
PECAN: A Product-Quantized Content Addressable Memory Network

Jie Ran, Rui Lin, Jason Chun Lok Li, Jiajun Zhou, Ngai Wong

Department of Electrical and Electronic Engineering,

The University of Hong Kong, Hong Kong

Email Address: {jieran, linrui, u3524157, jjzhou, nwong}@eee.hku.hk

Abstract

A novel deep neural network (DNN) architecture is proposed wherein the filtering and linear transform are realized solely with product quantization (PQ). This results in a natural implementation via content addressable memory (CAM), which transcends regular DNN layer operations and requires only simple table lookup. Two schemes are developed for the end-to-end PQ prototype training, namely, through angle- and distance-based similarities, which differ in their multiplicative and additive natures with different complexity-accuracy tradeoffs. Even more, the distance-based scheme constitutes a truly multiplier-free DNN solution. Experiments confirm the feasibility of such Product-Quantized Content Addressable Memory Network (PECAN), which has strong implication on hardware-efficient deployments especially for in-memory computing.

1 Introduction

Deep neural networks (DNNs) have achieved breakthroughs in various applications including classification [21], object detection [15] and semantic segmentation [29], etc. Nonetheless, the massive amount of parameters and computation make it difficult for both training and inference on edge devices with constrained hardware resources. Numerous efforts have been made to reduce the network complexity while preserving the output accuracy. Among various schemes, some are low-bitwidth neural networks using binary weights [4, 18, 27], replacing the expensive multiplications with cheaper sign flip operations during inference. Some approaches substitute multiplications with additions and bit-wise shifts. AdderNet [2] realizes convolution (in the sense of similarity matching) by l_1 -distance between the activation and weights, and maintains competitive output accuracy. ShiftCNN [8] is based on a power-of-two weight representation for converting convolutional neural networks (CNNs) without retraining. Among works that aim to improve the memory efficiency and performance of shift neural networks, DeepShift [7] is a framework for training low-bitwidth neural networks from scratch to replace multiplication with bit-wise shift and sign flip. All these works, despite specific implementations, still adhere to the traditional DNN architecture.

This work attempts to detach a neural network from its regular filtering operation and replace it with an associative memory, aka content addressable memory (CAM), whereby the content is derived from prototypes of product quantization [11]. Such framework, dubbed Product-Quantized Content Addressable Memory Network (PECAN), combines the storage and compute into one place, and is particularly suitable for the fast-emerging in-memory computing. The codebook/table lookup during inference also makes PECAN hardware-friendly and positions it as a strong candidate for edge artificial intelligence (AI). This is also warranted by the readiness in commodity platforms like FPGAs with CAM support, as well as next-generation memristive microelectronics like resistive random-access memory (RRAM) wherein a CAM is inherent to an RRAM crossbar [19, 12].

Our proposed PECAN is inspired by the lately proposed MADDNESS [1] that utilizes product quantization and table lookup to truly omit multipliers in matrix-matrix products. However, the main

contribution of MADDNESS, namely, the hash function for prototype matching, is heuristic and non-differentiable, thus making it incompatible with a learning framework. In fact, the authors also remark it will take several more papers to consolidate the framework for DNNs.

PECAN exactly fills this void by its end-to-end learnable PQ-based DNN architecture. The closest work to ours is differentiable product quantization (DPQ) [3], but *for the first time* we demonstrate its multi-layer feasibility and enrich DPQ prototype matching (viz. a similarity search) with an l_1 -distance metric. The latter comes from the lately proposed AdderNet [2] wherein the l_1 metric is utilized in a different context of CNN filtering, whereas our work is the first to show its feasibility for training prototypes in the DPQ setting. To our best knowledge, PECAN is a brand new architecture that transcends regular DNN filtering and uses similarity search and table lookup for inference. This allows it to be compatible with simple hardware without the need of dedicated neural engines, especially edge devices where compute and storage resources are limited. Our major contributions are: 1) A first-of-its-kind, end-to-end learnable CAM-based DNN. PECAN is hardware-generic and friendly to almost all hardware platforms especially those with built-in CAM support, and represents a strong candidate for edge AI deployment; 2) Two similarity measures in PECAN, based on angle and distance, to investigate the trade-offs between computation complexity and accuracy; 3) Joint fine-tuning and co-optimization of weight matrices and PQ prototypes, which permits PECAN to train from scratch; 4) A *totally* multiplier-free DNN via the distance-based PECAN.

2 Related Work

For efficient edge deployment, binary neural networks (BNNs) [10, 18, 16] exclusively make use of the logical XNOR operation that obviates regular multipliers, but in principle they are still doing 1-bit multiplication. Moreover, though BNNs have gone through major improvements in recent years, their top-1 accuracies measured on large-scale datasets are still noticeably lower than their full-precision counterparts. Indeed, most BNN implementations are only partial in the sense that the first and final layers are still using full-precision weights and activations [27, 25].

Other works replace multiplication with addition [2] or bit-shift operations [8, 30, 7], or both [28]. Specifically, AdderNet makes novel use of l_1 -norm difference and adders to do template matching required in a CNN. Yet it still employs multipliers for the necessary batch normalization to bring back signed pre-activations. Progressive kernel based knowledge distillation (PKKD) AdderNet [24] improves the performance of the vanilla AdderNet. AdderNet with Adaptive Weight Normalization (AWN) [6] further alleviates the curse of instability of running mean and variance in batch normalization layers. Applying bitwise shift on an element is mathematically equivalent to multiplying it by a power of two, and sign flipping is introduced to represent negative numbers. Although these works focus on largely multiplier-free DNNs, they still build on the traditional architectures.

The proposed PECAN is motivated by MADDNESS which realizes multiplier-free matrix-matrix product using hashing and table lookup rather than multiply-add operations. Although it achieves orders of speedups compared to existing approximate matrix multiplication (AMM) methods, the proposed hashing functions are not differentiable and not amenable to DNN training. DPQ [3] is proposed for end-to-end embedding, but it is only single-layer and targets word embedding, and still requires full-precision multiplication to obtain distances between the input and matching keys.

