attention module, to implement a channel mixer extremely efficient in terms of parameters. Ablations of the evolution of the fusion weight  $1 - \alpha$  learned in the CCA branch (see Fig. 2) over training show that it gradually decays to zero during training. This happens consistently across all layers of the model and across model architectures. Hence, while the CCA is important for the formation of good feature clusters during training, it can be removed at inference without any loss, as illustrated in Table 3 and Figure 2. As a result, the model accuracy improves over simply



117

128 129

130

131 132

133

134

135

136 137

138

139 140

141

Figure 2: Proposed SCHEME channel mixer. The channel mixer of the standard transformer consists of two MLP layers, performing dimensionality expansion and reduction by a factor of E. SCHEME uses a combi-118 nation of a block diagonal MLP (BD-MLP), which reduces the complexity of the MLP layers by using block 119 diagonal weights, and a channel covariance attention (CCA) mechanism that enables communication across 120 feature groups through feature-based attention. This, however, is only needed for training. The weights  $1 - \alpha$ decay to zero upon training convergence and CCA can be discarded during inference, as shown on the right. 121 Experiments show that CCA helps learn better feature clusters, but is not needed once these are formed. 122

using the BD-MLP mixer, but inference complexity does not. This leads to an extremely efficient 123 inference setup, both in terms of parameters and FLOPs. Overall, the paper makes the following 124 contributions, 125

- a study of the channel mixer of ViT MLPs, showing that dense feature mixing can be replaced by sparse feature mixing of higher internal feature dimensionality for improved accuracy, without increased complexity.
  - the SCHEME module, which combines 1) a BD-MLP to enable internal feature representations of larger dimensionality than previous MLP blocks, and 2) CCA to enable the learning of these representations without cost at inference.
- various models that combine SCHEME with previous transformer architectures to achieve SOTA trade-offs between accuracy and model size, FLOPS, or throughput, such as the SCHEMEformer of Figure 1. This is shown particularly effective for the design of fast transformers with small model size, of interest for edge devices, robotics, and low-power applications.
  - Experiments on image classification, object detection, and semantic segmentation, showing consistent gains in accuracy for fixed computation and size.

#### 2 **RELATED WORK**

142 **Vision Transformers:** Vision transformers advanced the SOTA in several vision tasks since (26; 11) 143 successfully applied the transformer-based self-attention NLP model of (41) to image generation 144 and classification tasks. These models rely on a spatial attention mechanism, based on the matrix of 145 dot-products between features extracted from image patches. This has quadratic complexity in the 146 number of patches and is quite intensive. Most follow up work (15; 33; 53; 38; 34; 46; 35; 54; 55; 40) 147 improved the spatial attention mechanism of ViT. DeiT (35) and subsequent works (36: 37) intro-148 duced a distillation token to distill information, typically from a CNN teacher, without large amounts of data or compute. PvT (42) proposed a progressive shrinking pyramid architecture with spatial-149 reduction attention that scales ViTs for dense prediction tasks beyond image classification. Swin 150 transformers (25) introduced a hierarchical shifted window attention mechanism, which reduces the 151 complexity to linear with respect to the number of windows. HaloNet (40) proposed two extensions 152 for local ViTs (27) with blocked local attention and relaxed translational equivariance for scaling 153 ViTs. A more extensive review of ViTs is given in (24). Recently, several works have shown that 154 spatial attention is not the critical ViT feature. Some works improved performance by relying on hybrid architectures, which augment or replace ViT layers with convolutions (45; 39). Efficient 156 transformer designs, for edge devices, frequently sacrifice spatial attention to achieve better trade-157 offs between FLOPs and accuracy (2; 21). Other works have questioned the need for spatial attention 158 altogether. Metaformer (52) argued that the fundamental trait of ViT is the mixing of information across patches, showing that competitive results can be obtained by simply replacing attention with 159 pooling or identity operations. Similarly, (23) showed that spatial attention is not critical for vision 160 transformers by proposing a spatial-gating MLP of comparable performance to ViT (11). DaViT (8) 161 showed that spatial attention is helpful by reusing its design for channel attention, building a cascade

162						SCH	EME					
163		Ratio	FLOPs	Top-1 Acc		44-e8	12-e8		CCA	Params	FLOPs	Top-1 Acc
105	Model	(E)	(G)	(%)	S12 P (M)	12	21	SCHEME	Used	(M)	(G)	(%)
164	Metaformer*	8	4.14	81.6 (+0.6)	S12 F (G)	1.8	3.3	44-e8-S12	$\checkmark$	11.83	2.16	79.74
165	Metaformer*	6	3.35	81.8 (+0.8)	S24 P (M)	21	40	44-e8-S12		11.83	1.77	79.72
100	Metaformer (52)	4	2.55	81.0 (+0.0)	S24 F (G)	33	6.5	12-e8-S24	1	40.0	73	82.80
166	Metaformer*	2	1.77	78.9 (-2.1)	S36 P (M)	55	59	12-e8-S24	•	40.0	6.5	82.00
167	Metaformer*	1	1.33	76.0 (-5.0)	\$36 F (G)	80	96	CoatNet-44-e8	.(	17.80	3.83	80.70
107	Metaformer-S18*	1	2.51	78.3 (-2.7)	T 11 0 00		<u> </u>	ContNet 44 e8	•	17.80	3 12	80.68
168	Table 1. Met	aForr	ner-S	12(52)	Table 2: SC	HEM	E param-	Swin 12 of T		26.02	7.00	81.60
100	ImagaNat 1K	voli	dation	12 (32)	eters (P) and	nd FL	OPs (F).	Swin-12-co-1	v	26.02	5.00	81.69
169	Inagenet-IK	vano	Jation	accu-	44-e8 (12-e	e8) de	ownscales	3wiii-12-00-1		30.93	5.69	81.09
170	racy vs MLP ex	rpans	ion ra	tio $(E)$ .	(umagala) th	$\sim M$	ato Dormor	Table 3: In	ipact	of rei	novin	g CCA
	*: results by au	ithor	code.		(upscale) u	le Mi	etarormer	branch durir	ng in	ference	э.	
1/1	· · · · · · · · · · · · · · · · · · ·				model.				0			

172 of alternating spatial and channel attention blocks. While this improves performance, it increases 173 the complexity of the transformer block, resulting in a model with many parameters and potentially 174 redundant channel mixer operations. Despite all this work on ViT architectures, little emphasis has 175 been devoted to the channel mixer module (MLP) that follows attention. This is surprising because the mixer dominates both the parameter count and complexity (FLOPs) of the standard transformer 176 block. (22; 53) modify the MLP block to mimic the inverted residual block of the MobileNetV2 177 (29), by adding a depthwise convolution. This improved performance but increases parameter and 178 computation costs. Switch Transformer (13) replaces the FFN with a sparse mixture of experts (32) 179 that dynamically routes the input tokens. This design allows scaling models to large sizes using 180 higher number of experts but is not effective for ViTs. XCiT (12) employs a cross-covariance at-181 tention (XCA) operator, which can be seen as a "transposed" version of self-attention that operates 182 across feature channels. The architecture of XCiT is composed of three primary components: the 183 core XCA operation, a local patch interaction (LPI) module, and a feedforward network (FFN). The XCA mechanism computes the covariance operation across different head groups, akin to multi-head 185 attention. In this work, we extend this concept by proposing a Cross-Covariance Channel Attention (CCA) operation to facilitate feature mixing across different channel groups. Unlike XCiT, which uses "heads" for interaction and projection matrices for queries, keys and values, CCA leverages 187 the full feature set to compute covariance without any projections, ensuring a more comprehensive 188 representation of inter-feature interactions. In this work, we propose an efficient and generic channel 189 mixer module (BD-MLP and CCA) that improves both the parameter and computational efficiency 190 of the transformer and allows for flexible scaling of ViTs. 191

192 193

194

# 3 THE SCHEME MODULE

Figure 2 depicts the proposed SCHEME module for feature mixing in ViTs. As shown on the left, the standard channel mixer consists of two MLP layers, which expand the dimensionality of the input features and then reduce it to the original size. Let  $\mathbf{x} \in \mathbb{R}^{d \times N}$  be the matrix containing the *N d*-dimensional input feature vectors extracted from *N* image patches. The mixer computes an intermediate representation  $\mathbf{z} \in \mathbb{R}^{Ed \times N}$  and an output representation  $\mathbf{y} \in \mathbb{R}^{d \times N}$  according to

$$\mathbf{z} = \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \mathbf{1}_N^T) \tag{1}$$

(2)

201 202 203

204

205

200

where  $\mathbf{W}_1 \in \mathbb{R}^{Ed \times d}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{d \times Ed}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{Ed}$ ,  $\mathbf{b}_2 \in \mathbb{R}^d$ ,  $\mathbf{1}_N$  is the *N*-dimensional vector containing ones as all entries,  $\sigma(.)$  is the activation function (typically GELU), and *E* is an expansion factor, typically 4.

