

# Deep Convolutional Tables: Deep Learning Without Convolutions

Shay Dekel, Yosi Keller<sup>ID</sup>, and Aharon Bar-Hillel<sup>ID</sup>

**Abstract**—We propose a novel formulation of deep networks that do not use dot-product neurons and rely on a hierarchy of voting tables instead, denoted as convolutional tables (CTs), to enable accelerated CPU-based inference. Convolutional layers are the most time-consuming bottleneck in contemporary deep learning techniques, severely limiting their use in the Internet of Things and CPU-based devices. The proposed CT performs a fern operation at each image location: it encodes the location environment into a binary index and uses the index to retrieve the desired local output from a table. The results of multiple tables are combined to derive the final output. The computational complexity of a CT transformation is independent of the patch (filter) size and grows gracefully with the number of channels, outperforming comparable convolutional layers. It is shown to have a better capacity:compute ratio than dot-product neurons, and that deep CT networks exhibit a universal approximation property similar to neural networks. As the transformation involves computing discrete indices, we derive a soft relaxation and gradient-based approach for training the CT hierarchy. Deep CT networks have been experimentally shown to have accuracy comparable to that of CNNs of similar architectures. In the low-compute regime, they enable an error:speed tradeoff superior to alternative efficient CNN architectures.

**Index Terms**—Convolutional tables (CTs), deep learning, efficient computation.

## I. INTRODUCTION

DEEP learning techniques in general, and convolutional neural networks (CNNs) in particular, have become the core computational approach in computer vision. Currently, “deep learning techniques” and “deep neural networks” are synonyms, despite relating to different concepts. By definition [9], deep learning is the “use of multiple layers to progressively extract higher level features from raw input,” i.e., learning a useful feature hierarchy. Artificial neural networks, on the other hand, are graphs of dot-products/convolutions computing elements. In this article, we argue that it is possible to unravel this dualism and that this is a worthwhile endeavor.

It is the need for accelerated and more efficient CPU-based inference that has motivated us to depart from the dot-product/convolution-based paradigm. CNNs consist of millions of “neurons”: computing elements applying dot products to their inputs. The resulting computational load is significant and is usually handled by GPUs and other specialized

hardware. Due to these computational demands, deep learning applications are often restricted to the server-side, where powerful dedicated hardware carries the computational load. Alternatively, applications running on standard CPUs or low-computing platforms such as cellphones are limited to small, low-accuracy networks, or avoid using networks altogether. Such applications are natural user interfaces (NUIs), face recognition, mobile robotics, and the Internet of Things (IoT) devices, to name a few. Low-compute platforms are not equipped with a GPU and are expected to handle the inference of simple deep learning tasks. Such devices are unable to train CNNs or run ImageNet-scale CNNs. Other devices of interest, such as most laptops and desktop computers, are unequipped with GPUs but have gigabytes of memory.

In this article, we propose a novel deep learning paradigm based on a hierarchy of deep voting tables. This allows the embedding of input signals such as images, similar to backbone CNNs, which is shown to be preferable in low-compute and GPU-less domains. Our approach uses CPUs that are capable of efficiently accessing large memory domains (random access). For instance, for a CPU with ideally unbounded memory access, the fastest classifier is a single huge table: a series of simple binary queries is applied to the input, an index is built from the resulting bits, and an answer is retrieved from the table entry. Such an approach is impractical due to the huge size of the resulting table. A practical alternative is to replace the single table with an ensemble of smaller hierarchical tables. Rather than relying on dot products, we propose a two-step, computationally efficient alternative. First, a set of  $K$  simple binary queries is applied to the input patch, thus encoding it with a binary codeword of  $K$ -bit. Second, the computed codeword is used as an index into a table to retrieve the corresponding output representation. The proposed transformation is called a convolutional table (CT). A CT layer consists of multiple CTs whose outputs are summed. Similar to convolutional layers, CT layers are stacked to form a CT network and derive a hierarchy of features. We detail this construction in Section III.

A fern is a tree in which all split criteria at a fixed tree depth are identical, such that the leaf identity is determined by using fixed  $K$  split criteria termed “bit functions.” The CT operation is the application of a single fern to a single input location, and the operation of a CT layer (sum of CTs) corresponds to applying a fern ensemble. In Section IV, we compare the computational complexity of the CNN and CT operations. By design, the complexity of a single CT layer is independent of the size of the patch used (“filter size” in a convolutional layer), and its dependence on quadratic depth

Manuscript received 1 April 2021; revised 17 October 2021, 15 July 2022, and 12 January 2023; accepted 22 April 2023. (All authors contributed equally to this work.) (Corresponding author: Yosi Keller.)

Shay Dekel and Yosi Keller are with Bar Ilan University, Ramat-Gan, Israel (e-mail: yosi.keller@gmail.com).