3 PECAN

The convolution operation in a CNN is conceptually illustrated as a window sliding across the c_{in} -channel input feature (cf. Fig. 1(a)). Actual implementations often unfold the convolution into a matrix-matrix product (cf. Fig. 1(b)). Specifically, the `im2col` command stretches the input entries covered in each filter stride into a column and concatenates the columns into a matrix X , whereas the kernel tensors are reshaped into a filter matrix F , such that PQ can be used to approximate FX . For an intermediate CNN layer, consider the flattened feature matrix $X \in \mathbb{R}^{c_{in}k^2 \times H_{out}W_{out}}$, where c_{in} and k are the number of input channels and the kernel size, H_{out} and W_{out} are height and width of the output feature, codebooks $C \in \mathbb{R}^{c_{in}k^2 \times p}$ are assigned with parameters to construct an embedding table for the features, where p is the number of choices for each *codebook* $C^{(j)}$, $j = 1, 2, \dots, D$. $C_m^{(j)} \in \mathbb{R}^d$ are called *prototypes*, $m = 1, 2, \dots, p$ (cf. Fig. 1(c)). It is natural to set each prototype in PQ to be a $k^2 \times 1$ subvector (viz. same size as a vectorized kernel), with p prototypes in each

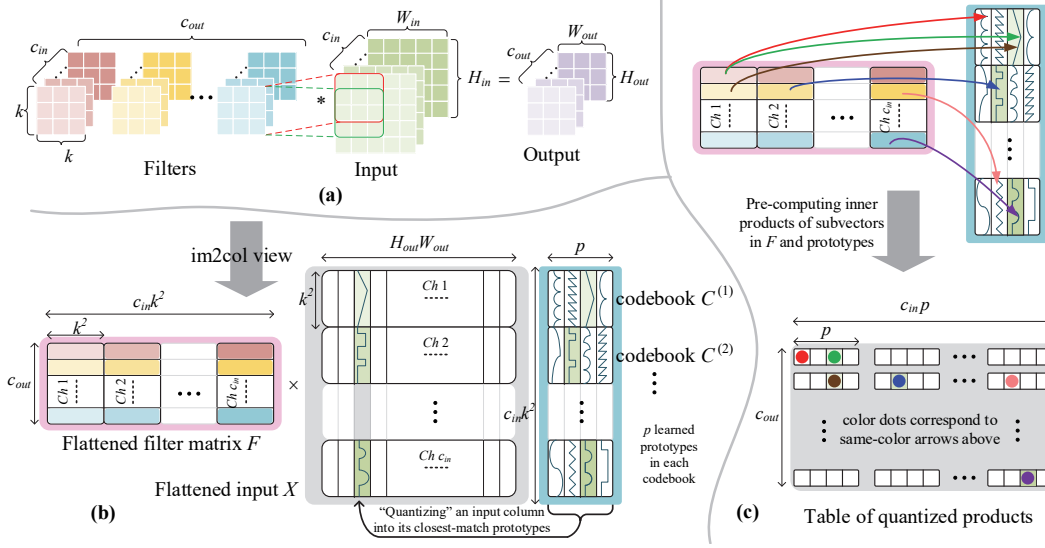


Figure 1: CNN convolution in its (a) conceptual form; (b) equivalent matrix-matrix-multiply by flattening the filters and input features via `im2col`, with mapping of input data sub-columns onto the closest prototypes in different codebooks; (c) Precomputing inner products of F -subvectors and prototypes to form a lookup table.

of the c_{in} input channels according to the patterns of flattened matrices. With this setting, there are two main components in a trained PECAN that require memory storage in each layer, namely, i) pc_{in} prototypes for “quantizing” the input subvectors; ii) $c_{out}c_{in}p$ inner product values between the (sub)rows in F and each prototype.

In short, PECAN is mapping (quantizing) the original input features onto prototype patterns in compact codebooks, then multiplication between weights (F) and features (X) can be approximated by lookup table operation during inference. Below we elaborate two content addressing techniques (i.e. similarity matching) approaches based respectively on angle (dot product) and distance (l_1 -norm) which are both end-to-end learnable. Accordingly, these two schemes are dubbed PECAN-A and PECAN-D, which cover both ends of complexity-accuracy spectrum: The angle-based scheme uses multiplicative operations and generally leads to higher output accuracy, whereas the distance-based one uses additive operations and is much more lightweight at the expense of slight accuracy loss.

3.1 PECAN-A: Angle-Based Similarity Measure

A scaled dot-product attention module [23], widely used in Transformers, computes the dot products of queries with keys and applies a softmax function to obtain the weights on the values.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (1)$$

where d_k is the dimension of keys, which serves as a scaling factor. Generally, Q , K and V are obtained from three distinct learned projection matrices. However, different from self-attention, we learn the keys K (viz. prototypes in PQ) directly without the intermediate linear transforms, and make V equal to K . For PECAN-A, we compute the approximated matrix \tilde{X} by splitting its rows into $D = c_{in}$ groups, each with subvectors of dimension $d = k^2$, and get the attention scores $K_i^{(j)}$ to formulate the combination of prototypes $C_m^{(j)}$:

$$K_i^{(j)} = \text{softmax}((C^{(j)})^T X_i^{(j)}), \quad \tilde{X}_i^{(j)} = C^{(j)} K_i^{(j)}, \quad (2)$$

where $i = 1, 2, \dots, H_{out}W_{out}$. Since the dot product distance function with softmax is differentiable, mapping features to prototypes can be learned end-to-end. It is worth noting that all intermediate features are replaced with the combination of learned prototypes after training.

