

A Detailed notation

Products: The *Hadamard* product of $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{I \times N}$ is defined as $\mathbf{A} * \mathbf{B}$ and is equal to $a_{(i,j)} b_{(i,j)}$ for the (i, j) element. The *Khatri-Rao* product of matrices $\mathbf{A} \in \mathbb{R}^{I \times N}$ and $\mathbf{B} \in \mathbb{R}^{J \times N}$ is denoted by $\mathbf{A} \odot \mathbf{B}$ and yields a matrix of dimensions $(IJ) \times N$. The Khatri-Rao product for a set of matrices $\{\mathbf{A}_{[m]} \in \mathbb{R}^{I_m \times N}\}_{m=1}^M$ is abbreviated by $\mathbf{A}_{[1]} \odot \mathbf{A}_{[2]} \odot \dots \odot \mathbf{A}_{[M]} \doteq \bigodot_{m=1}^M \mathbf{A}_{[m]}$.

Tensors: Each element of an M^{th} order tensor \mathcal{X} is addressed by M indices, i.e., $(\mathcal{X})_{i_1, i_2, \dots, i_M} \doteq x_{i_1, i_2, \dots, i_M}$. An M^{th} -order tensor \mathcal{X} is defined over the tensor space $\mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}$, where $I_m \in \mathbb{Z}$ for $m = 1, 2, \dots, M$. The *mode- m unfolding* of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}$ maps \mathcal{X} to a matrix $\mathbf{X}_{(m)} \in \mathbb{R}^{I_m \times \bar{I}_m}$ with $\bar{I}_m = \prod_{k \neq m}^M I_k$ such that the tensor element x_{i_1, i_2, \dots, i_M} is mapped to the matrix element $x_{i_m, j}$ where $j = 1 + \sum_{k \neq m}^M (i_k - 1) J_k$ with $J_k = \prod_{n \neq m}^{k-1} I_n$. The *mode- m vector product* of \mathcal{X} with a vector $\mathbf{c} \in \mathbb{R}^{I_m}$, denoted by $\mathcal{X} \times_m \mathbf{c} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_{m-1} \times I_{m+1} \times \dots \times I_M}$, results in a tensor of order $M - 1$:

$$(\mathcal{X} \times_m \mathbf{c})_{i_1, \dots, i_{m-1}, i_{m+1}, \dots, i_M} = \sum_{i_m=1}^{I_m} x_{i_1, i_2, \dots, i_M} u_{i_m}. \quad (13)$$

The *CP decomposition* [32] factorizes a tensor into a sum of component rank-one tensors. The rank- R CP decomposition of an M^{th} -order tensor \mathcal{X} is written as:

$$\mathcal{X} \doteq \llbracket \mathbf{C}_{[1]}, \mathbf{C}_{[2]}, \dots, \mathbf{C}_{[M]} \rrbracket = \sum_{r=1}^R \mathbf{c}_r^{(1)} \circ \mathbf{c}_r^{(2)} \circ \dots \circ \mathbf{c}_r^{(M)}, \quad (14)$$

where \circ is the vector outer product. The factor matrices $\{\mathbf{C}_{[m]} = [\mathbf{c}_1^{(m)}, \mathbf{c}_2^{(m)}, \dots, \mathbf{c}_R^{(m)}] \in \mathbb{R}^{I_m \times R}\}_{m=1}^M$ collect the vectors from the rank-one components. By considering the mode-1 unfolding of \mathcal{X} , the CP decomposition can be written in matrix form as:

$$\mathbf{X}_{(1)} \doteq \mathbf{C}_{[1]} \left(\bigodot_{m=2}^M \mathbf{C}_{[m]} \right)^T \quad (15)$$

The following lemma is useful in our method:

Lemma 1 ([9]). *For a set of N matrices $\{\mathbf{A}_{[\nu]} \in \mathbb{R}^{I_\nu \times K}\}_{\nu=1}^N$ and $\{\mathbf{B}_{[\nu]} \in \mathbb{R}^{I_\nu \times L}\}_{\nu=1}^N$, the following equality holds:*

$$\left(\bigodot_{\nu=1}^N \mathbf{A}_{[\nu]} \right)^T \cdot \left(\bigodot_{\nu=1}^N \mathbf{B}_{[\nu]} \right) = (\mathbf{A}_{[1]}^T \cdot \mathbf{B}_{[1]}) * \dots * (\mathbf{A}_{[N]}^T \cdot \mathbf{B}_{[N]}) \quad (16)$$