Aharon Bar-Hillel is with Ben-Gurion University, Negev, Israel.

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TNNLS.2023.3270402>.

Digital Object Identifier 10.1109/TNNLS.2023.3270402

92 terms is with significantly lower constants. This allows us to  
 93 derive conditions for which a CT network can be considerably  
 94 faster than a corresponding CNN.

95 In Section V, we derive rigorous results regarding the  
 96 capacity of a single fern and the expressiveness of a two-layer  
 97 CT network. Specifically, the VC dimension of a single fern  
 98 with  $K$  bits is shown to be  $\Theta(2^K)$ . This implies that a fern has  
 99 an advantage over a dot-product operation: it has a signifi-  
 100 cantly higher capacity-to-computing-effort ratio. Our results  
 101 show that a two-layer CT network has the same universal  
 102 approximation capabilities as a two-layer neural network and  
 103 that this also holds for ferns with a single bit function per fern.

104 A core challenge of the proposed deep-table approach is  
 105 the need for a proper training process. The challenge arises  
 106 because inference requires computing discrete indices that are  
 107 unnameable for gradient-based learning. Thus, we introduce  
 108 in Section VI a “soft formulation of the calculation of the  
 109  $K$  bit codeword,” as a Cartesian product of the  $K$  soft bit  
 110 functions. Following [23], [37], bit functions are computed  
 111 by comparing only two pixels in the region of interest.  
 112 However, in their soft version, the pixel indices parameters  
 113 are noninteger, and their gradient-based optimization relies on  
 114 the horizontal and vertical spatial gradients of the input maps.  
 115 During training, soft indices are used for table voting using a  
 116 weighted combination of possible votes. We use an annealing  
 117 mechanism and an incremental learning schedule to gradually  
 118 turn trainable but inefficient soft voting into efficient, hard  
 119 single-word voting at the inference time.

120 We developed a CPU-based implementation of the pro-  
 121 posed scheme and applied it to multiple standard datasets,  
 122 consisting of up to 250-K images. The results, presented  
 123 in Section VII, show that CT-networks achieve an accuracy  
 124 comparable to CNNs with similar structural parameters while  
 125 outperforming the corresponding CNNs trained using weights  
 126 and activations binarization. More importantly, we consider the  
 127 speed:error and speed:error:memory tradeoffs, where speed is  
 128 estimated via operation count. CT networks have been shown  
 129 to provide a better speed:error tradeoff than efficient CNN  
 130 architectures such as MobileNet [35] and ShuffleNet [29] in  
 131 the tested domain. The improved speed:error tradeoff requires  
 132 a higher memory consumption, as CT networks essentially  
 133 trade speed for memory. Overall, CT networks can provide  
 134 3.1–12× acceleration over CNNs with memory expansions  
 135 of 2.3×–10.6×, which are applicable to laptops and IoT  
 136 applications.

137 In summary, our contribution in this article is twofold: first,  
 138 we present an alternative to convolution-based networks and  
 139 show for the first time that useful deep learning (i.e., learning  
 140 of a feature hierarchy) can be achieved by other means. Sec-  
 141 ond, from an applicable viewpoint, we suggest a framework  
 142 enabling accelerated CPU inference for low-compute domains,  
 143 with moderate costs of accuracy degradation and memory  
 144 consumption.

## II. RELATED WORK

145 Computationally efficient CNN schemes were extensively  
 146 studied using a plethora of approaches, such as low-rank  
 147 or tensor decomposition [18], weight quantization [2], [20],

149 conditional computation [1], [3], [8], using FFT for effi-  
 150 cient convolution computations [32], or pruning of networks’  
 151 weights and filters [6], [13], [26], [40]. In the following,  
 152 we focus on the methods most related to our work.

153 Binarization schemes aim to optimize CNNs by binarizing  
 154 network weights [7], [28], [31] and internal network  
 155 activations [27], [34]. Recent approaches apply weights-  
 156 activations-based binarization [4], to reduce memory footprint  
 157 and computational complexity by utilizing only binary oper-  
 158 ations. Methods that binarize the activations are related to  
 159 the proposed scheme as they encode each local environment  
 160 into a binary vector. However, their encoding is dense, and  
 161 the output is computed using a (binary) dot product. The  
 162 proposed scheme also utilizes binary activations, but differs  
 163 significantly, as convolution/dot-product operations are not  
 164 applied. In contrast, we apply binary comparisons between  
 165 random samples in the activation maps.