 $\mathbf{y} = \mathbf{W}_2 \mathbf{z} + \mathbf{b}_2 \mathbf{1}_N^T$ 

206 Table 1 details the impact of the mixer on overall transformer performance by evaluating the role 207 of the expansion factor E on the performance of Metaformer-PPAA-S12 (52) (Pooling, Pooling, 208 Attention, Attention) architecture on ImageNet-1K. Classification accuracy increases from 76.0 to 209 81.8%, as E ranges from 1 to 6, decreasing for E = 8, which suggests overfitting. The table 210 also shows that these gains are not trivial. The S18 model, which has no expansion (E = 1) but 211 more transformer layers and complexity comparable to that of the S12 model with E = 4, has an 212 accuracy 2.7 points lower than the latter. In summary, for a given computation budget, it is beneficial 213 to trade off transformer depth for dimensionality expansion in the channel mixer. This shows that this expansion is a critical component of the transformer architecture. On the other hand, naively 214 scaling E beyond 6 severely increases parameters and computation, leading to models that over-fit 215 and are impractical for many applications.

# 216 3.1 SCALABLE CHANNEL MIXER (SCHEME)217

224

225 226

244

245 246

247

248 249

**Block Diagonal MLP (BD-MLP):** Block diagonal matrices have been previously used to efficiently approximate dense matrices (1; 6). In CNNs, group channel operations are frequently used to design lightweight mobile models with improved accuracy-computation trade-off (19; 3; 29). This consists of splitting the feature vectors of (1)-(2) into disjoint groups, e.g. x into a set of  $g_1$  disjoint features  $\{x_k\}_{k=1}^{g_1}$  where  $x_k \in \mathbb{R}^{d/g_1 \times N}$ , and y into a set  $\{y_k\}_{k=1}^{g_2}$  where  $y_k \in \mathbb{R}^{Ed/g_2 \times N}$ . As illustrated in Figure 2, the MLPs of (1)-(2) are then implemented independently for each group, according to

$$\mathbf{z}_k = \sigma(\mathbf{W}_{1,k}\mathbf{x}_k + \mathbf{b}_{1,k}\mathbf{1}_N^T)$$
(3)

$$\mathbf{y}_k = \mathbf{W}_{2,k} \mathbf{z}_k + \mathbf{b}_{2,k} \mathbf{1}_N^T \tag{4}$$

227 where  $\mathbf{W}_{1,k} \in \mathbb{R}^{Ed/g_1 \times N/g_1}, \mathbf{W}_{2,k} \in \mathbb{R}^{d/g_2 \times Ed/g_2}, \mathbf{b}_1 \in \mathbb{R}^{Ed/g_1}, \mathbf{b}_2 \in \mathbb{R}^{d/g_2}$  and  $\mathbf{z}$  is decom-228 posed into a set  $\{\mathbf{z}_k\}_{k=1}^G$  where  $\mathbf{z}_k \in \mathbb{R}^{Ed/G \times N}$ , with  $G = g_1$  in (3) and  $G = g_2$  in (4). Since 229 the complexity of (3) is  $g_1^2$  times smaller than that of (1) and there are  $g_1$  groups, the complexity 230 of the first MLP is  $1/g_1$  times that of standard MLP. Similarly, the complexity of the second MLP 231 is  $1/g_2$  times that of the standard MLP. Hence, a transformer equipped with the BD-MLP and ex-232 pansion factor  $\frac{2g_1g_2}{g_1+g_2}E$  has identical complexity to a standard transformer of factor E. For example, 233 when  $g_1 = g_2 = g$  this allows growing the expansion factor by a factor of g without computational 234 increase. 235

Channel Covariance Attention (CCA): While the introduction of groups enables accuracy gains 236 due to the increased expansion factor by  $\frac{2g_1g_2}{g_1+g_2}$ , it results in sub-optimal features. This is because 237 the features in the different groups of (3)-(4) are processed *independently*, i.e. there is no inter-group 238 feature fusion. This reduces the efficiency of the BD-MLP. To enable feature mixing between all 239 feature channels and thus induce the formation of better feature clusters, we introduce a covariance 240 attention mechanism in a parallel branch, as illustrated in Figure 2. The input features are first 241 transposed to obtain the  $d \times d$  covariance matrix<sup>1</sup>  $\mathbf{x}\mathbf{x}^T$ . This is then used to re-weigh the input 242 features by their covariance with other feature channels, using 243

$$CCA(\mathbf{x}) = \operatorname{softmax}\left(\frac{\mathbf{x}\mathbf{x}^{T}}{\tau}\right)\mathbf{x}$$
 (5)

where the softmax operation is applied across the matrix rows and  $\tau$  is a smoothing factor. The output of the channel mixer block is the weighted summation of the BD-MLP and CCA branches

$$\mathbf{y}_{out} = \alpha \mathbf{y} + (1 - \alpha) CCA(\mathbf{x}),\tag{6}$$

where  $\alpha$  is a mixing weight learned across all samples. Various other design choices are discussed in Section 4.2.

253 CCA as a Regularizer: The introduction of a parameter free attention branch and a learnable weight  $\alpha$  allows the model to form better feature clusters during training and gradually decay the contri-254 bution from CCA branch once the feature clusters are formed. This can be seen in Figure 3, which 255 plots the value of the learned mixing weight  $1 - \alpha$  as a function of training epochs on ImageNet-1K, 256 for all transformer layers. These plots are typical of the behavior we observed with all transformer 257 backbones and architectures we considered. Clearly,  $1 - \alpha$  starts with high to intermediate values, 258 indicating that information flows through both branches of the mixer, but decays to  $1 - \alpha \approx 0$ 259 as training converges. Hence, as shown in Table 3, there is no degradation if the CCA branch is 260 removed during inference. This eliminates a substantial amount of computation during inference, 261 leading to the training and inference setup of Figure 2, where CCA is not used at inference. 262

We explain this behavior by conjecturing that the downside of the computational efficiency of the BD-MLP is a more difficult learning problem, due to the independent processing of channel groups. This creates symmetries in the cost function, e.g. the order of the feature groups is not important, and requires a feature clustering operation that is likely to produce more local minima. The CCA branch helps to smooth out this cost function during training, while the feature groups are not established, by allowing inter-group communication. However, once the right feature groupings are found, CCA

<sup>&</sup>lt;sup>1</sup>Since the features are normalized before the mixer, i.e centered such that  $\mathbf{x} \mathbf{1}_N = 0$ ,  $\mathbf{x} \mathbf{x}^T$  is the covariance matrix of features  $\mathbf{x}$ .



277 Figure 3: Impact of CCA (SCHEMEformer-44-e8- Figure 4: SCHEME tradeoffs. Left: Accuracy vs S12). Left: Evolution of weight  $1 - \alpha$  across model FLOPs of various SCHEME models with different 278 layers. Right: Class separability of output features y MLP configurations. Right: SCHEME mixer im-279 (over 50 random classes of ImageNet-1K validation proves accuracy for fixed throughput or vice-versa for 280 set) for model trained with and without CCA. See 6 various popular ViT architectures. See 7 for zoomed 281 for zoomed version.

version.

282 is no longer needed, and a simple BD-MLP mixer has no loss of performance over the standard MLP. 283 Note that the operation of (5) is basically a projection of x into canonical subspaces of features that 284 are correlated in the input image. This is likely to be informative to guide the group formation, but 285 less useful when the features are already clustered. While this hypothesis is not trivial to test, since 286 (5) varies from example to example, we confirmed that using CCA during training enhances class 287 separability, which likely reduces overfitting for large expansion ratios. See the section 4.2 for more 288 details.