B Polynomials as a single tensor product

As mentioned in the main paper polynomials and tensors are closely related. To illustrate the differences between the proposed variant of sec. 4.1 and the proposed taxonomy we can formulate them as a single tensor product. We assume a second-degree polynomial expansion of (1). The tensors are then up to third-order, which enables a visualization (as in the Fig.1). The initial equation is:

$$\mathbf{y} = \boldsymbol{\beta} + \left(\mathbf{W}^{[1]}\right)^T \mathbf{z} + \left(\mathcal{W}^{[2]} \times_2 \mathbf{z} \times_3 \mathbf{z}\right) \quad (17)$$

The τ^{th} output of (17) can be written in element-wise form as:

$$y_\tau = \beta_\tau + \sum_{k=1}^{\delta} w_{\tau,k}^{[1]} z_k + \sum_{k,m=1}^{\delta} w_{\tau,k,m}^{[2]} z_k z_m \quad (18)$$

We can collect all the parameters of (17) under a single tensor by padding the input $\mathbf{z} \in \mathbb{R}^\delta$. Specifically, if we consider the padded version $\tilde{\mathbf{z}} = [z_1, \dots, z_\delta, 1]^T$, then (17) can be written in the format $\mathbf{y} = \tilde{\mathcal{W}} \times_2 \tilde{\mathbf{z}} \times_3 \tilde{\mathbf{z}}$ as we demonstrate below.

The τ^{th} output of $\tilde{\mathcal{W}} \times_2 \tilde{\mathbf{z}} \times_3 \tilde{\mathbf{z}}$ is:

$$\begin{aligned} y_\tau = & \sum_{k,m=1}^{\delta+1} \tilde{w}_{\tau,k,m} \tilde{z}_k \tilde{z}_m = \underbrace{\tilde{w}_{\tau,\delta+1,\delta+1}}_{\text{constant term}} + \\ & \underbrace{\sum_{m=1}^{\delta} \tilde{w}_{\tau,\delta+1,m} z_m}_{\text{first-degree term}} + \underbrace{\sum_{k=1}^{\delta} \tilde{w}_{\tau,k,\delta+1} z_k + \sum_{k,m=1}^{\delta} \tilde{w}_{\tau,k,m} z_k z_m}_{\text{second-degree term}} \end{aligned} \quad (19)$$

If we set:

$$\begin{cases} \beta_\tau = \tilde{w}_{\tau,\delta+1,\delta+1} \\ w_{\tau,k}^{[1]} = \tilde{w}_{\tau,\delta+1,k} + \tilde{w}_{\tau,k,\delta+1} & \text{for } k = 1, \dots, \delta \\ w_{\tau,k,m}^{[2]} = \tilde{w}_{\tau,k,m} & \text{for } k, m = 1, \dots, \delta \end{cases} \quad (20)$$

then (19) becomes the polynomial expansion of (18).

This enables us to express different degree polynomial expansions with a third-order tensor. The first-degree methods, e.g., ResNet [22], have $w_{\tau,k,m}^{[2]} = 0$, while SENet [25] assumes $w_{\tau,k}^{[1]} = 0$. The Π -net family assumes low-rank decomposition with shared factors, i.e., the low-rank decompositions of $\mathcal{W}_{n=1}^{[n]N}$ share factor matrices. On the contrary, our proposed PDC does not assume a sharing pattern, thus it can express independently the terms $\mathbf{W}^{[1]}$, $\mathcal{W}^{[2]}$.

C Proofs

Claim. The Squeeze-and-excitation block of (3) is a special form of second-degree polynomial term.

Proof. The global pooling function on a matrix \mathbf{C} can be expressed as $\frac{1}{hw} \vec{\mathbf{1}}^T \mathbf{C}$. The r function that replicates the channels acts on a vector \mathbf{c} and results in the expression $\vec{\mathbf{1}} \mathbf{c}^T$.

The identity $\mathbf{X} * \mathbf{a} \mathbf{b}^T = \text{diag}(\mathbf{a}) \mathbf{X} \text{diag}(\mathbf{b})$ can be used to convert the Hadamard product of (3) into a matrix multiplication [55]. Then, (3) becomes:

$$\begin{aligned} \mathbf{Y}_s = (\mathbf{Z} \mathbf{C}_1) * \vec{\mathbf{1}} \left(\left(\frac{1}{hw} \vec{\mathbf{1}}^T \mathbf{Z} \mathbf{C}_1 \right) \mathbf{C}_2 \right)^T &= (\mathbf{Z} \mathbf{C}_1) \frac{1}{hw} \text{diag}(\mathbf{C}_2^T \mathbf{C}_1^T \mathbf{Z}^T \vec{\mathbf{1}}) = \\ &= (\mathbf{Z} \mathbf{C}_1) \frac{1}{hw} \mathcal{I} \times_3 (\mathbf{C}_2^T \mathbf{C}_1^T \mathbf{Z}^T \vec{\mathbf{1}}) \end{aligned} \tag{21}$$

where as a reminder \mathcal{I} is a third-order super-diagonal unit tensor. The last equation is a second-degree term with $\Phi_1^{[2]}(\mathbf{Z}) = \mathbf{Z} \mathbf{C}_1$ and $\Phi_2^{[2]}(\mathbf{Z}) = \frac{1}{hw} \mathcal{I} \times_3 (\mathbf{C}_2^T \mathbf{C}_1^T \mathbf{Z}^T \vec{\mathbf{1}})$.