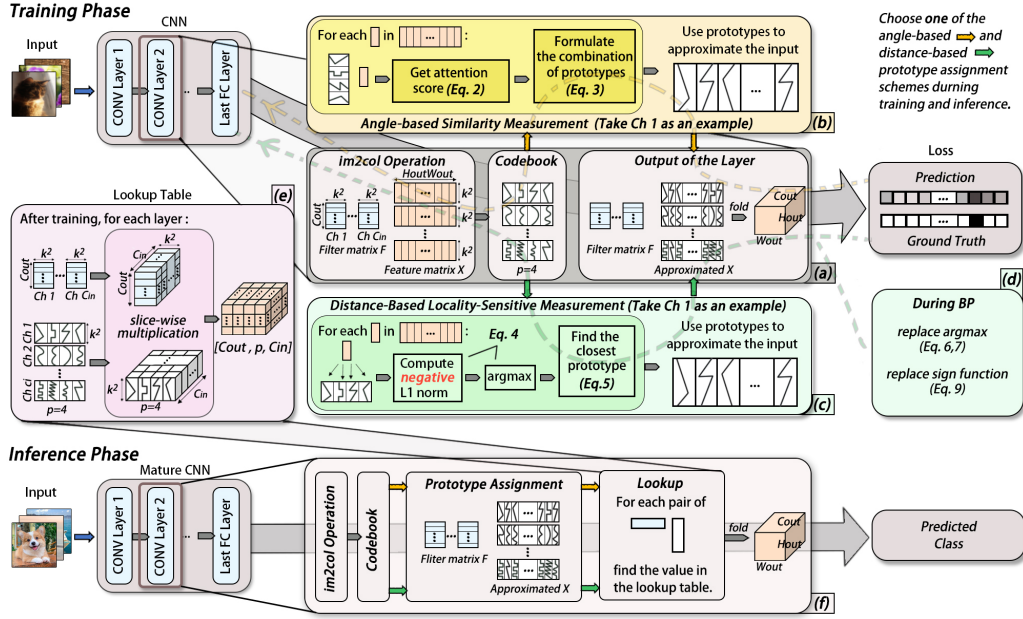


Figure 2: The proposed PECAN architecture. (a) The training phase is mainly composed of template matching for each subvector in the flattened feature map matrices after `im2col` operation. When approximating subvectors with the closest prototypes, PECAN-A and PECAN-D adopt different assignment schemes. (b) For PECAN-A, an attention module compares the subvectors with each of the prototypes in the same group. Subsequently, the resulting scores are subjected to the weighted sum to produce the approximate feature matrix. (c) For PECAN-D, the similarity is measured with a sign flip l_1 -norm and the approximation is selected with `argmax` function. (d) Since the `argmax` is not differentiable and the gradient of l_1 -norm is discrete $(1, -1, 0)$, we propose Eq. (4, 5) and (6) to do the backpropagation. (e) After getting the converged neural network, we calculate the slice-wise product between convolution filters and prototypes, and store the results in the memory. (f) In the inference phase, we only need to calculate the distance of feature maps with a small number of prototypes and look up in the stored memory to get the quantized output.

3.2 PECAN-D: Distance-Based Similarity Measure

Now we attempt to get rid of all multipliers. To achieve this, we make use of only l_1 -norm difference for the so-called template matching, namely, finding the closest match through absolute difference which involves only subtraction. Specifically, in this distance-based framework, l_1 -norm is applied in order to discard multiplication:

$$k_i^{(j)} = \arg \max_m -\|X_i^{(j)} - C_m^{(j)}\|_1, \tilde{X}_i^{(j)} = C^{(j)} \text{one_hot}(k_i^{(j)}), \quad (3)$$

where $K_i^{(j)} = \text{one_hot}(k_i^{(j)})$ denotes a p -dimensional vector with the $k_i^{(j)}$ -th entry as 1 and others 0. To enable optimization for prototypes with the non-differentiable function `argmax`, we approximate it with a differentiable softmax function defined as follows:

$$\tilde{K}_i^{(j)} = \frac{\exp(-\|X_i^{(j)} - C_m^{(j)}\|_1/\tau)}{\sum_{m'} \exp(-\|X_i^{(j)} - C_{m'}^{(j)}\|_1/\tau)}, \quad (4)$$

where τ is the temperature to relax the softmax function. Note that Eq. (4) can be considered as the proportion of Laplacian kernels when $\tau \neq 0$. It relies on the observation that the positive definite function $k(X_i^{(j)}, C_m^{(j)}) = \exp(-\|X_i^{(j)} - C_m^{(j)}\|_1/\tau)$ here defines an inner product and a lifting function ϕ such that the inner product $\langle \phi(X_i^{(j)}), \phi(C_m^{(j)}) \rangle$ can be computed quickly using the kernel trick [17].

Now the approximated index $\tilde{K}_i^{(j)}$ is fully differentiable when $\tau \neq 0$. However, this yields the combination of prototypes for $\tilde{X}_i^{(j)}$ again, while we need $\tau \rightarrow 0$ to get discrete indices during the

Table 1: Inference complexities of PECAN-A and PECAN-D.