289 **Computational Complexity:** The complexity of the BD-MLP block is controlled by the group 290 numbers  $g_1$ ,  $g_2$  and the expansion factor E, with a total cost  $\mathcal{O}(Ed^2/g_1 + Ed^2/g_2)$ , where d is 291 the channel dimension. The computational cost of CCA is  $\mathcal{O}(Nd^2)$  where N is the number of 292 tokens. Since CCA is not used during inference, it only adds to the computations during training. 293 The SCHEME framework provides a systematic way to control the trade-off of transformer width vs depth, by controlling the block size and expansion hyperparameters. 294

295 296 297

# 3.2 THE SCHEMEFORMER FAMILY

The proposed SCHEME module enables efficient control of model complexity via the mixer hy-298 perparameters  $g_1, g_2$ , and E. Table 1 shows that naively scaling down the ViT model by simply 299 reducing E causes a significant accuracy loss. The SCHEME module allows much more effective 300 control of the accuracy/complexity trade-off, producing models of better performance for a fixed computational budget. This is demonstrated by the introduction of a new family of models, denoted 302 as SCHEMEformer, obtained by replacing the channel mixer of the Metaformer-PPAA (52) archi-303 tecture with the SCHEME module. Two such configurations are shown in Table 2 where the naming 304 follows the convention {model-name}-{ $g_1g_2$ }-e{E} where the model name is skipped for brevity.

305 306 307

308

309

301

#### 4 **EXPERIMENTAL RESULTS**

4.1 COMPARISONS TO THE STATE OF THE ART

310 **Image Classification:** Image classification is evaluated on Imagenet-1K, without using extra data. 311 We report the results with single crop top-1 accuracy at  $224 \times 224$  input resolution. We evaluate the 312 SCHEMEformer family of models based on the Metaformer-PPAA-S12 (52) obtained by replacing 313 the MLP of the latter with the SCHEME module. Refer to Appendix for implementation details.

314 We start by evaluating how this improves the trade-off between model accuracy and complexity. 315 As discussed in Section 3.1, when  $q_1 = q_2$ , a SCHEME transformer of expansion factor qE has 316 identical complexity to a standard transformer of expansion factor E. Hence, for fixed FLOPS, 317 SCHEME allows an increase of the expansion factor by g. Figure 4 a) compares the performance 318 of the Metaformer-PPAA-S12 with expansion ratios  $E \in \{1, 2, 4, 8\}$  to comparable variants of the 319 SCHEMEformer-PPAA-S12, with SCHEME mixers of either g = 2 (green curve) or g = 4 (blue 320 curve) groups. The SCHEMEformer models have a better trade-off between accuracy and FLOPS, 321 achieving higher accuracies for all complexity levels. Among these, the one with more feature groups (g = 4) has the best performance. While SCHEMEformer gains are observed for all FLOP 322 levels, they are larger for lower complexity models. This makes SCHEME particularly attractive for 323 the design of low complexity transformers, e.g. for edge devices or equivalent applications.

324 325	Model	#Params	FLOPs (G)	Thru (im/s)	Top-1 Acc
000		(111)	(0)	(111/3)	(,0)
326	gMLP-Ti (23)	6	1.4	-	72.3
327	ViT-L/16 (11)	307	63.6	37	76.1
200	Meta-11-e2-S12 (52)	12	1.8	133	78.9
320	MogaNet-T (20)	5	1.10	44	79.0
329	XCiT-T24 (12)	12	2.3	-	79.4
220	ViT-B/16 (11)	86	17.6	112	79.7
330	SCHEME-44-S12	12	1.77	133	79.7
331	$S^2$ -MLP-deep (51)	51	10.5	-	80.7
222	Meta-S12 (52)	17	2.6	87	81.0
332	Swin-Tiny (25)	29	4.5	100	81.3
333	T2T-ViT t-14 (54)	22	6.1	70	81.7
22/	DeiT-B (35)	86	17.5	114	81.8
334	ViL-Small (56)	25	5.1	-	82.0
335	SCHEME-12-S12	21	3.35	130	82.0
226	Focal-Tiny (50)	29	4.9	29	82.2
330	CPVT-Base (4)	88	17.6	-	82.3
337	DaViT-Tiny (8)	23.0	4.3	61	82.8
338	CSWin-Tiny (9)	23.0	4.3	20	82.8
550	SCHEME-12-S24	40	6.47	69	82.8
339	XCiT-L24 (12)	189	36.1	-	82.9
3/10	Swin-Small (25)	50	8.7	31	83.0
340	ViL-Base (56)	56	13.4	-	83.2
341	CSWin-Small (9)	35	6.9	11	83.6
342	SCHEME-12-S36	58	9.58	38	84.0

Backbone	#P (M)	F (G)	mIoU (%)
Semantic FI	PN		<u> </u>
ResNet-18(16)	16	32.2	32.9
PVT-Tiny(42)	17	33.2	35.7
ResNet-50(16)	29	45.6	36.7
PoolFormer-S12(52)	16	30.9	37.2
ResNet-101(16)	48	65.1	38.8
ResNeXt-101-32x4d(49)	47	64.7	39.7
PVT-Small(42)	28	44.5	39.8
XCiT-T12/8(12)	8.4	-	39.9
PoolFormer-S24(52)	23	39.3	40.3
SCHEMEformer-44-S12	15.5	34.3	40.9
PVT-Medium(42)	48	61.0	41.6
PoolFormer-S36(52)	35	47.5	42.0
PVT-Large(42)	65	79.6	42.1
PoolFormer-M36(52)	60	67.6	42.4
SCHEMEformer-44-S24	24.8	45.7	42.5
UperNet			
Swin-Tiny (25)	60	945	44.5
PVT-Large (43)	65	318	44.8
Focal-Tiny (50)	62	998	45.8
XCiT-S12/16 (12)	52	-	45.9
DaViT-Tiny (8)	60	940	46.3
SCHEME-DaViT-12-Tinv	68	969	47.1

Table 4: Image Classification on ImageNet-1K. Comparison Table 5: Semantic Segmentation results 343 with SOTA ViTs grouped by accuracy. Proposed SCHEME on ADE20K. FLOPs calculated at  $512 \times$ 344 models use expansion ratio 8. SCHEMEformer family has 512 resolution for Semantic FPN and 345 higher throughput and accuracy than SOTA models.

 $1024 \times 1024$  input resolution for UperNet.

346 We next compare the SCHEMEformer family against the SOTA transformers in the literature. This is 347 not an easy comparison, since models vary in size, FLOPS, and throughput. Because it is difficult to 348 make any of these variables exactly the same for two different architectures, the comparison is only 349 possible in terms of how the different architectures trade-off accuracy for any of the other factors. 350 For a given pair of variables, e.g. FLOPS vs accuracy, the model is said to be on the Pareto frontier 351 of the two variables if it achieves the best trade-off between the two. Table 4 presents a comparison of the SCHEMEformer family against various SOTA transformers in the literature. Each section 352 of the table compares a SCHEMEformer model to a group of SOTA transformers of equivalent 353 size or complexity. Note that, in each section, the remaining models have *both* lower throughput 354 and accuracy than the SCHEMEformer model. In many cases they also have more parameters 355 and FLOPs. Figure 1 provides a broader visualization of how SCHEME models establish Pareto 356 frontiers for accuracy vs FLOPS, accuracy vs throughput, and accuracy vs model size (parameters). 357 Figure 1 a) illustrates the trade-off between accuracy and FLOPS of many SOTA transformers. The 358 dashed line connects the SCHEMEformer model results, summarizing the accuracy-FLOPs trade-359 off of the family. It can be seen that the SCHEME models lie on the Pareto frontier for these 360 two objectives. This illustrates the fine control that SCHEME allows over the accuracy/complexity 361 trade-off of transformer models. Figure 1 b) presents a similar comparison for model sizes. Like for 362 FLOPS, the SCHEME models lie on the Pareto frontier for accuracy vs model size. In fact the two plots are quite similar, showing that in general there is a good correlation between model size and 363 FLOPS. 364

This is not the case for throughput, which for transformers is known to not necessarily correlate 366 with FLOPS, due to GPU parallelism. For example, Table 4 shows that, CSWin models have lower 367 throughput despite having lower FLOPs while ViT has higher throughput despite having higher 368 FLOPs. SCHEMEformer controls this trade-off by controlling the expansion ratio and block diagonal structure, which enables higher FLOPs utilization for a given throughput (14). We compute 369 throughput on a NVIDIA-Titan-X GPU with a batch size of 1 with input size 224x224 averaged 370 over 1000 runs. While comparisons could be made for larger batch sizes, we consider the setting for 371 live/streaming applications, where speed is most critical. In these applications, the concern is usually 372 inference throughput, which requires batch size of 1. Figure 1 c) illustrates the trade-off between 373 accuracy and throughput of various models, including very fast ResNet models of low accuracy. It 374 can bee seen that the SCHEMEformer again achieves the best trade-off between these two variables, 375 thus lying on the Pareto frontier for accuracy vs throughput. Its performance is particularly dominant 376 in the range of throughputs between 75 and 150 images/sec, where it significantly outperforms the 377 other methods. These results demonstrate how the SCHEME module endows transformer designers