166 Knowledge distillation approaches [15], [44] use a large  
 167 “Teacher” CNN to train a smaller computationally efficient  
 168 “Student” CNN. We show in our experiments that this tech-  
 169 nique can be combined with our approach to improve the accu-  
 170 racy. Regarding network design, several architectures were  
 171 suggested for low-compute platforms [11], [29], [35], [41],  
 172 utilizing group convolutions [35], efficient depthwise convolu-  
 173 tions [29], [35], map shifts [41], dubbed Bi-Real net [28],  
 174 and sparse convolutions [11]. Differentiable soft quantization  
 175 (DSQ) was proposed by Gong et al. [12] to quantize the  
 176 weights of a CNN to a given number of bits. An article by  
 177 Huang et al. [16] proposed a method for efficient quantiza-  
 178 tion using a multilevel binarization technique and a linear  
 179 or logarithmic quantizer to simultaneously binarize weights  
 180 and quantize activations to a low bit width. The authors  
 181 of this article then compare the tradeoffs of this method to  
 182 low-compute versions of MobileNetV2 [35], ShuffleNet [29],  
 183 and Bi-Real Net [28]. The use of trees or ferns applied using  
 184 convolutions for fast predictors was thoroughly studied in  
 185 computer vision [23], [24], [36], [37]. A convolutional random  
 186 forest enabled the estimation of human pose in real time in the  
 187 Kinect console [37]. A flat CT ensemble classifier for hand  
 188 poses, that was proposed in [23], enabled inference in less  
 189 than 1 ms on a CPU. It was extended [24] to a full-hand  
 190 pose estimation system. However, in all these works, flat  
 191 predictors were used, whereas we extend the notion to a deep  
 192 network with gradient-based training. Other works merged  
 193 tree ensembles with MLPs or CNNs for classification [17],  
 194 [19], or semantic segmentation [5]. Specifically, “conditional  
 195 networks” were presented in [17] as CNNs with a tree struc-  
 196 ture, enabling conditional computation to improve the speed-  
 197 accuracy ratio. Unlike our approach, where ferns replace the  
 198 convolutions, these networks use standard convolutions and  
 199 layers, and the tree/forest is a high-level routing mechanism.  
 200 Yang et al. [42] proposed a probability distribution approach  
 201 for quantization by treating the discrete weights in a quantized  
 202 neural network as searchable variables and using a differential  
 203 method to search for them accurately. They represent each  
 204 weight as a probability distribution over the set of discrete  
 205 values. Although this method was able to produce quantized  
 206 neural networks with higher performance than other state-of-

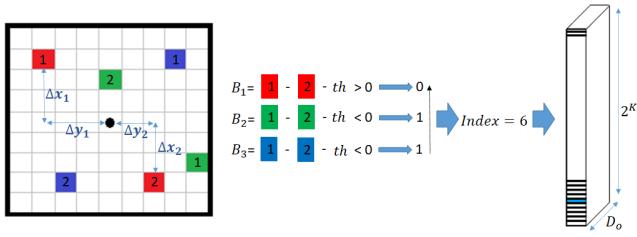


Fig. 1. CT with a word calculator of three bits operating on a 2-D input tensor. Each bit function is given by two pixels whose difference is compared to a threshold (colored in red, green, and blue). The word computed by the three bits, “110” = 6, is used to access the 6th entree in the voting table, thus computing an output vector  $\in \mathbb{R}^{D_o}$ .

the-art quantized methods in image classification, it is still reliant on standard convolution layers with dot products, which could be a bottleneck in low-power systems. In Section VII, we compare our method with that of the authors and report superior results in terms of accuracy. An effective channel-pruning method was suggested by Su et al. [38], using a bilaterally coupled network to determine the optimal width of a neural network, resulting in a more compact network. A bilaterally coupled network is a network that has been trained to approximate the optimal width of another network, to achieve a more compact network while maintaining or improving its performance. A new approach for dynamically removing redundant filters from a neural network was proposed by Tang et al. [39]. The channels are pruned by embedding the manifold information of all instances into the space of pruned networks. The performance of deep convolutional neural networks is improved by removing redundant filters while preserving the important features of the input data. The authors show that their method is an alternative to traditional channel-pruning methods, which may not be able to achieve optimal results in terms of both accuracy and efficiency. A scheme related to our approach was proposed by Zhou et al. [47], wherein the authors suggest a cascade of forests applied using convolutions. Each forest is trained to solve the classification task, and the class scores of a lower layer forest are used as features for the next cascade level. While the classifiers’ structure is deep, they are not trained end to end, but the trees are trained independently using conventional random forest optimization. Hence, the classifier uses thousands of trees and is only able to reach the accuracy of a two-layer CNN.

### III. CT TRANSFORM

A CT is a transformation accepting a representation tensor  $T^{\text{in}} \in \mathbb{R}^{H_i \times W_i \times D_i}$  and outputting a representation tensor  $T^{\text{out}} \in \mathbb{R}^{H_o \times W_o \times D_o}$ , where  $H_i, H_o$  and  $W_i, W_o$  are the spatial dimensions, whereas  $D_i, D_o$  are the representation’s dimensionality for the input and output tensors, respectively. Formally, it is the tuple  $(B, W)$ , where  $B$  is a word calculator, and  $W$  is a voting table. A word calculator  $B$  is a feature extractor applied to a patch, returning a  $K$ -bit index, such that  $B : \mathbb{R}^{l \times l \times D_i} \rightarrow \{0, 1\}^K$ , where  $l$  is the patch size. The voting table  $W \in M^{2^K \times D_o}$  contains all possible  $2^K$  outputs as rows, each being a vector  $\in \mathbb{R}^{D_o}$ .