Method	Layer	#Add.	#Mul.
Baseline	CONV	$c_{in}H_{out}W_{out}k^2c_{out}$	$c_{in}H_{out}W_{out}k^2c_{out}$
	FC	$c_{in}c_{out}$	$c_{in}c_{out}$
PECAN-A	CONV	$pDH_{out}W_{out}(d+c_{out})$	$pDH_{out}W_{out}(d+c_{out})$
	FC	$pD(d+c_{out})$	$pD(d+c_{out})$
PECAN-D	CONV	$DH_{out}W_{out}(2pd+c_{out})$	0
	FC	$D(2pd+c_{out})$	0

forward inference. To this end, we follow [3] and define a new index to solve both non-differentiable and discrete problems in one go. Specifically, in the forward and backward passes during training, we adopt

$$\tilde{K}_i^{(j)}(\tau \neq 0) - sg \left(\tilde{K}_i^{(j)}(\tau \neq 0) - \tilde{K}_i^{(j)}(\tau = 0) \right), \quad (5)$$

where sg is *stop gradient*, which takes the identity function in the forward pass and drops the gradient inside it in the backward pass. Based on this, we can now use the argmax function in the forward pass and softmax function during backpropagation. However, the partial derivative of the distance $d_{im}^{(j)} = -\|X_i^{(j)} - C_m^{(j)}\|_1$ with respect to codebook subvector $C_m^{(j)}$ is a sign function:

$$\frac{\partial d_{im}^{(j)}}{\partial C_m^{(j)}} = \text{sgn}(X_i^{(j)} - C_m^{(j)}), \quad \frac{\partial d_{im}^{(j)}}{\partial C_m^{(j)}} = \tanh \left(a(X_i^{(j)} - C_m^{(j)}) \right) \text{ where } a = \exp\left(\frac{4e}{E}\right), \quad (6)$$

where $\text{sgn}(\cdot)$ is the sign function and takes the values of $\{+1, 0, -1\}$. Such zero gradient almost everywhere makes it impossible to train a neural network. In this regard, we adopt the right term in Eq. (6) to replace the gradient, where e is the current epoch and E the total number of training epochs. Fig. 3 shows this epoch-aware approximation to the sign function *w.r.t.* values of $\frac{e}{E}$ as epoch increases during training. In the early stage, the function is smoother for stable training. As the training progresses, the approximation gradually turns into the sign-like function.

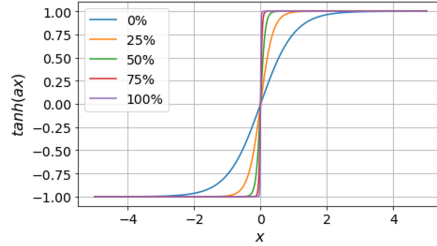


Figure 3: Approximation to the sign gradient for different $\frac{e}{E}$.

3.3 Inference Details and Complexity

For the original `im2col` convolution, the computation complexity is $O(c_{in}H_{out}W_{out}k^2c_{out})$. During inference, our method includes two stages, the first is to get the indices by computing the distance between the flattened features and prototypes, while the second is to retrieve the product between weights and prototypes computed in advance, i.e., a simple table lookup. The inference algorithm for both PECAN variants is given in Algorithm 1.

Table 1 illustrates the number of multiplication and addition operations in convolution and fully-connected layers for the traditional CNNs, angle-based and distance-based PECAN during the inference phase. Note that the fully-connected layer can be regarded as a convolution layer when $k = H_{out} = W_{out} = 1$. Instead of using the specialized setting of $D = c_{in}$ and $d = k^2$, we further consider the more general case in Table 1 where the group number D and dimension of prototypes d satisfy $Dd = c_{in}k^2$. Choosing smaller p and D will reduce the computation complexity for both PECAN-A and PECAN-D. Specifically, in order to limit multiplication complexity in PECAN-A to be smaller than the baseline, we need $p \leq \min(\lambda c_{out}, (1 - \lambda)d)$ with $\lambda \in (0, 1)$. This constraint is also taken into consideration in the experiment section. Note that by design, PECAN-D needs *no multiplication during inference*, thus making it genuinely totally multiplier-less.

Algorithm 1 Inference Algorithm of PECAN

Input: Codebook $C \in \mathbb{R}^{c_{in}k^2 \times p}$, 4-D learned kernel tensor $\mathcal{K} \in \mathbb{R}^{c_{out} \times c_{in} \times k \times k}$, unfolded features $X \in \mathbb{R}^{c_{in}k^2 \times H_{out}W_{out}}$.

Output: The approximated convolution output $\tilde{Y} \in \mathbb{R}^{c_{out} \times H_{out}W_{out}}$

```
1: Permute and reshape weights to  $W_1 \in \mathbb{R}^{D \times c_{out} \times d}$ , codebooks to  $C_1 \in \mathbb{R}^{D \times d \times p}$ 
2: for  $j$  in  $\{1, 2, \dots, D\}$  do
3:    $Y^{(j)} = W_1^{(j)} C_1^{(j)} \in \mathbb{R}^{c_{out} \times p}$ 
4: end for
5: for  $i$  in  $\{1, 2, \dots, H_{out}W_{out}\}$  do
6:   if PECAN-A then
7:      $\tilde{Y}_i = \sum_{j=1}^D Y^{(j)} \text{softmax}(C^{(j)T} X_i^{(j)})$ 
8:   end if
9:   if PECAN-D then
10:     $k_i^{(j)} = \arg \max_m -\|X_i^{(j)} - C_m^{(j)}\|_1$ 
11:     $\tilde{Y}_i = \sum_{j=1}^D Y_{k_i^{(j)}}^{(j)}$ 
12:   end if
13: end for
14: return Concatenate  $(\tilde{Y}_1, \tilde{Y}_2, \dots, \tilde{Y}_{H_{out}W_{out}})$ 
```

4 Experiments

To demonstrate the effectiveness of the proposed PECAN and further benchmark the differences between its two variants (PECAN-A and PECAN-D), we apply PECAN to the classification tasks, taking MNIST [5], CIFAR-10 and CIFAR-100 [13] as datasets. The models employed in this section include modified LeNet5, VGG-Small [27], ResNet20 and ResNet32 [9]. The experiments explore two PECAN training strategies, namely, co-optimization on weights and codebooks and uni-optimization on codebooks. We also compare PECAN with state-of-the-art (SOTA) approaches that aim to reduce the amount of multiplication operations, namely, XNOR-Net [18], IR-Net [16], FDA-BNN [25], ReCU [26], SD-BNN [27] and AdderNet [2]. Additionally, we investigate how the dimension and number of groups of the codebook affect the performance of PECAN, validate the necessity of our training strategies, and provide visual results to confirm the approximation capability of the prototypes.