378									
379	RetinaNet 1×	#Par	Т	AP	$AP_{50}$	$AP_{75}$	$AP_S$	$AP_M$	$AP_L$
380	PoolFormer-S12(52)	21.7	13.1	36.2	56.2	38.2	20.8	39.1	48.0
381	ResNet-50(16) SCHEME-44-e8-S12	37.7	15.7 6.1	36.3 38.3	55.3 58.0	38.6 <b>40.4</b>	19.3 21.0	40.0 <b>41.4</b>	48.8 52.3
382	PoolFormer-S24(52)	31.1	8.9	<b>38.9</b>	<b>59.7</b>	41.3	23.3	42.1	51.8
383	SCHEME-44-e8-S24	30.7	3.5	38.5	57.8 58.7	41.2	21.4	42.6	<b>53.5</b>
384	DAT-T (47)	38	-	42.8	64.4	45.2	28.0	45.8	<b>57.8</b>
385	DaViT-Tiny (8)	41 39	8.2	44.4 44.0	55.5 65.6	47.3	19.5 29.6	40.0 47.9	48.8 57.3
386	SCHEME-DaViT	47	7.8	44.7	66.2	48.3	30.0	48.8	57.2
387	Mask R-CNN 1×	#Par	Т	$AP^{o}$	$AP_{50}^{o}$	$AP_{75}^{o}$	$AP^m$	$AP_{50}^m$	$AP_{75}^{m}$
388	PoolFormer-S12(52) ResNet 50(16)	31.6	9.9 15.4	37.3	59.0 58.6	40.1	34.6 34.4	55.8	36.9 36.7
389	SCHEME-44-e8-S12	31	6.0	<b>39.8</b>	<b>61.9</b>	42.9	37.1	<b>59.2</b>	<b>39.4</b>
390	PoolFormer-S24(52) ResNet-101(16)	41.0 63.2	7.9 12.1	40.1 40.4	62.2 61.1	43.4 44.2	37.0 36.4	59.1 57.7	39.6 38.8
391	SCHEME-44-e8-S24	41	3.4	40.9	62.5	44.6	37.8	59.7	40.4
392	DAT-T (47)	48	-	44.4	67.6	48.5	40.4	64.2	43.1
392 393	DAT-T (47) CrossFormer-S (44) DaViT-Tiny (8)	48 50 48	- - 7.8	44.4 45.4 45.0	67.6 68.0 68.1	48.5 49.7 49.4	40.4 41.4 41.1	64.2 64.8 64.9	43.1 <b>44.6</b> 44.2

Mod	lel	BDM	CCA	Acc (%)									
Base	eline	1		78.9									
44-е	8-S12	$\checkmark$		79.1									
44-е	8-S12	$\checkmark$	$\checkmark$	79.7									
Table 7	7: C	ontrib	nution	1 of	BD-								
MIP	and	CCA	hra	nches	in								
CLIEN		cen	014	nenes									
SCHEN	/IE.												
Module	#Par	FLOP	T	Top	<b>b</b> -1								
	(M)	(G)	(img/	s) Acc	(%)								
Shuffle	11.83	1.77	108	79	.1								
SE	11.98	1.77	86	79	.3								
Conv	13.31	1.97	95	79	.6								
DyCCA	11.90	1.83	75	79	.6								
CCA	11.83	1.77	133	79	.7								
Table 8	B: Al	ternat	ive f	eature	at-								
tention	tention designs in SCHEME.												
CCA	G1   G	2 G3	G4	Ensemb	ole								
	(%) (%	6) (%)	(%)	(%)									
5	51.1 54	.3 55.4	54.7	73.2									
$\checkmark$	47.6 50	.0 49.2	55.1	73.8									

Table 6: COCO-17 Object Detection and Instance Segmentation. All backbones pretrained on ImageNet-1K (1x learning schedule).  $(AP^b, AP^m)$ : (bounding box AP, mask AP). T: throughput (images/sec). Table 9: **Ablation study** on the formation of feature clusters in the BD-MLP branch of the SCHEME module.

with the ability to produce multiple models at different points of the Pareto frontiers of accuracy vsmodel size, FLOPS, or throughput.

SCHEME with other ViTs: Figure 4 shows the Accuracy-Throughput curves for various SCHEME
 models obtained by replacing the MLP blocks of popular ViT architectures. It shows that SCHEME
 improves the accuracy for a fixed throughput, the throughput for a fixed accuracy, or both. The gains
 can be substantial, e.g. about 1% accuracy gain (constant throughput) for the fastest transformer
 or speed gains of up to 20% at constant accuracy. This shows that SCHEME benefits various ViT
 backbones, not just the MetaFormer.

406 Semantic Segmentation: Table 5 compares the semantic segmentation performance of two 407 SCHEMEformer models using the semantic FPN framework (18) to various SOTA models of sim-408 ilar complexity, on ADE20K. Since we do not not have access to the throughput of most models, 409 we report only parameter sizes and FLOPS. In each section of the table, the remaining models have 410 comparable or larger FLOPs and model size but lower accuracy than the corresponding SCHEMEformer. For example, SCHEMEformer-44-e8-S12 achieves 40.9% mIoU, which is 5.2/3.7 points 411 higher than the PvT-Tiny/PoolFormer-S12, which both have comparable size and FLOPs. Similarly, 412 the S24 model outperforms PoolFormer-M36 using only 41% of its parameters. To demonstrate the 413 applicability of SCHEME to larger models, we also present a comparison of the SCHEME version 414 of the DaViT-Tiny using the UperNet framework (48). While the DaViT-Tiny is already the best 415 model in the table, the use of the SCHEME mixer improves its performance by an additional 0.8 416 points. 417

**Object Detection:** Table 6 compares SCHEMEformer models to models of similar complexity 418 on the COCO-17 object detection and instance segmentation benchmark, for both Retinanet and 419 Mask-RCNN detection heads. Again, the SCHEMEformer models outperform most other models 420 of the same or smaller size. The only exception is the PoolFormer-S24, which slightly outperforms 421 (0.1 points) the comparable SCHEME-former-44-S24, for the RetinaNet head. However, with the 422 stronger Mask R-CNN head, the SCHEMEformer-44-S24 beats the PoolFormer-S24 by 0.8 points. 423 For the top performing models in the bottom third of each section of Table 6, SCHEME-DaViT-Tiny 424 improves over DaViT-Tiny by an additional 0.7% and 0.9%, for RetinaNet and MaskRCNN heads 425 respectively, while maintaining a comparable throughput. 426

427 4.2 ABLATION STUDIES

428

**Contribution of BD-MLP and CCA branch:** Table 7 shows an ablation of the contribution of the BD-MLP and CCA branches. Starting from the Metaformer-PPAA-11-e2-S12, with expansion ratio E = 2 and dense MLP, we replace the channel mixer by SCHEME to obtain the SCHEMEformer-PPAA-44-e8-S12 ( $g_1 = g_2 = 4$  and E = 8), which maintains the number of parameters and FLOPs 432 constant. The SCHEME model with only BD-MLP improves on the baseline by 0.2%. The addition
433 of the CCA branch provides an additional gain of 0.6%, showing the gains of better feature clusters.
434 Since CCA is not used at inference, its gains are *free* in terms of additional parameters/FLOPs.

**Regularizing effect of CCA:** Fig. 3 shows the evolution of the weight  $1 - \alpha$  of the CCA branch in (6) during training. While initially large, it gradually decays to zero as training progresses. This holds for all network layers. Hence, CCA can be discarded at inference. Fig. 1 in supplementary plots the weights  $1 - \alpha$  upon training convergence, for the family of SCHEMEformer-PPAA-44-e8 models, confirming that the weights are indeed very close to zero across all layers.

Effect of large expansion ratios: The left of Figure 4 shows the effect of simultaneously increasing the expansion ratio *E* and adjusting groups to realize different models of similar size and complexity. All models are based on the Metaformer-PPAA-S12. For fixed parameters/FLOPs, SCHEMEformer models achieve a gain of 1.0% to 1.4% over the baseline by increasing *E* from 1 to 4. The gains increase with larger expansions and saturate at larger FLOPs. This shows that higher internal feature dimensions are important for obtaining better accuracy with smaller ViT models.