The index is computed by sequentially applying  $K$  bit functions  $\{B^k\}_{k=1}^K$ , each computing a single bit of the index. Let  $P = (p_x, p_y) \in \mathbb{R}^{H_i \times W_i}$  be a location in  $T^{\text{in}}$ . We denote by  $B(P; \Theta) = (B^1(P; \Theta^1), \dots, B^K(P, \Theta^K))$  the index computed for the patch centered at location  $P$ . Here,  $\Theta = (\Theta^1, \dots, \Theta^K)$  are the word calculator parameters, with  $\Theta^k$  explicitly characterized below. The functions  $B^k(P; \Theta^k) \in \{0, 1\}$  (dependence on  $T^{\text{in}}$  is omitted for notation brevity) are computed by thresholding at 0 the simple smooth functions  $b^k(P; \Theta^k)$  of a few input pixels (i.e.,  $B^k(P; \Theta) = Q(b^k(P; \Theta^k))$  with  $Q(\cdot)$  being the Heaviside function). Bit functions  $b^k$  are the thresholded differences between two locations

$$\begin{aligned} b^k(P; \Theta_k) = & T^{\text{in}}(p_x + \Delta x_1^k, p_y + \Delta y_1^k, c^k) \\ & - T^{\text{in}}(p_x + \Delta x_2^k, p_y + \Delta y_2^k, c^k) - \text{th}^k \end{aligned} \quad (1)$$

where the learned parameters  $\Theta^k$  define the compared locations  $(\Delta x_1^k, \Delta y_1^k, c^k)$  and  $(\Delta x_2^k, \Delta y_2^k, c^k)$ , and the margin threshold  $\text{th}^k$ .  $\text{th}^k$  is the learnable scalar margin threshold for comparing two locations within the patch corresponding to the  $k$ th bit. The codeword computed by the word calculator is used as an index to the voting table to retrieve  $W(B(P; \Theta), :) \in \mathbb{R}^{D_o}$  for location  $P$ . Fig. 1 depicts the CT operation at a single location. A CT (using padding = “same” and stride = 1) applies the operation defined above at each position  $P \in \mathbb{R}^{H_i \times W_i}$  of the input tensor to obtain the corresponding output tensor column. The same as in convolutional layers, padding can be applied to the input for computations at border pixels, and strides may be used to reduce the spatial dimensions of the output tensor. A CT layer includes  $M$  CTs  $\{(B_m, W_m)\}_{m=1}^M$  whose outputs are summed

$$\begin{aligned} T^{\text{out}}[P, :] = & \sum_{m=1}^M W_m[B_m(T^{\text{in}}_{N(P)}), :] \\ B_m = & \{B^k\}_{k=1}^K. \end{aligned} \quad (2)$$

The operation of a CT layer is summarized in Algorithm 1. Essentially, it applies a fern ensemble at each location, summing the output vectors from the ensemble’s ferns. A CT network stacks multiple CT layers in a cascade or graph structure, similar to stacking convolutional layers in standard CNNs.

### IV. COMPUTATIONAL COMPLEXITY

Let the input and output tensors be  $T^{\text{in}} \in \mathbb{R}^{H_i \times W_i \times D_i}$  and  $T^{\text{out}} \in \mathbb{R}^{H_o \times W_o \times D_o}$ , respectively, where the dimensions are the same as in Section III, and let  $l \times l$  be the filter’s support. For a convolution layer, the number of operations per location is  $C_{\text{CNN}} = l^2 D_i D_o$ , since it computes  $D_o$  inner products at the cost of  $l^2 D_i$  each. For a CT layer with  $M$  CTs and  $K$  bits in each CT, the cost of bits computation is  $M K C_b$ , with  $C_b$  being the cost of computing a single bit function. The voting cost is  $M D_o$ , since we add the vectors  $\in \mathbb{R}^{MD_o}$  to get the result. Thus, the complexity of a CT layer is  $C_{\text{CT}} = M(C_b \cdot K + D_o)$ .