Implementation Details. When using the PECAN framework to train the modified LeNet5 with kernels of size 3×3 on MNIST, we employ the uni-optimization strategy that only updates the prototypes with the trained weights being frozen. The prototypes are trained for 150 epochs. The learning rate is set to 0.01 initially, decaying every 50 epochs. To implement the PECAN framework for the CIFAR-10 and CIFAR-100 tasks, we use the co-optimization strategy and set the training epochs for PECAN-A and PECAN-D as 150 and 300, respectively. The learning rate for PECAN-A follows the LeNet5 scheme, while that of PECAN-D is initialized as 0.001, decaying at epoch 200. For both datasets, we employ `softmax` function and set the temperature τ at 1 and 0.5 for PECAN-A and PECAN-D, respectively. We set the batch size to 64, and use cross-entropy as the loss function, which is optimized by Adam. All experiments are run on a machine equipped with four NVIDIA Tesla V100 GPU with 24GB frame buffer, and all codes are implemented by PyTorch.

4.1 Modified LeNet5 on MNIST

To quickly zoom into the superiority of PECAN, the amount of required addition and multiplication operations and the detailed codebook information for each layer in the modified LeNet5 are shown in Appendix Table A2. Focusing on the second and third columns, it is noticeable that PECAN-A has fewer multiplications and additions compared with the baseline, and PECAN-D needs no multiplication at all. For the codebook settings, it is seen that the number of prototypes p used in PECAN-A is much fewer than that of PECAN-D for all five layers. We adopt this setting considering the gaps between the representation capabilities of PECAN-A and PECAN-D. By adjusting the

weights assigned to prototypes, PECAN-A is expected to better approximate the features with limited choices, i.e., a smaller p .

Table 2: Experiment results of LeNet on MNIST.

Model	#Add.	#Mul.	Acc.(%)
Baseline	248.10K	248.10K	99.41
PECAN-A	196.88K	196.88K	99.25
PECAN-D	2.00M	0	99.01

vs 99.41%). The accuracy of PECAN-A achieves 99.25%, which is merely 0.16% lower than the original LeNet5. However, PECAN-A performs fewer operations. To this end, the LeNet5 example demonstrates the effectiveness of the PECAN framework, and shows the advantages of PECAN-A and PECAN-D from different perspectives.

Table 3: Experiment results on CIFAR10.

Model	Method	#Add.	#Mul.	Accuracy (%)
VGG-Small	Baseline	0.61G	0.61G	91.21
	PECAN-A	0.54G	0.54G	91.82
	PECAN-D	0.37G	0	90.19
ResNet20	Baseline	40.55M	40.55M	92.55
	PECAN-A	38.12M	38.12M	90.32
	PECAN-D	211.71M	0	87.88
ResNet32	Baseline	68.86M	68.86M	92.85
	PECAN-A	64.20M	64.20M	90.53
	PECAN-D	353.26M	0	88.46

The performance and the required number of addition and multiplication operations of the whole modified LeNet5 employing PECAN-A and PECAN-D schemes are summarized in Table 2. It is worth noting that LeNet with PECAN-D is multiplier-free, and maintains the good performance compared with the baseline (99.01%

Table 4: Experiment results on CIFAR100.

Model	Method	#Add.	#Mul.	Accuracy (%)
VGG-Small	Baseline	0.61G	0.61G	67.84
	PECAN-A	0.54G	0.54G	69.21
	PECAN-D	0.37G	0	60.43
ResNet20	Baseline	40.56M	40.56M	69.55
	PECAN-A	38.12M	38.12M	63.15
	PECAN-D	211.71M	0	58.01
ResNet32	Baseline	68.86M	68.86M	70.57
	PECAN-A	64.20M	64.20M	64.13
	PECAN-D	353.27M	0	58.26

4.2 VGG and ResNet on CIFAR-10/100

After the proof-of-concept on LeNet5, we proceed to VGG-Small and ResNet20/32 on CIFAR-10 and CIFAR-100. VGG-Small is a simplified VGGNet [20] with only one fully-connected layer. The size of the output feature maps and the corresponding codebook information for each layer are provided in Appendix Table A3. We remark that the bottom row of each block in the table represents the FC layer, while the rows above represent the CONV layers. The number of required addition and multiplication operations and the accuracy of the models are summarized in Table 3, where it can be seen that the VGG-Small baseline has 0.61G multiplication and addition operations with 91.21% accuracy. Since batch normalization can be folded into convolution layers in the inference stage, we do not count FLOPs for both baseline and PECAN. We find that PECAN-A only performs 0.54G multiplications while reaching 91.82% accuracy on CIFAR-10, which is even higher than the baseline, similar performance can be obtained on CIFAR-100 in Table 4. A possible reason is that PECAN experiences less information loss for shallower CNNs, and bigger input channels allow more groups of prototypes to improve the representation capability. This assumption is also validated by the experiments on ResNet20/32 that are deeper than VGG-Small but with smaller input channels.

Although PECAN-D has an accuracy drop compared with the baseline, it eliminates all multiplications during inference. For VGG-Small, the number of additions reduces to 0.37G while accuracy drops only around 1% compared with the baseline on CIFAR10. To boost the performance, we use smaller size for subvectors in PECAN-D as shown in the last column of Appendix Table A3, at the expense of more computation.

4.3 Comparison with AdderNet

We compare PECAN-D with AdderNet on VGG-Small in Table 5. It should be emphasized that batch normalization is not taken into consideration in this table, it can not be folded into AdderNet layer so multiplication is indispensable. For VGG-Small, the memory cost is so high that even four NVIDIA Tesla V100 GPUs are not able to train successfully. As shown in the table, the proposed PECAN-D with only 0.37G additions achieves a 90.19% accuracy on VGG-Small. Here we showcase the efficacy of PECAN on larger models from a hardware perspective. In the Intel VIA Nano 2000 CPU (used in the AdderNet paper), the latency cycles of float multiplication and addition are 4 and 2, respectively. PECAN-D of VGG-Small model will incur $\sim 720M$ (cycles) while that of a CNN is

~3660M. The power consumption ratio of 32bit multiplication and addition units is 4:1. Power-wise and latency-wise, PECAN-D network is more efficient than both AdderNet and regular CNN.