Alternative designs of Channel Mixer: We investigate whether alternative choices to the CCA 447 branch could accomplish this goal more effectively. Table 8 compares models that replace CCA 448 with other feature mixing operations: the channel shuffling operation of ShuffleNet (57), a squeeze 449 and excitation (17) network (SENet), a single layer of convolution, and a dynamic version of CCA 450 (DyCCA), where the weight  $\alpha$  of (6) is predicted dynamically, using GCT attention (28). CCA 451 obtains the best result. While CCA is computationally heavier than some of these alternatives, it is 452 not needed at inference, as shown in Table 3. This is not true for the alternatives, which produce 453 much more balanced weights  $\alpha$  after training convergence, and cannot be discarded at inference 454 without performance drop. We conjecture that, because the alternatives have learnable parameters, 455 the network learns to use them to extract complementary features, which must be used at inference.

456 Feature Clustering: We conjectured above that CCA helps training because it facilitates feature 457 clustering into naturally independent groups that do not require inter-group communication. We 458 tested this hypothesis by studying the intermediate feature vectors  $\mathbf{v}_l$  (obtained after max-pooling 459 the features y of (4)) extracted from four randomly selected layers l of two ImageNet pretrained 460 models, trained with and without CCA. We split the features into the 4 groups used in the model  $\mathbf{v}_{l,q}, g \in \{1, \dots, 4\}$  and concatenated the features of all layers in the same group. This produced 461 four vectors  $\mathbf{u}_g = \operatorname{concat}(\{v_{l,g}\}_l)$  containing the features of each group g extracted throughout 462 the network. A linear classifier was then learned over each vector  $\mathbf{u}_q$ . Table 9 shows the top-1 463 accuracy per feature group and model. To evaluate whether groups learn different class clusters, 464 we also average the outputs from the four group classifiers to obtain the final accuracy. Without 465 CCA, i.e. no group communication during training, the network produces feature groups individu-466 ally more predictive of the image class, but less predictive when combined. This suggests that there 467 is redundancy between the features of the different groups. By introducing inter-group communi-468 cation, CCA enables the groups to learn more diverse features, that complement each other. Fig. 3 469 shows the class separability of the intermediate features  $\mathbf{y}$  of (4) of a randomly chosen layer of the 470 SCHEMEformer-PPAA-44-e8-S12. Class separability was measured as in (30), with a final value 471 obtained by averaging the class separability across all classes. The model trained with CCA has higher class separability than that without it. This confirms that CCA is helpful in forming feature 472 clusters that increase class separability during training. Conversely, we tested if CCA is helpful 473 when channel shuffling is inserted in between the two mixer MLP layers, which destroys the group 474 structure. This variant of the SCHEMEformer-PPAA-44-e8-S12 model achieved an accuracy of 475 79.1% for both training with and without CCA (Table 7, row 2). This shows that CCA is not help-476 ful when feature groups are mixed. Similarly, CCA did not provide any gains when applied to the 477 standard MLP branch with full feature mixing. These results suggest that CCA indeed helps to form 478 the independent feature clusters needed to achieve the computational efficiency of channel groups 479 without performance degradation. Table 9 shows that, despite the higher accuracy of the individual 480 group features of the model trained without CCA, the model trained with CCA has 0.6% higher 481 ensemble accuracy.

482

483 484

# <sup>486</sup> 5 LIMITATION AND CONCLUSION

487 488

505 506

507 508

509

510 511

512

513 514

515

516

517

518

519

520 521

522

523 524

525

526

527

528

529

530 531

532

533

534 535

536

538

Although SCHEME module improves the accuracy-throughput curves of popular ViTs, it incurs a slight overhead in memory (see appendix section A.3.2) during training due to the channel attention operation which limits the batch size on smaller GPUs. However, Figure 3 shows that the contribution of the CCA branch becomes negligible after 150 epochs. Therefore, the CCA branch can be removed beyond this point, significantly reducing the training overhead.

In this work, we proposed the SCHEME module for improving the performance of ViTs. SCHEME 493 494 leverages a block diagonal feature mixing structure to enable MLPs with larger expansion ratios, a property that is shown to improve transformer performance, without increase of model parameters 495 or computation. It uses a weighted fusion of a BD-MLP branch, which abstracts existing MLPs with 496 block diagonal structure, and a parameter-free CCA branch that helps to cluster features into groups 497 during training. The CCA branch was shown to improve training but not be needed at inference. 498 Experiments showed that it indeed improves the class separability of the internal feature representa-499 tion of the BD-MLP branch, helping create feature clusters that are informative of the image class. 500 The standard transformer MLP was replaced with the SCHEME module to obtain a new family of 501 SCHEMEformer models with improved performance for classification, detection, and segmentation, 502 for fixed parameters and FLOPs with favorable latency. SCHEME was shown to be effective for var-503 ious ViT architectures and to provide a flexible way to scale models, always outperforming models 504 with smaller MLP expansion ratios having the same complexity.

### References

- [1] Chen, B., Dao, T., Liang, K., Yang, J., Song, Z., Rudra, A., Ré, C.: Pixelated butterfly: Simple and efficient sparse training for neural network models. International Conference on Learning Representations (2021)
- [2] Chen, Y., Dai, X., Chen, D., Liu, M., Dong, X., Yuan, L., Liu, Z.: Mobile-former: Bridging mobilenet and transformer. ECCV (2021)
- [3] Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1800–1807 (2017). https://doi.org/10.1109/CVPR.2017.195
- [4] Chu, X., Tian, Z., Zhang, B., Wang, X., Shen, C.: Conditional positional encodings for vision transformers. In: The Eleventh International Conference on Learning Representations (2023), https://openreview.net/forum?id=3KWnuT-R1bh
- [5] Dai, Z., Liu, H., Le, Q.V., Tan, M.: Coatnet: Marrying convolution and attention for all data sizes. arXiv preprint arXiv:2106.04803 (2021)
- [6] Dao, T., Chen, B., Sohoni, N.S., Desai, A., Poli, M., Grogan, J., Liu, A., Rao, A., Rudra, A., Ré, C.: Monarch: Expressive structured matrices for efficient and accurate training. In: International Conference on Machine Learning. PMLR (2022)
- [7] Darcet, T., Oquab, M., Mairal, J., Bojanowski, P.: Vision transformers need registers. In: The Twelfth International Conference on Learning Representations (2024), https: //openreview.net/forum?id=2dn03LLiJ1
- [8] Ding, M., Xiao, B., Codella, N., Luo, P., Wang, J., Yuan, L.: Davit: Dual attention vision transformers. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIV. pp. 74–92. Springer (2022)
- [9] Dong, X., Bao, J., Chen, D., Zhang, W., Yu, N., Yuan, L., Chen, D., Guo, B.: Cswin transformer: A general vision transformer backbone with cross-shaped windows (2021)
- [10] Dong, Y., Cordonnier, J.B., Loukas, A.: Attention is not all you need, pure attention loses rank doubly exponentially with depth. ICML (2021), https://arxiv.org/abs/2103. 03404