It follows that the computational complexity of a CT is independent of the filter size  $l$ . Hence, it allows agglomerating evidence from large spatial regions with no additional

**Algorithm 1** CT Layer

---

1: **INPUT:**  
A tensor  $T^{\text{in}} \in \mathbb{R}^{H_i \times W_i \times D_i}$ ,  
transformation parameters  $\{(\Theta^m, W^m)\}_{m=1}^M$

2: **OUTPUT:**  
An Output tensor  $T^{\text{out}} \in \mathbb{R}^{H_o \times W_o \times D_o}$

3: Initialization:  
 $T^{\text{out}} = 0$ ,  
 $\{(\Theta^m, W^m)\}_{m=1}^M \sim \mathcal{N}(0, \sigma^2 = 1)$

4: **for**  $P \in \{1, \dots, H_i\} \times \{1, \dots, W_i\}$  **do** ▷ All pixel locations

5:   **for**  $m = 1, \dots, M$  **do** ▷ All CT tables

6:     Compute  $B_m(P; \Theta^{k,m}) \in \{0, 1\}^K$  ▷ Eq. 1

7:      $T^{\text{out}}[P, :] = T^{\text{out}}[P, :] + W^m[B_m(P; \Theta^{k,m}), :]$  ▷ Eq. 2

8:   **end for**

9: **end for**

10: Return  $T^{\text{out}}$

---

computational cost. To compare CNN and CT complexities, assume that both use the same representation dimension  $D_i = D_o = D$ . The relationship between the total number of bit functions  $MK$  and  $D$  depends on whether each bit function uses a separate dimension, or the dimensions are reused by multiple bit functions. For a general analysis, assume that each dimension is reused by  $R$  bit functions, such that  $D = MK/R$  (in our experiments  $R \in [1, 3]$ ). Thus, the ratio of complexities between CNNs and CT layers with the same representation dimension  $D$  is

$$\frac{C_{\text{CNN}}}{C_{\text{CT}}} = \frac{l^2 D^2}{M(C_b \cdot K + D)} = \frac{l^2 D^2}{C_b DR + \frac{D^2 R}{K}} \approx \frac{K l^2}{R} \quad (3)$$

where the last approximation is true for large  $D$  and the quadratic terms in  $D$  are dominant. The computational cost of computing a single bit function,  $C_b$ , is defined as the number of operations required to perform a single comparison operation between two locations within the patch size. This computation is described in (1) and involves four additions of the offsets of  $\Delta_x$ ,  $\Delta_y$ , and a subtraction of the threshold,  $\text{th}^k$ . This results in a total of five operations. In addition, load operations are required for the two relevant memory locations and their two base addresses, bringing the total number of operations to nine. Although the most stringent count does not exceed ten operations, for simplicity, we set  $C_b$  to 10. The ratio of complexities between CNNs and CT layers in (3) suggests that there is the potential for acceleration by more than an order of magnitude for reasonable choices of parameter values. For example, using  $l = 3$ ,  $K = 8$ , and  $R = 2$  results in a  $36 \times$  acceleration. Although the result of this analysis is promising, several lower level considerations are essential to achieve actual acceleration. First, operation counts translate into actual improved speed only if the operations can be efficiently parallelized and vectorized. In particular, vectorization requires contiguous memory access in most computations. The CT transformation adheres to these constraints. Contrary to trees, bit computations in ferns can be vectorized over adjacent input locations as all locations use the same bit functions, and the memory access is contiguous. The voting operation can

be vectorized over the output dimensions. This implementation was already built and has been shown to be highly efficient for flat CT classifiers [24]. Although a CT transformation uses a “random memory access” operation when the table is accessed, this random access is rare: there is a single random access operation included in the  $(C_b \cdot K + D_o)$  operations required for a single-fern computation. Another important issue is the need to keep the model’s total storage size within reasonable bounds to enable its memory to be efficiently accessed in a typical  $L_2$  cache. For example, a model consisting of 50 layers, wherein  $M = 10$  ferns,  $K = 6$ , and  $D = 60$ , has a total size of 1.92 MB when the parameters are stored in 8-bit precision.

## V. CT CAPACITY AND EXPRESSIVENESS

In this section, we present two analytic insights that indicate the feasibility of deep CT networks and their potential advantages. In Section V-A, the capacity of a fern is shown to be  $\Omega(2^K)$  with  $O(K)$  computational effort, providing a much better capacity:compute ratio than a linear convolution. In Section V-B, a network of two fern layers is shown to be a universal approximation, similar to known results on neural networks. In the analysis, we consider a simplified nonconvolutional setting, with an input vector  $X \in [0, 1]^d$ , and a nonconvolutional fern ensemble classifier. Furthermore, simple bit functions of the form  $B^k(X) = Q((-1)^{s^k}(X[I^k] - t^k))$  are considered, where  $s^k \in \{0, 1\}$ ,  $I^k \in 1, \dots, d$ , and  $t^k \in [0, 1]$ . These simple bit functions just compare a single entry to a threshold.

## A. Capacity:Compute Ratio

The following lemma characterizes the capacity of a single fern in VC-dimension terms.