Table 5: Comparison with AdderNet.

Model	Method	# Mul.	# Add.	Accuracy (%)	Normalized Power	Latency(cycles)
VGG-Small	CNN	0.61G	0.61G	93.80	8.24	3.66G
	AdderNet	0	1.22G	N.A.	3.30	2.44G
	PECAN-D	0	0.37G	90.19	1	0.72G

4.4 Ablation Study

4.4.1 Prototypes Dimension

To investigate the effect of the prototype dimension on the performance of both PECAN-A and PECAN-D, we conduct comparative experiments on ResNet20 (on CIFAR-10 dataset) by reducing the dimension from k^2 to k , and then increasing up to c_{in} . Subsequently, the number of groups changes from c_{in} to kc_{in} and k^2 correspondingly. Except for the prototype dimension, we keep other settings (e.g., number of prototypes p in each layer, and learning rate decay scheme) the same as the experiments in Table 3. The results are visualized in Fig. 4, wherein the prototype dimension increases from left to right for each set of bar charts. It is observed that the performances of PECAN-A and PECAN-D have different trends when the prototype dimension increases. Compared with PECAN-D, the angle-based PECAN-A is more robust, and its accuracy does not change sharply under the three cases. It achieves the best performance when the prototype dimension is k^2 , while still maintaining decent accuracy when the dimension changes to k and c_{in} . On the contrary, the performance of PECAN-D is more sensitive, which is inversely proportional to the prototype dimension. It is intuitive since approximation in the fine-scale (viz. larger groups with smaller dimensions) is expected to be more accurate. We remark that PECAN-A, which has weighted prototypes, trades for robustness at all scales by the higher computational complexity compared with PECAN-D.

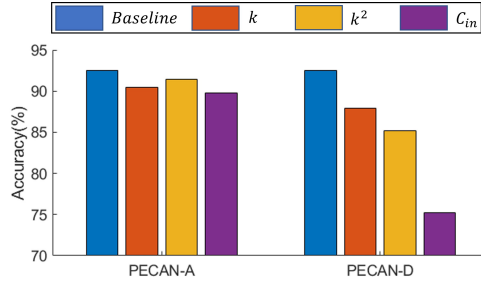


Figure 4: Accuracy of ResNet20 on CIFAR10 using different dimensions of subvectors including k , k^2 and c_{in} for both PECAN-A and PECAN-D.

4.4.2 Freezing Weights during Training

For both angle- and distance-based measures, we freeze the pretrained weights and only train the matched codebooks in the MNIST example, but when moving to a larger dataset, both weights and prototypes are not frozen but are updated during training from scratch. To further clarify the reason of this choice, we also train VGG-Small which has the same setting as in Table 3 on CIFAR10.

Table 6: Effects of training strategies on PECAN accuracy.

Model	From Scratch	Freeze Weights	Accuracy(%)
Baseline	✓	✗	91.21
PECAN-A/D	✓	✗	91.82/90.19
PECAN-A/D	✗	✓	91.76/87.43

However, different from training from scratch, we only unfreeze prototypes and start from the pretrained mature CNN. Specifically, we initialize all convolution filters with the pretrained model and only learn those prototypes in each layer. Experimental results are shown in Table 6. As seen from the table, updating prototypes only still has accuracy gap especially for PECAN-D compared with the one learning both weights and prototypes from scratch. A possible reason is that the convolution weights in the pretrained model are not well-matching with the templates.

4.4.3 Visualization of Prototypes

To visually inspect the effectiveness of PECAN-D in CNNs, we take the intermediate convolution layers of VGG-Small and plot the patterns of the feature maps before and after replacement. In Fig. 5,

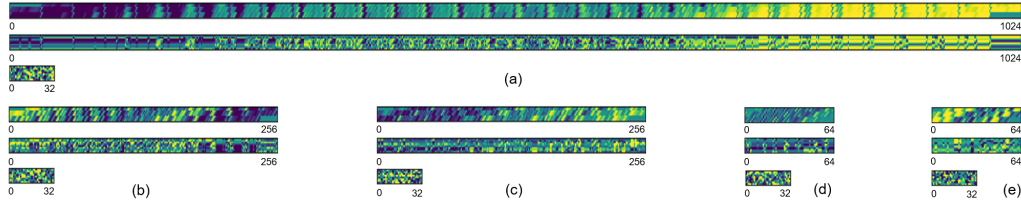


Figure 5: The flattened features and codebooks for five different layers in VGG-Small, (a)-(e) for conv1-conv5. For each subfigure, the upper image is the input feature after `im2col` operation, the second image shows the approximation matrix after substitution with PECAN-D which is composed of the corresponding codebook shown in the third row. The y -axis is the dimension of each subvector k^2 . The x -axis represents the size of output feature maps $H_{out}W_{out}$ for the first two rows, and denotes the number of prototypes for the third row.

we select the first channel of the flattened feature maps and visualize the matrices. The dimension of all subvectors is set as $k^2 = 9$. As can be seen, though the number of prototypes is limited for each convolution layer, the quantized feature maps can still preserve the basic patterns after training.