546

547

548

549

550 551

552

553

554

555

556

557 558

559

560

561

562 563

564

565

566

567 568

569 570

571

572 573

574

575

576

577

578 579

580

581

582

583

584 585

586

588

589

590

592

- [11] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021), https://openreview.net/forum?id=YicbFdNTTy
  - [12] El-Nouby, A., Touvron, H., Caron, M., Bojanowski, P., Douze, M., Joulin, A., Laptev, I., Neverova, N., Synnaeve, G., Verbeek, J., et al.: Xcit: Cross-covariance image transformers. arXiv preprint arXiv:2106.09681 (2021)
  - [13] Fedus, W., Zoph, B., Shazeer, N.: Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. Journal of Machine Learning Research (2022)
    - [14] Fu, D.Y., Arora, S., Grogan, J., Johnson, I., Eyuboglu, S., Thomas, A.W., Spector, B., Poli, M., Rudra, A., Ré, C.: Monarch mixer: A simple sub-quadratic gemm-based architecture. In: Advances in Neural Information Processing Systems (2023)
  - [15] Graham, B., El-Nouby, A., Touvron, H., Stock, P., Joulin, A., Jegou, H., Douze, M.: Levit: A vision transformer in convnet's clothing for faster inference. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 12259–12269 (October 2021)
  - [16] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 770–778 (2015)
  - [17] Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
  - [18] Kirillov, A., Girshick, R., He, K., Dollar, P.: Panoptic feature pyramid networks. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6392–6401. IEEE Computer Society, Los Alamitos, CA, USA (jun 2019). https://doi.org/10.1109/CVPR.2019.00656, https://doi.ieeecomputersociety. org/10.1109/CVPR.2019.00656
  - [19] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C., Bottou, L., Weinberger, K. (eds.) Advances in Neural Information Processing Systems. vol. 25. Curran Associates, Inc. (2012), https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf
  - [20] Li, S., Wang, Z., Liu, Z., Tan, C., Lin, H., Wu, D., Chen, Z., Zheng, J., Li, S.Z.: Moganet: Multi-order gated aggregation network. In: International Conference on Learning Representations (2024)
  - [21] Li, Y., Yuan, G., Wen, Y., Hu, J., Evangelidis, G., Tulyakov, S., Wang, Y., Ren, J.: Efficientformer: Vision transformers at mobilenet speed. Neurips (2022)
  - [22] Li, Y., Zhang, K., Cao, J., Timofte, R., Van Gool, L.: Localvit: Bringing locality to vision transformers. arXiv preprint arXiv:2104.05707 (2021)
  - [23] Liu, H., Dai, Z., So, D.R., Le, Q.V.: Pay attention to mlps (2021). https://doi.org/10.48550/ARXIV.2105.08050, https://arxiv.org/abs/2105.08050
  - [24] Liu, Y., Zhang, Y., Wang, Y., Hou, F., Yuan, J., Tian, J., Zhang, Y., Shi, Z., Fan, J., He, Z.: A survey of visual transformers (2021). https://doi.org/10.48550/ARXIV.2111.06091, https: //arxiv.org/abs/2111.06091
  - [25] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021)
  - [26] Parmar, N.J., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., Tran, D.: Image transformer. In: International Conference on Machine Learning (ICML) (2018), http:// proceedings.mlr.press/v80/parmar18a.html

601

602 603

604

605

606

607

608

609 610

611

612

613

614

615

616

617

618

619

620

621 622

623 624

625

626 627

628

629 630

631

632

633

634 635

636

637 638

639 640

641

642

643

644

645

646

- [27] Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., Shlens, J.: Standalone self-attention in vision models. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 32. Curran Associates, Inc. (2019), https://proceedings.neurips.cc/ paper/2019/file/3416a75f4cea9109507cacd8e2f2aefc-Paper.pdf
  - [28] Ruan, D., Wang, D., Zheng, Y., Zheng, N., Zheng, M.: Gaussian context transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 15129–15138 (June 2021)
  - [29] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. CVPR (2018)
  - [30] Schilling, A., Maier, A., Gerum, R., Metzner, C., Krauss, P.: Quantifying the separability of data classes in neural networks. Neural Networks 139, 278– 293 (2021). https://doi.org/https://doi.org/10.1016/j.neunet.2021.03.035, https://www. sciencedirect.com/science/article/pii/S0893608021001234
  - [31] Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Gradcam: Visual explanations from deep networks via gradient-based localization. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 618–626 (2017). https://doi.org/10.1109/ICCV.2017.74
  - [32] Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., Dean, J.: Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In: International Conference on Learning Representations (2017), https://openreview.net/forum? id=BlckMDqlg
  - [33] Srinivas, A., Lin, T.Y., Parmar, N., Shlens, J., Abbeel, P., Vaswani, A.: Bottleneck transformers for visual recognition. In: CVPR (2021). https://doi.org/10.48550/ARXIV.2101.11605, https://arxiv.org/abs/2101.11605
  - [34] Tang, S., Zhang, J., Zhu, S., Tan, P.: Quadtree attention for vision transformers. ICLR (2022)
  - [35] Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H.: Training dataefficient image transformers and distillation through attention. In: International Conference on Machine Learning. vol. 139, pp. 10347–10357 (July 2021)
  - [36] Touvron, H., Cord, M., El-Nouby, A., Verbeek, J., Jegou, H.: Three things everyone should know about vision transformers. ECCV (2022)
  - [37] Touvron, H., Cord, M., Jegou, H.: Deit iii: Revenge of the vit. ECCV (2022)
  - [38] Touvron, H., Cord, M., Sablayrolles, A., Synnaeve, G., Jégou, H.: Going deeper with image transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 32–42 (October 2021)
  - [39] Tu, Z., Talebi, H., Zhang, H., Yang, F., Milanfar, P., Bovik, A., Li, Y.: Maxvit: Multi-axis vision transformer. ECCV (2022)
  - [40] Vaswani, A., Ramachandran, P., Srinivas, A., Parmar, N., Hechtman, B., Shlens, J.: Scaling local self-attention for parameter efficient visual backbones. In: CVPR (2021)
  - [41] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 30. Curran Associates, Inc. (2017), https://proceedings.neurips. cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
  - [42] Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 568–578 (2021)

655

656

657

658

659

660

661

662 663

664

665

666

667

668 669

670

671

672 673

674 675

676

677 678

679

680

681

682

683

684 685

686

687

688

689

690 691

692

693

694

695

696

697

698 699

700

- [43] Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: PVT v2: Improved baselines with pyramid vision transformer. Computational Visual Media (2022)
- [44] Wang, W., Yao, L., Chen, L., Lin, B., Cai, D., He, X., Liu, W.: Crossformer: A versatile vision transformer hinging on cross-scale attention. In: International Conference on Learning Representations, ICLR (2022), https://openreview.net/forum?id=\_PHymLIxuI
  - [45] Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., Zhang, L.: Cvt: Introducing convolutions to vision transformers. IEEE/CVF International Conference on Computer Vision (ICCV) (2021)
  - [46] Wu, K., Peng, H., Chen, M., Fu, J., Chao, H.: Rethinking and improving relative position encoding for vision transformer. In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 10013–10021. IEEE Computer Society, Los Alamitos, CA, USA (oct 2021). https://doi.org/10.1109/ICCV48922.2021.00988, https://doi. ieeecomputersociety.org/10.1109/ICCV48922.2021.00988
  - [47] Xia, Z., Pan, X., Song, S., Li, L.E., Huang, G.: Vision transformer with deformable attention. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4794–4803 (June 2022)
  - [48] Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: European Conference on Computer Vision. Springer (2018)
  - [49] Xie, S., Girshick, R., Dollar, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5987–5995 (07 2017). https://doi.org/10.1109/CVPR.2017.634
  - [50] Yang, J., Li, C., Zhang, P., Dai, X., Xiao, B., Yuan, L., Gao, J.: Focal self-attention for localglobal interactions in vision transformers (2021)
  - [51] Yu, T., Li, X., Cai, Y., Sun, M., Li, P.: S<sup>2</sup>-mlp: Spatial-shift mlp architecture for vision. WACV (2021)
  - [52] Yu, W., Luo, M., Zhou, P., Si, C., Zhou, Y., Wang, X., Feng, J., Yan, S.: Metaformer is actually what you need for vision. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10819–10829 (2022)
  - [53] Yuan, K., Guo, S., Liu, Z., Zhou, A., Yu, F., Wu, W.: Incorporating convolution designs into visual transformers (2021). https://doi.org/10.48550/ARXIV.2103.11816, https: //arxiv.org/abs/2103.11816
  - [54] Yuan, L., Chen, Y., Wang, T., Yu, W., Shi, Y., Jiang, Z.H., Tay, F.E., Feng, J., Yan, S.: Tokensto-token vit: Training vision transformers from scratch on imagenet. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 558–567 (October 2021)
  - [55] Zhang, P., Dai, X., Yang, J., Xiao, B., Yuan, L., Zhang, L., Gao, J.: Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. ICCV 2021 (2021)
  - [56] Zhang, P., Dai, X., Yang, J., Xiao, B., Yuan, L., Zhang, L., Gao, J.: Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 2998–3008 (October 2021)
  - [57] Zhang, X., Zhou, X., Lin, M., Sun, J.: Shufflenet: An extremely efficient convolutional neural network for mobile devices. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6848–6856 (06 2018). https://doi.org/10.1109/CVPR.2018.00716
  - [58] Zhu, L., Wang, X., Ke, Z., Zhang, W., Lau, R.: Biformer: Vision transformer with bi-level routing attention. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
    - 13