*Lemma 1:* Define the hypothesis family of binary classifiers based on a single  $K$ -bit fern  $H = \{f(X) = \text{sign}(W(B(X))) : B = \{(s^k, I^k, t^k)\}_{k=1}^K, W \in \mathbb{R}^{2^K}\}$ . For  $K \geq 4$  and  $\log_2(d) < K \leq d$ , its VC dimension is

$$2^K \leq \text{VC-dim}(H) \leq K 2^K. \quad (4)$$

*Proof:* Since we only consider the sign of the chosen table entry  $W(B(X))$ , we may equivalently consider the family of classifiers  $f(X) = W(B(X))$  for  $W \in \{-1, 1\}^{2^K}$ .

1) *Lower Bound:* Consider a sample containing points on the  $d$ -dimensional cube

$$S = \{P_c = (-1^{c_1}, \dots, -1^{c_K}) : c = (c_1, \dots, c_K) \in \{0, 1\}^K\}. \quad (5)$$

It follows that this sample can be shattered by the binary fern family. We may choose the bit functions  $b^k(X) = Q(X[k])$ . Assume an arbitrary label assignment  $l$ , i.e., for each example  $P_c$ , we have  $l(P_c) = y_c$  with  $y_c$  arbitrarily chosen in  $\{-1, 1\}$ . By definition, for each  $c$  and  $k$ , we have  $b^k(P_c) = \{-1\}^{c_k} = P_c[k]$ . Since any two points  $P_{c_1}, P_{c_2}$  differ by at least a single dimension, they will have at least different bits in a single bit function. Hence, the  $2^K$  points are mapped into the  $2^K$  different cells of  $W$ . By choosing  $W[c] = y_c$ , we may get the desired labeling  $l$ . Therefore,  $S$  is shattered, showing  $\text{VCdim}(H) \geq 2^K$ .

391    2) *Upper Bound:* Assume that  $K \geq 4$  and  $\log_2(d) < K \leq$   
 392     $d$ . We compute an upper bound on the number of possible  
 393    labels (label vectors) enabled by the binary fern family on  
 394    a sample of size  $n$  and show that it is smaller than  $2^n$  for  
 395     $n = K2^K$ .

396    For a  $K$ -bit fern, there are  $d^K$  ways to choose the input  
 397    dimension of the  $K$  bit functions (one dimension per function).  
 398    For each bit function, once the input dimension  $I$  is chosen,  
 399    we may reorder the examples according to their value in  
 400    dimension  $I$ , i.e.,  $X_1[I] \leq X_2[I], \dots, \leq X_n[I]$ . Nontrivial  
 401    partitions of the sample, i.e., partitions into two nonempty sets,  
 402    are introduced by choosing the threshold  $t$  in  $(X_1[I], X_n[I])$ ,  
 403    and there are at most  $n - 1$  different options. Including the  
 404    trivial option of partitioning  $S$  into  $S$  and  $\Phi$ , there are  $n$   
 405    options.

406    Consider the number of partitions of the examples into  $2^K$   
 407    cells, which are made possible using  $K$  bit functions. Two  
 408    different partitions must differ w.r.t. the induced partition by  
 409    a single bit function choice, at least. Hence, their number is  
 410    bounded by the number of ways to choose the bit function  
 411    partitions, which is  $d^K n^K$  at most, according to the above  
 412    considerations. For each such partition into  $2^K$  cells,  $2^K$   
 413    different classifiers can be defined by choosing the table entries  
 414     $\{W[c]\}_{c=0}^{2^K-1}$ . Hence, the total number of possible different fern  
 415    classifiers for  $n$  points is bounded by  $(dn)^K 2^{2^K}$ . When this  
 416    bound is lower than  $2^n$ , the sample cannot be shattered. Thus,  
 417    we solve the problem

$$(dn)^K 2^{2^K} < 2^n \quad (6)$$

419    and

$$2^K + K \log_2 d + K \log_2 n < n. \quad (7)$$

421    By choosing  $n = K2^K$ , it follows that (7) is satisfied since

$$\begin{aligned} 422 \quad 2^K + K \log_2 d + K \log_2 K + K^2 \\ 423 \quad < 2^K + 3K^2 \leq 2^K + (K-1)2^K \leq K2^K. \end{aligned} \quad (8)$$

424    For the second inequality, we used  $\log_2 d < K$  (assumed),  
 425    and for the third, we used the inequality  $3K^2 \leq (K-1)2^K$ ,  
 426    which holds for  $K \geq 4$ . ■

427    Results similar to Lemma 1 were derived for trees [30], [45],  
 428    and the lemma implies that the capacity of ferns is not signifi-  
 429    cantly lower than that of trees. To understand its implications,  
 430    compare such a single-fern classifier to a classifier based on a  
 431    single dot-product neuron, i.e., a linear classifier. A binary dot-  
 432    product neuron with  $d = K$  input performs  $O(K)$  operations  
 433    to obtain a VC dimension of  $O(K)$ . The fern also performs  
 434     $O(K)$  computations, but the response is chosen from a table  
 435    of  $2^K$  leaves, and the VC dimension is  $\Omega(2^K)$ . The higher  
 436    (capacity):(computational-complexity) ratio of ferns indicates  
 437    that they can compute significantly more complex functions  
 438    than linear neurons using the same computational budget.