5 Discussion

To boost the performance for PECAN, we can choose small size of prototypes in our experiments which might incur a challenge for memory footprint. Reducing the number of prototypes and increasing the subvector lengths can lower the memory cost, but this may harm the model accuracy. To this end, We propose further means to save memory: For PECAN-A we can exploit the projectors inherent to attention to shrink dimensions of query and key (corresponding to weights & prototypes in PECAN-A). Even more interesting, for PECAN-D, take the 2nd CONV layer of ResNet20 on CIFAR10 as an example, only 26 out of 64 prototypes are used in the inference stage, meaning all other prototypes and lookup entries can be pruned without affecting accuracy. Fig. 6 illustrates the sparse usage count of each prototype in the first group of codebooks for 18 intermediate CNN layers. We will report these exciting results in the follow-up work of PECAN, as these are beyond the central theme of this paper.

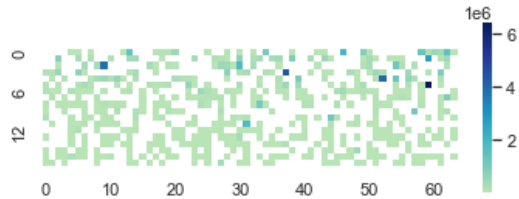


Figure 6: Call frequencies of 64 prototypes in the middle 18 CNN layers @ ResNet20. The x -axis represents the indices of prototypes and the y -axis is the order of the 18 layers. White grid cells denote 0 times of usage.

6 Conclusion

A brand new DNN architecture called PECAN is proposed which transcends the regular DNN linear transform, and replaces it by product quantization and table lookup. Both angle- and distance-based measures are developed for similarity matching of prototypes in product quantization for different complexity-accuracy tradeoffs. The distance-based PECAN, to our knowledge, is the *first* neural network that is multiplier-less and uses only adders all over. PECAN is end-to-end trainable and infers only through a content addressable memory (CAM)-like, similarity search protocol. It facilitates a lightweight and hardware-generic solution favorable for edge AI, and fits perfectly into the in-memory-computing regime. Experiments have shown that PECAN exhibits accuracies on par with multi-bit networks even without using multipliers. We expect more advancement on top of this interesting PECAN framework will follow after this debut.

References

- [1] Davis Blalock and John Guttag. Multiplying matrices without multiplying. *arXiv preprint arXiv:2106.10860*, 2021.
- [2] Hanqing Chen, Yunhe Wang, Chunjing Xu, Boxin Shi, Chao Xu, Qi Tian, and Chang Xu. Addernet: Do we really need multiplications in deep learning? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1468–1477, 2020.
- [3] Ting Chen, Lala Li, and Yizhou Sun. Differentiable product quantization for end-to-end embedding compression. In *International Conference on Machine Learning*, pages 1617–1626. PMLR, 2020.
- [4] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pages 3123–3131, 2015.
- [5] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [6] Minjing Dong, Yunhe Wang, Xinghao Chen, and Chang Xu. Towards stable and robust addernets. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [7] Mostafa Elhoushi, Zihao Chen, Farhan Shafiq, Ye Henry Tian, and Joey Yiwei Li. Deepshift: Towards multiplication-less neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2359–2368, 2021.
- [8] Denis A Gudovskiy and Luca Rigazio. Shiftcnn: Generalized low-precision architecture for inference of convolutional neural networks. *arXiv preprint arXiv:1706.02393*, 2017.
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [10] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. *Advances in neural information processing systems*, 29, 2016.
- [11] Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2010.
- [12] Geethan Karunaratne, Manuel Schmuck, Manuel Le Gallo, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas Rahimi. Robust high-dimensional memory-augmented neural networks. *Nat. Commun.*, 12, 2021.
- [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [14] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [15] Yudong Liu, Yongtao Wang, Siwei Wang, TingTing Liang, Qijie Zhao, Zhi Tang, and Haibin Ling. Cbnet: A novel composite backbone network architecture for object detection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11653–11660, 2020.
- [16] Haotong Qin, Ruihao Gong, Xianglong Liu, Mingzhu Shen, Ziran Wei, Fengwei Yu, and Jingkuan Song. Forward and backward information retention for accurate binary neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2250–2259, 2020.
- [17] Ali Rahimi, Benjamin Recht, et al. Random features for large-scale kernel machines. In *NIPS*, volume 3, page 5. Citeseer, 2007.
- [18] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, pages 525–542. Springer, 2016.

- [19] Yuan Ren, Rui Lin, Jie Ran, Chang Liu, Chaofan Tao, Zhongrui Wang, Can Li, and Ngai Wong. Batmann: A binarized-all-through memory-augmented neural network for efficient in-memory computing. In *2021 IEEE 14th International Conference on ASIC (ASICON)*, pages 1–4, 2021. doi: 10.1109/ASICON52560.2021.9620292.
- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [21] Yanan Sun, Bing Xue, Mengjie Zhang, Gary G Yen, and Jiancheng Lv. Automatically designing cnn architectures using the genetic algorithm for image classification. *IEEE transactions on cybernetics*, 50(9):3840–3854, 2020.
- [22] Asher Trockman and J Zico Kolter. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [24] Yixing Xu, Chang Xu, Xinghao Chen, Wei Zhang, Chunjing Xu, and Yunhe Wang. Kernel based progressive distillation for adder neural networks. *arXiv preprint arXiv:2009.13044*, 2020.
- [25] Yixing Xu, Kai Han, Chang Xu, Yehui Tang, Chunjing Xu, and Yunhe Wang. Learning frequency domain approximation for binary neural networks. *arXiv preprint arXiv:2103.00841*, 2021.
- [26] Zihan Xu, Mingbao Lin, Jianzhuang Liu, Jie Chen, Ling Shao, Yue Gao, Yonghong Tian, and Rongrong Ji. Recu: Reviving the dead weights in binary neural networks. *arXiv preprint arXiv:2103.12369*, 2021.
- [27] Ping Xue, Yang Lu, Jingfei Chang, Xing Wei, and Zhen Wei. Self-distribution binary neural networks. *arXiv preprint arXiv:2103.02394*, 2021.
- [28] Haoran You, Xiaohan Chen, Yongan Zhang, Chaojian Li, Sicheng Li, Zihao Liu, Zhangyang Wang, and Yingyan Lin. Shiftaddnet: A hardware-inspired deep network. *arXiv preprint arXiv:2010.12785*, 2020.
- [29] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6881–6890, 2021.
- [30] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.