702	Α ΑΡΡΕΝΟΙΧ
703	
704	A.1 CODE RELEASE
705	
707	Code and trained models will be released upon acceptance of the paper.
708	
709	A.2 IMPLEMENTATION DETAILS
710	A.2.1 Hyperparameter Settings
711	Table 11 shows the detailed hyperparameter settings of the family of SCHEMEformer models re-
712	ported in the main paper.
714	
715	A.3 ABLATION STUDIES
716	A 2.1 DE OT OF THE LEADNED 1 - A MERCHITE IN THE SCHEME MODILE FOR
717	A.5.1 PLOT OF THE LEARNED 1 – $\alpha$ weights in the SCHEME MODULE FOR SCHEMEFORMER MODELS
718	beneficiel okalik mobiles
719	The learned weights $1 - \alpha$ for the SCHEME former model family are shown in Fig. 8. Interestingly,
721	the learned weights coarsely approximate the snape of a gaussian distribution. The learned weights reach a peak value in the middle layers of the network and drop to zero for all the other layers. The
722	middle of the network typically corresponds to the initial few layers of the third stage of the model
723	that contains the maximum number of transformer blocks for all the models shown in Fig. 8. We
724	conjecture that the weights for these layers have not fully converged and that feature mixing can still
725	be useful for these layers and so training for more epochs will allow the $1 - \alpha$ weights of these layers to converge to zero. To test this hypothesis, we trained the SCHEMEformer DPAA 12 e8 S12 model
726	for an additional 200 epochs beyond the standard 300 epochs and observed that the peak value of
727	the $1 - \alpha$ weights decreased further by 0.11 as compared to the model trained for 300 epochs and
720	the accuracy improved by <b>0.4</b> %. The weight norm (of all layers) decreased from 0.18 for 300 epoch
730	model to 0.05 for 500 epoch model. This confirms our hypothesis and training SCHEMEformer
731	models for larger epochs can further improve the accuracy.
732	A 3.2 TRAINING OVERHEAD OF CCA
733	
734	Table 12 compares the training GPU memory and throughput for SCHEME mixer with and without
735	using UCA. UCA improves the accuracy by $0.6\%$ with only a slight increase in the GPU mem- ory ( $\pm 12.5\%$ ) and training time ( $\pm 16.7\%$ ). Further CCA is not needed during inference thereby
730	providing gains for "free" without additional computational cost at inference.
101	

739 A.3.3 TRAINING FOR LONGER EPOCHS

Table 13 shows the comparison of training for longer epochs for SCHEMEformer with the base-line Metaformer model. We train for an additional 200 epochs from the standard 300 epochs.
SCHEMEformer-PPAA-44-e8-S12 trained for 300 epochs even outperforms the baseline model trained for 500 epochs. On continuing the training from 300 to 500 epochs, SCHEMEformer continues to improve the performance without saturation suggesting that it is beneficial to train with CCA for longer epochs. 300 to 500 epochs is a much larger increase of training time (67%) than the 16.7% increase in training time required by SCHEME mixer (see Table 12).

- 747
- 748 A.3.4 IMPACT OF REMOVING CCA AT INFERENCE

Table 14 shows the impact of removing CCA at inference, for various backbones. While the number of FLOPS decreases, the top-1 accuracy changes very little ( $\approx 0.02$  difference). Hence, there is no advantage in using CCA at inference. This is unlike training, where the use of CCA makes a non-negligible difference, as shown in Table 6 of the main paper.

Table 15 and 16 show the results of removing CCA at inference for object detection and semantic
 segmentation models, respectively. For both tasks, the results are identical to the model using CCA showing that the CCA also generalizes to downstream tasks.

# 756 A.3.5 CCA

In the main paper, feature groups across different layers of the model was used to demonstrate the learning of feature clusters by CCA. Here, the feature from the final layer of the model is only considered. Table 18 shows the top-1 accuracy per feature group and model. The average of the outputs from the four group classifiers is reported in the final column of the table. The effect is more pronounced when using a single layer feature with +1.09% accuracy difference between the model with and without CCA. This further reinforces that by introducing inter-group communication, CCA enables the groups to learn more diverse sets of features, that complement each other.

A.3.6 LARGER EXPANSION RATIOS

Table 17 shows additional results of using larger expansion ratios with SCHEME mixer (illustrated in Fig. 3 of the main paper) for the same number of parameters and FLOPs using the MetaformerPPAA-S12 baseline model. We observe that SCHEMEformer-PPAA consistently outperforms the baseline for larger expansions ratios with larger gains at lower FLOPs. The performance saturates as the model size and FLOPs increases.

- 773 A.4 QUALITATIVE ANALYSIS

# 774 A.4.1 GRAD-CAM (31) VISUALIZATION

Fig. 9 shows the results of class activation maps for SCHEMEformer-PPAA-44-e8-S12 model for
a few examples from the validation set of ImageNet-1K dataset. The stronger heatmap responses
around the salient features of an object (e.g., body of a bird, cat) shows that the model ignores the
background and attends to more discriminative spatial regions. Fig. 9 also shows the qualitative
comparison with a few existing methods such as ResNet-50, DeiT-S, Poolformer etc. of similar
complexity. It demonstrates that SCHEMEformer-PPAA-44-e8-S12 attends to the complete object
class and less spurious features showing that it is better than the competing methods.

### 784 A.4.2 Attention Visualization

Figure 10 shows the results of attention maps with and without CCA on Imagenet dataset. We find that there are differences in attention between using or not using CCA. CCA increases the tendency of the model to attend to the regions where the class is present, which is denoted by green ellipses.



Figure 5: [Zoomed version] Comparison of the proposed SCHEMEformer family, derived from the Metaformer-PPAA-S12 model (52) with higher expansion ratios in the MLP blocks, and many SOTA transformers from the literature. The SCHEMEFormer family establishes a new Pareto frontier (optimal trade-off)
for a) accuracy vs. FLOPs, b) accuracy vs model size, and c) accuracy vs, throughput. SCHEMEformer models are particularly effective for the design of fast transformers (throughput between 75 and 150 images/s) with small model size.



Figure 6: [Zoomed version] Impact of CCA Figure 7: (SCHEMEformer-44-e8-S12). Top: Evolution of Top: Accuracy vs FLOPs of various SCHEME arability of output features y (over 50 random classes tom: SCHEME mixer improves accuracy for fixed of ImageNet-1K validation set) for model trained throughput or vice-versa for various popular ViT arwith and without CCA.

[Zoomed version] SCHEME tradeoffs. weight  $1-\alpha$  across model layers. **Bottom:** Class sep- models with different MLP configurations. **Bot**chitectures.

Model	#Par (M)	FLOPs (G)	T (img/s)	Acc (%)
ViT-Base (11)	86	17.6	112	79.7
SCHEME-ViT-12-e8-Base	77.1	15.5	130	79.9
DeiT-Tiny (35)	6	1.3	117	74.5
SCHEME-DeiT-12-e8-Tiny	7.5	1.6	117	76.0
Poolformer-S12	12.0	1.8	166	77.2
SCHEME-Poolformer-12-e8-S12	16.7	2.6	171	78.5
Metaformer-S12	12.0	1.8	133	78.9
SCHEMEformer-PPAA-12-e8-S12	11.8	1.8	133	79.7
CAformer-S12	25	4.20	74	82.9
SCHEME-CAformer-12-e8-S12	23.9	3.60	80	82.9
CoAtNet-0	25.0	4.2	88	81.6
SCHEME-CoAtNet-12-e8-0	24.0	4.1	110	81.6
CSWin-Tiny	23.0	4.3	20	82.8
SCHEME-CSWin-12-e8-Tiny	29.1	5.6	21	83.2
DaViT-Tiny	23.0	4.3	61	82.8
SCHEME-DaViT-12-e8-Tiny	37.0	6.6	62	83.0
T2T-ViT-14	21.5	6.1	70	81.7
SCHEME-T2T-ViT-12-e8-14	27.7	8.1	71	82.1
BiFormer-Tiny(58)	13	2.2	57	81.4
SCHEME-BiFormer-12-e3-Tiny	11.5	1.9	59	81.4

Table 10: Comparison with state-of-the-art ViT models on the ImageNet-1K dataset. The SCHEME module enhances the accuracy of existing ViTs while maintaining or achieving higher throughput.