### 439    B. CT Networks Are Universal Approximators

440    The following theorem states the universal approximation  
 441    capabilities of a two-layer CT network. It resembles the known  
 442    result for two-stage neural networks [14].

Theorem 1: Any continuous function  $f : [0, 1]^d \rightarrow [0, 1]$   
 can be approximated in  $L_\infty$  using a two-layer fern network  
 with  $K = 1$  in all ferns.

*Proof:* The core of the proof is to show that by using  $2K$   
 1-bit ferns at the first layer, a second layer fern can create  
 a function with an arbitrary value inside a hyperrectangle of  
 choice, and zero outside the rectangle. Since sums of such  
 “step functions” are dense in  $C[0, 1]$ , so are two-layer fern  
 networks.

Let  $R = [a_1, b_1] \times [a_2, b_2], \dots, \times [a_d, b_d]$  be a hyperrect-  
 angle in  $R^d$ , and let  $v \in R$  be a scalar value. A rectangle  
 function  $s(x; v, R) : [0, 1]^d \rightarrow R$  is defined as  $v \cdot 1_{x \in R}$ , i.e.,  
 a function whose value is  $f(x) = v$  for  $x \in R$  and 0 otherwise.  
 Define the family of step functions

$$G = \{g : [0, 1]^d \rightarrow [0, 1] : g(x) = \sum_{p=1}^P s(x; v_p, R_p)\}. \quad (9)$$

It is known that  $G$  is dense in  $(C[0, 1]^d)$  (from the Stone-  
 Weierstrass theorem; see, e.g., [10]). We will show that the  
 family of two-layer CT networks includes this set.

Let  $s(v, R)$  be an arbitrary rectangle function. For each  
 dimension  $i = 1, \dots, d$ , define the following two bit func-  
 tions:  $L_i(x) = Q(x - a_k)$  and  $R_i(x) = Q(-(x - b_k))$ . Denote  
 these bit functions by  $B_{(R)}^j$  for  $j = 1, \dots, 2d$ . By construction,  
 $\{B_R^j\}_{j=1}^{2d}$  characterize the rectangle

$$x \in R \iff \forall j = 1, \dots, 2d, \quad B_R^j(x) = 1. \quad (10)$$

Equivalently, we have  $x \in R$  iff  $\sum_{j=1}^{2d} B_R^j(x) > 2d - (1/2)$ .  
 Given a function  $g(x) = \sum_{p=1}^P s(x; v_p, R_p)$  to implement,  
 we define at the first layer  $P$  sets of ferns. Fern set  $p$  contains  
 $2d$  ferns, denoted by  $F_p^j$ , with a single bit function each, and  
 the bit function of  $F_p^j$  is  $B_{R_p}^j$ . The output dimension of the first  
 layer is  $P$ . The weight table  $W_p^j \in M_{2 \times P}$  of  $F_p^j$  is defined by

$$W_p^j[i, l] = \begin{cases} 1, & i = 1, l = p, \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

The output vector of a layer is the sum of all ferns  $Y =$   
 $\sum_{p'=1}^P \sum_{j=1}^{2d} W_p^j[B_{R_p}^j(x), :]$ . By construction, we have

$$Y[p] = \sum_{p'=1}^P \sum_{j=1}^{2d} W_p^j[B_{R_p}^j(x), p] \quad (12)$$

$$= \sum_{j=1}^{2d} W_p^j[B_{R_p}^j(x), p] = \sum_{j=1}^{2d} B_{R_p}^j(x).$$

The second equality holds since  $W_p^j[i, p]$  only differs from  
 zero for  $p = p'$ . The last equality is valid since  $W_p^j[i, p] = i$   
 for  $i = 0, 1$ . We, hence, have  $Y[p] = 2d > 2d - (1/2)$  for  
 $x \in R_p$ , which has lower values anywhere else.