D Auxiliary experiments

Table 9: Image classification on CIFAR100 with variants of ResNet34.

Model	# param ($\times 10^6$)	Accuracy
ResNet34	21.3	0.769
H-net-ResNet	14.7	0.769
PDC-channels	36.3	0.774
PDC	10.5	0.770

D.1 Image classification with limited data

A number of experiments is performed by progressively reducing the number of training samples per class. The number of samples is reduced uniformly from the original 5,000 down to 50 per class, i.e., a $100\times$ reduction, in CIFAR10. The architectures of Table 3 (similar to ResNet18) are used unchanged; only the number of training samples is progressively reduced. The resulting Fig. 4 visualizes the performance as we decrease the training samples. The accuracy of ResNet18 decreases fast for limited training samples. SENet deteriorates at a slower pace,

steadily increasing the difference from ResNet18 (note that both share similar number of parameters). Π -net-ResNet improves upon SENet and performs better even under limited data. However, the proposed PDC-comp outperforms all the compared methods for 50 training samples per class. The difference in the accuracy between PDC and Π -net-ResNet increases as we reduce the number of training samples. Indicatively, with 50 samples per class, ResNet18 attains accuracy of 0.347, SENet scores 0.355, Π -net-ResNet scores 0.397 and PDC-comp scores 0.426, which is a 22% increase over the ResNet18 baseline.

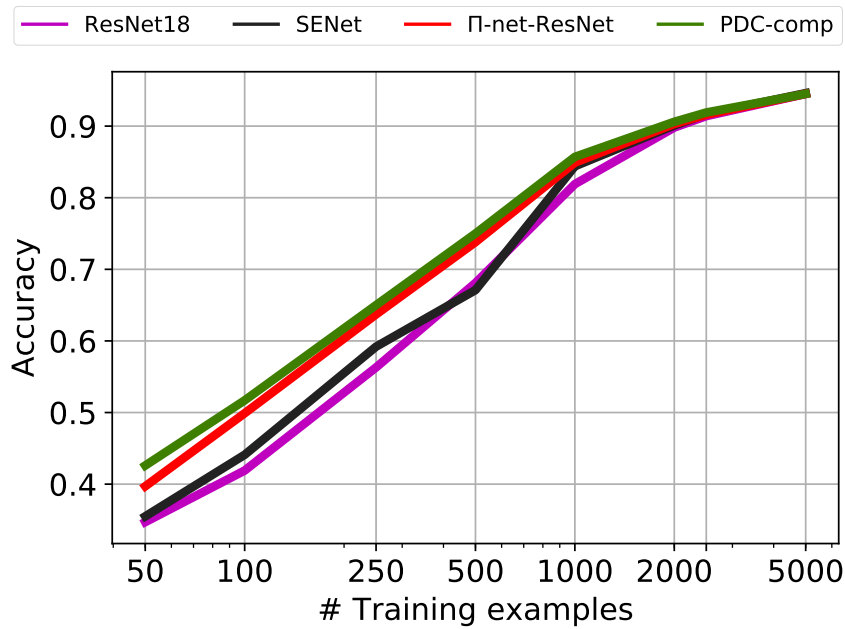


Fig. 4: Image classification with limited data. The x-axis declares the number of training samples per class (log-axis). As the number of samples is reduced (i.e., moving from right to the left), the performance gap between Π -net-ResNet and ResNet18 increases. Similarly, PDC-comp performs better than Π -net-ResNet, especially in the limited data regimes on the left.

D.2 Classification without activation functions

Typical feed-forward neural networks, such as CNNs, require activation functions to learn complex functions [23]. However, the proposed view of polynomial expansion enables capturing higher-order correlations even in the absence of activation functions. That is, the expressivity of higher-degree polynomials can be assessed without activation functions. We conduct a series of experiments on all three datasets with higher-degree polynomials. Our core experiments study

the higher-degree polynomials of Π -nets [8], versus the proposed model of (10). We also implement the ResNet without activation functions to assess how first-degree polynomials perform.

For the first experiment, we utilize ResNet18 as the backbone and test the baselines on CIFAR100. Three variations of Pi -net are considered as the compared methods: one with second-degree, one with third-degree and one with fourth-degree residual blocks. The same polynomial expansions are used for the proposed PDC. The accuracy of each method is reported in Table 10. All the variants of Π -net-ResNet and PDC exhibit a high accuracy based solely on the high-degree polynomial expansion. However, Π -net-ResNet saturates when the residual block is a third or fourth degree polynomial, while the PDC does not suffer from the same issue. On the contrary, the performance of the PDC variant with third and fourth degree residual block outperforms the second-degree residual block.