A Computation Complexity of PECAN

In this section, we analyze the complexity of PECAN-A and PECAN-D during the inference stage, which is shown in Table 1. For both PECAN-A and PECAN-D, there are two stages of calculation: 1) computing the distance between the flattened features and prototypes, and 2) a simple table lookup to retrieve the product between weights and prototypes computed in advance.

Specifically, the first stage requires $H_{out}W_{out}$ subvectors in each group to compare with p prototypes. Therefore, PECAN-A needs $H_{out}W_{out}Dp \cdot d$ multiplications and additions, respectively, while PECAN-D has $H_{out}W_{out}Dp \cdot 2d$ additions. During the second stage, a lookup table of $c_{out} \times Dp$ is available to address the quantized product. It takes the weighted sum for PECAN-A or summation for PECAN-D in D groups. We can get $H_{out}W_{out}Dpc_{out}$ additions and multiplications for PECAN-A and $H_{out}W_{out}Dc_{out}$ for PECAN-D.

Since the FC layer can be regarded as a CONV layer when $k = H_{out} = W_{out} = 1$, the computation complexity of an FC layer can be obtained accordingly.

B Details of PECAN on MNIST

Table A1: Detail structure of LeNet used in PECAN.

LeNet	kernel size	Output	Flattened Weights
	$k \times k$	$[c_{out}, H_{out}, W_{out}]$	$[c_{out}, k^2 c_{in}]$
CONV1	3×3	[8, 26, 26]	[8, 9]
ReLU1	-	[8, 26, 26]	-
MaxPooling1	2×2	[8, 13, 13]	-
CONV2	3×3	[16, 11, 11]	[16, 72]
ReLU2	-	[16, 11, 11]	-
MaxPooling2	2×2	[16, 5, 5]	-
FC1	1×1	[128, 1, 1]	[128, 400]
ReLU3	-	[128, 1, 1]	-
FC2	1×1	[64, 1, 1]	[64, 128]
ReLU4	-	[64, 1, 1]	-
FC3	1×1	[10, 1, 1]	[10, 64]

Table A2: PECAN settings of LeNet on MNIST.

Layer	#Add.	#Mul.	p	D	d
CONV1	48.67K	48.67K	-	-	-
CONV1(PECAN-A)	45.97K	45.97K	4	1	9
CONV1(PECAN-D)	784.16K	0	64	1	9
CONV2	139.39K	139.39K	-	-	-
CONV2(PECAN-A)	116.16K	116.16K	8	3	24
CONV2(PECAN-D)	1.13M	0	64	8	9
FC1	51.2K	51.2K	-	-	-
FC1(PECAN-A)	28.8K	28.8K	8	25	16
FC1(PECAN-D)	57.60K	0	64	50	8
FC2	8.19K	8.19K	-	-	-
FC2(PECAN-A)	5.12K	5.12K	8	8	16
FC2(PECAN-D)	17.41K	0	64	16	8
FC3	0.64K	0.64K	-	-	-
FC3(PECAN-A)	0.83K	0.83K	8	4	16
FC3(PECAN-D)	8.27K	0	64	8	8

C Details of PECAN on CIFAR10

For PECAN, specialized settings are employed for different models. Table A3 describes detailed information for each layer in VGG-Small and ResNet20/32.

Table A3: The settings of prototype numbers and dimensions for each layer in different models for PECAN on CIFAR10.

Model	#Layers	Output map size	p/d (PECAN-A)	p/d (PECAN-D)
VGG-Small	2	32×32	16/9	32/3
	2	16×16	16/32	32/3
	2	8×8	16/32	32/3
	1	1×1	16/16	32/16
ResNet20	1	32×32	8/9	128/3
	6	32×32	8/9	64/3
	6	16×16	8/16	64/3
	6	8×8	8/16	64/3
	1	1×1	8/16	64/4
ResNet32	1	32×32	8/9	128/3
	10	32×32	8/9	64/3
	10	16×16	8/16	64/3
	10	8×8	8/16	64/3
	1	1×1	8/16	64/4

Table A4: Experiment results on TinyImageNet.

Model	Method	#Add.	#Mul.	Accuracy (%)
Modified ConvMixer	Baseline	3.36G	3.36G	56.76
	PECAN-A	2.36G	2.36G	59.42
	PECAN-D	0.98G	0	50.48

D Additional Experiments on Tiny-ImageNet

Due to space limitations, we put the experimental results on a larger dataset TinyImageNet [14] here. Different from the main paper, we choose ConvMixer [22] replacing all pointwise and depthwise convolution layers with conventional convolution layers. Besides, we keep the first convolution layer and the last fully-connected layer uncompressed. The depth of ConvMixer is 8 and kernel sizes in all blocks are $k = 5$. We set $p = 16$, $d = 25$ for PECAN-A and $p = 32$, $d = 25$ for PECAN-D. As shown in Table A4, PECAN-A achieves 59.42% accuracy which is much higher than the baseline when reducing around 1G multiplications and additions.

E Codes Instruction

In order to make it easier to verify the experimental results, we provide codes and the running commands for readers to reproduce the results in the tables.

For dataset MNIST, CIFAR-10 and CIFAR-100, we train our models using the following commands.

```
python train.py \
-log_dir [directory of the saved logs and models] \
-data_dir [directory to training data] \
-dataset [MNIST/CIFAR10/CIFAR100] \
-arch [resnet20_pecan_a/resnet20_pecan_d] \
-batch_size [training batch] \
-epochs [training epochs] \
-learning_rate [training learning rate] \
-lr_decay_step [learning rate decay step] \
-query_metric [dot/adder] \
-gpu [index of the GPU that will be used]
```