#### 920 SCHEMEformer-PPAA-44-e8 SCHEMEformer-PPAA-12-e8 SCHEME-CAformer Model 12-e8-S12 S12 S24 \$36 S12 S24 \$36 44-e8-S18 921 Peak drop rate of stoch. depth $d_r$ 0.1 0.2 0.4 0.1 0.2 0.4 0.15 0.15 $10^{-5}$ $10^{-5}$ $10^{-6}$ $10^{-5}$ $10^{-5}$ $10^{-6}$ $10^{-5}$ 922 $10^{-5}$ LayerScale initialization $\epsilon$ Data augmentation AutoAugmen 923 Repeated Augmentation off 924 Input resolution 224 Epochs 300 925 Hidden dropout 0 926 GELU dropout 0 Classification dropout 0 927 Random erasing prob 0.25 EMA decay 0 928 Cutmix $\alpha$ 1.0 0.8 929 Mixup $\alpha$ Cutmix-Mixup switch prob 0.5 930 Label smoothing 0.1 Batch size used in the paper 1024 931 Learning rate decay Weight decay cosine 932 0.05 Gradient clipping None 933 Warmup epochs 20 Relation between peak learning rate and batch size $lr = \frac{batch size}{1024} \times$ $\frac{\text{batch size}}{1024} \times 8 \times e$ lr 934 AdamW Optimizer LAMB 935 Adam $\epsilon$ $1e^{i}$ None (0.9, 0.999) $\mathrm{Adam}\,(\beta_1,\beta_2)$ None 936 937 938 939 940 941 942 943 944 (d) SCHEME-CoAtNet-44-e8-945 (a) SCHEMEformer-44-e8-S12 (b) SCHEMEformer-44-e8-S24 (c) SCHEMEformer-44-e8-S36 Tiny 946 947 948 949 0.20 950 951 952 (e) SCHEMEformer-12-e8-S12 (f) SCHEMEformer-12-e8-S24 (g) SCHEMEformer-12-e8-S36 (h) SCHEME-Swin-12-e8-Tiny 953

Table 11: Hyperparameter Settings for the family of SCHEMEformer models trained on ImageNet-1K dataset.

Figure 8: Plot of the learned weight  $(1-\alpha)$  values across different layers of a network for the family of SCHEMEformer models. The weights reach a peak value near the middle of the network. We demonstrate that these peak weights are not yet converged and training the network for more epochs decays these weights to zero while also improving the accuracy. For example, training the SCHEMEformer-PPAA-12-e8-S12 model for 200 additional epochs reduced the weight norm of the vector of  $1 - \alpha$  weights from 0.18 to 0.05 showing that these weights gradually approach zero as the training progresses while improving the accuracy further by 0.4%.

960 961 962

963

970 971

954 955

956

957

958

959

Table 12: Training Overhead of CCA. CCA adds only a small overhead in GPU memory and training time.

	CCA	#Par	Train FLOPs	Val Acc	GPU Mem.	Train Throughput
Model		(M)↓	(G)↓	(%)↑	$(G)\downarrow$	(iters/s) ↑
SCHEMEformer-PPAA-44-e8-S12		11.8	1.77	79.1	8	215
SCHEMEformer-PPAA-44-e8-S12	$\checkmark$	11.8	2.16	79.7	9	180

Table 13: Ablation study	of longer	r training	for SCI	IEMEfo	rmer-	PPAA	<b>\</b> -44-	e8-S12.
	1000 0 0				000 1	( 24 )	<b>#</b> 00.1	( 444 )

Model	#P (M)	FLOPs (G)	Throughput (img/s)	300-Acc (%)	500-Acc (%)
Metaformer-11-e2-S12 (Baseline)	11.8	1.77	133	78.9	79.6
SCHEMEformer-44-e8-S12	11.8	1.77	133	79.7	80.1

Table 14: Impact of removing CCA branch during inference.

	CCA	#Params	FLOPs	Top-1 Acc
Model	Used	(M)	(G)	(%)
SCHEMEformer-PPAA-44-e8-S12	$\checkmark$	11.83	2.16	79.74
SCHEMEformer-PPAA-44-e8-S12		11.83	1.77	79.72
SCHEMEformer-PPAA-12-e8-S24	$\checkmark$	40.0	7.3	82.80
SCHEMEformer-PPAA-12-e8-S24		40.0	6.5	82.76
SCHEMEformer-PPAA-12-e8-S36	$\checkmark$	58.8	10.8	84.00
SCHEMEformer-PPAA-12-e8-S36		58.8	9.6	83.95
SCHEME-CoatNet-44-e8-0	$\checkmark$	17.80	3.83	80.70
SCHEME-CoatNet-44-e8-0		17.80	3.42	80.68
SCHEME-Swin-12-e8-T	$\checkmark$	36.93	7.00	81.69
SCHEME-Swin-12-e8-T		36.93	5.89	81.69

 
 Table 15: Ablation study of removing CCA for COCO-17 Object Detection and Table
 16:
 Instance Segmentation. Removing CCA at inference does not impact the AP val- study of removing CCA ues as they are identical to the model using CCA at inference.  $AP^{b}$  and  $AP^{m}$  for Semantic Segmentadenote bounding box AP and mask AP, respectively. Backbone models denote tion results on ADE20K SCHEMEformer-PPAA-44-e8- variants.

Ablation validation dataset.

Backbo	ne CCA	A d #P	AP	$AP_{50}^b$	RetinaN $AP_{75}^b$	let $1 \times AP_S^b$	$AP_M$	$AP_L$	<b>#P</b> A	$AP^b$	$AP_{50}^b$	lask R- $O$ $AP_{75}^b$	CNN 1> AP <sup>m</sup>	$AP_{50}^m$	$AP_{75}^m$	CCA	A #Par	Semanti FLOPs	c FPN mIoU (%)
S12	$\checkmark$	21	38.3	58.0	40.4	21.0	41.4	52.3	31 3	9.8	61.9	42.9	24.1	53.0	42.3	$\checkmark$	15.5	36.4	40.9
S12		21	38.3	58.0	40.4	21.0	41.4	52.3	31 3	9.8	61.9	42.9	24.1	53.0	42.3		15.5	34.3	40.9
S24	$\checkmark$	31	38.8	58.7	41.2	22.5	41.5	53.5	41 4	0.9	62.5	44.6	24.6	55.6	43.8	$\checkmark$	24.8	49.8	42.5
S24		31	38.8	58.7	41.2	22.5	41.5	53.5	41 4	0.9	62.5	44.6	24.6	55.6	43.8		24.8	45.7	42.5

Table 17: Ablation study on the effect of larger expansion ratios in BD-MLP block of SCHEME on ImageNet-1K validation dataset. 

	#Par	FLOPs	Top-1
Model	(M)	(G)	Acc (%)
Metaformer-PPAA-11-e1-S12	9.6	1.37	76.0
SCHEMEformer-PPAA-22-e2-S12	9.6	1.37	77.0
SCHEMEformer-PPAA-44-e4-S12	9.6	1.37	77.4
SCHEMEformer-PPAA-66-e6-S12	10.0	1.45	77.9
Metaformer-PPAA-11-e2-S12	11.8	1.77	78.9
SCHEMEformer-PPAA-22-e4-S12	11.8	1.77	79.4
SCHEMEformer-PPAA-44-e8-S12	11.8	1.77	79.7
SCHEMEformer-PPAA-33-e6-S12	12.5	1.87	79.8
SCHEMEformer-PPAA-22-e6-S12	14.1	2.16	80.4
Metaformer-PPAA-11-e4-S12	16.5	2.56	81.0
SCHEMEformer-PPAA-22-e8-S12	16.5	2.56	81.1
SCHEMEformer-PPAA-44-e16-S12	16.5	2.56	81.2

Table 18: Ablation study on the formation of feature clusters in the BD-MLP branch of the SCHEME mod-ule. We train a linear classifier on top of the four feature groups extracted from the final MLP mixer of the transformer block of the network. The model trained with CCA forms feature clusters that learn diverse and complementary set of features that can obtain 1% higher validation accuracy than the model trained without CCA. 

	CCA	Group1	Group2	Group3	Group4	Ensemble
Model		(%)	(%)	(%)	(%)	(%)
SCHEMEformer-44-e8-S12		29.63	37.59	37.57	35.89	73.95
SCHEMEformer-44-e8-S12	$\checkmark$	27.89	37.64	30.02	37.67	75.04



Figure 9: GRAD-CAM (31) visualization for a few validation samples on ImageNet-1K dataset for SCHEMEformer-PPAA-44-e8-S12 model and comparison with other competing methods.



by using registers (7).