Table 10: Image classification on CIFAR100 without activation functions. Both Π -net-ResNet and PDC use high-degree polynomial expansion to achieve high accuracy even in the absence of activation functions. The proposed PDC achieves both increased performance and improves its performance when each residual block has third or fourth degree polynomial instead of second.

Model	# param ($\times 10^6$)	Accuracy
ResNet18	11.2	0.168
Π -net-ResNet	11.9	0.667
Π -net-ResNet ^[3]	11.2	0.648
Π -net-ResNet ^[4]	11.2	0.626
PDC	5.46	0.689
PDC ^[3]	11.2	0.703
PDC ^[4]	18.8	0.699

The models are also evaluated on CIFAR10 with ResNet18 and three variants of Π -nets as the backbone. Three variants of PDC with different expansion degrees are designed. The results are tabulated on Table 11. Each variant of Π -net-ResNet and PDC surpasses the 0.87 accuracy and outperform the ResNet18 by a wide margin. In contrast to Π -net-ResNet, the performance of PDC does not decrease when the degree of the residual block increases, i.e., from second to fourth-degree. Overall, PDC outperforms Π -net.

The last experiment is conducted on the Speech Commands dataset. The baseline of ResNet18 is selected, while the Π -net-ResNet is the compared method. The results in Table 12 depict the same motif: the two polynomial expansions are very expressive. Impressively, in this dataset the result without activation functions is only 0.007 decreased when compared to the respective results with activation functions. This highlights that simple datasets might not always demand activation functions to achieve high-accuracy.

Table 11: Image classification on CIFAR10 without activation functions. The results illustrate the expressiveness of the proposed model even in the absence of activation functions. Notice that PDC^[3] improves upon PDC with second-degree blocks. On the contrary, this does not happen to the compared *II*-net-ResNet.

Model	# param ($\times 10^6$)	Accuracy
ResNet18	11.2	0.391
<i>II</i> -net-ResNet	11.9	0.907
<i>II</i> -net-ResNet ^[3]	11.2	0.891
<i>II</i> -net-ResNet ^[4]	11.2	0.877
PDC	5.4	0.909
PDC ^[3]	11.2	0.918
PDC ^[4]	18.8	0.918

Table 12: Audio classification without activation functions.

Model	# param ($\times 10^6$)	Accuracy
ResNet18	11.2	0.464
<i>II</i> -net-ResNet	11.9	0.971
PDC	5.4	0.972

Table 13: **COCO object detection and segmentation results** using Mask-RCNN and Cascade Mask-RCNN. The backbone models are pre-trained ResNet18 and PDC-ResNet18 models on ImageNet-1K. We employ MMDetection with $1\times$ schedule.

backbone	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	AP ^{mask}	AP ₅₀ ^{mask}	AP ₇₅ ^{mask}
Mask-RCNN $1\times$ schedule						
ResNet18	33.9	53.9	36.2	31.0	50.9	33.0
PDC-ResNet18	34.8	55.2	37.4	31.8	52.2	34.1
Cascade Mask-RCNN $1\times$ schedule						
ResNet18	37.3	54.8	40.4	32.6	52.2	34.9
PDC-ResNet18	38.1	55.9	41.7	33.2	53.3	35.7

D.3 Object detection and segmentation

We adopt MS COCO 2017 [39] as the primary benchmark for the experiments of object detection and segmentation. We use the train split (118k images) for training and report the performance on the val split (5k images). We employ standard evaluation metrics for COCO dataset, where multiple IoU thresholds from 0.5 to 0.95 are applied. The detection results are evaluated with mAP.

We use the final model weights from ImageNet-1K pre-training as network initializations and fine-tune Mask R-CNN [21] and Cascade Mask R-CNN [4] on

the COCO dataset. Following default settings in MMDetection, we use the $1\times$ schedule (i.e., 12 epochs).

Table 13 shows object detection and instance segmentation results comparing ResNet18 and the proposed PDC-ResNet18. As we can see from the results, the proposed PDC-ResNet18 achieves an obvious better performance than the baseline ResNet18 in terms of the box and mask AP, confirming the effectiveness of the proposed polynomial learning scheme.

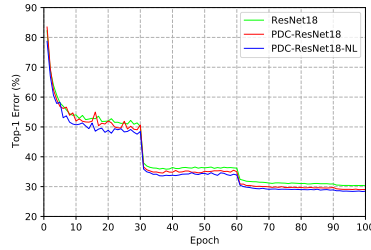


Fig. 5: Top-1 validation error on ImageNet with proposed PDC and NL methods throughout the training.