The second layer of ferns contains  $P$  ferns, each with a  
 single bit and a single output dimension. For fern  $p$ , we define  
 the bit function  $B_p^{(2)}(Y) = Q(Y[p] - (2d - (1/2)))$  which fires  
 only for  $x \in R_p$  and the table  $W_p^{(2)} = [0, v(p)]$ . The output

486 unit  $U$  computes

$$\begin{aligned} 487 \quad U &= \sum_{p=1}^P W_p^{(2)}[B_p^{(2)}(x)] \\ 488 \quad &= \sum_{p=1}^P v(p) \cdot B_p^{(2)}(x) = \sum_{p=1}^P v(p) \cdot 1_{x \in R_p}. \end{aligned} \quad (13)$$

489 Therefore,  $U$  implements the function  $g(x)$ , which completes the proof. ■

## VI. TRAINING WITH SOFT CTS

492 Following (1),  $B^k(P; \Theta_k) = Q(b^k(P; \Theta_k))$  does not enable  
493 gradient-based optimization, as the derivative of the Heaviside  
494 function  $Q(\cdot)$  is zero almost anywhere. We suggest a soft  
495 version of the CT in Section VI-A that enables gradient-based  
496 optimization and discusses its gradient in Section VI-B.

### A. Soft CT Version

498 To enable gradient-based learning, we suggest replacing the  
499 Heaviside function during optimization with a linear sigmoid  
500  $q(x; t)$  for  $t > 0$

$$501 \quad q(x; t) = \min(\max((t+x)/2t, 0), 1) \quad (14)$$

502 where  $q(x; t)$  is identical to the Heaviside function for  
503  $x \gg 0$  and is a linear function in the vicinity of  $x = 0$ .  
504 It has a nonzero gradient in  $[-t, t]$ , where  $t$  is a hyperparameter  
505 controlling its smoothness. For low  $t$  values, we have  
506  $q(x; t) \xrightarrow[t \rightarrow 0]{} Q(x)$ . Moreover, we have  $q(x; t) + q(-x; t) = 1$ .  
507 Thus, it can be interpreted as a pseudoprobability, with  $q(x)$   
508 estimating the probability of a bit being 1 and  $q(-x)$  the  
509 probability of being 0.

510 Following Section III, a word calculator  $B$  maps a patch  $P$   
511 to a single index or equivalently to a one-hot vector (row) in  
512  $R^{2^K}$ . Denote the word calculator in the latter view (mapping  
513 into  $R^{2^K}$ ) by  $\vec{B}(P; \Theta)$ . Given this notation, the CT output  
514 is given by the product  $\vec{B} \cdot W$ . Hence, extending the word  
515 calculator to a soft function in  $R^{2^K}$ , with dimension  $b \in$   
516  $\{0, \dots, 2^K - 1\}$  measuring the activity of the word  $b$ , provides  
517 a natural differentiable extension of the CT formulation. For  
518 an index  $b$ , define the sign of its  $k$ th bit by  $s(b, k) \triangleq$   
519  $(-1)^{(1+u(b,k))}$ , where  $u(b, k)$  denotes the  $k$ th bit of the index  
520  $b$  in its standard binary expansion, such that

$$521 \quad s(b, k) = \begin{cases} 1, & \text{bit } k = 1 \\ -1, & \text{bit } k = 0. \end{cases} \quad (15)$$

522 The activity level of the word  $b$  in a soft-word calculator is  
523 defined by

$$524 \quad \vec{B}^s(P; \Theta)[b] = \prod_{k=1}^K q(s(b, k) \cdot b_k(P; \Theta_k); t). \quad (16)$$

525 Intuitively, the activity level of each possible codeword  
526 index  $b$  is a product of the activity levels of the individual  
527 bit functions in the directions implied by the codeword bits.  
528 Marginalization implies that  $\sum_{b=0}^{2^K-1} \vec{B}^s(P)[b] = 1$ , and hence  
529 the soft-word calculator defines a probability distribution over

530 the possible  $K$ -bit words. Soft CT is defined as a natural  
531 extension  $\vec{B}^s W$ . Note that as  $t \rightarrow 0$ , the sigmoid  $q(x; t)$   
532 becomes sharp,  $\vec{B}^s$  becomes a one-hot vector, and the soft  
533 CT becomes a standard CT as defined in Section III, without  
534 any inner products.

535 Soft CT can be considered as a consecutive application of  
536 two layers: a soft indexing layer  $\vec{B}^s$  followed by a plain linear  
537 layer (though sparse if most of the entries in  $\vec{B}^s$  are zero).  
538 Since CT is applied to all spatial locations on the activation  
539 map, the linear layer  $W$  corresponds to a  $1 \times 1$  convolution.  
540 We implemented a sparse  $1 \times 1$  convolutional layer operating  
541 on a sparse tensor  $\mathbf{x} \in R^{H \times W \times 2^K}$  and outputting a dense tensor  
542  $\mathbf{y} \in R^{H \times W \times D_o}$ .

### B. Gradient Computation and Training

543 Since applying the soft CT in a single position is  $W \cdot$   
544  $B^s(P; \Theta)$ , with  $W$  being a linear layer, it suffices to consider  
545 the gradient of  $\vec{B}^s(P; \Theta)$ . Denote the neighborhood  $l \times l$  of  
546 the location  $P$  in  $T^{\text{in}}$  by  $T_{N(P)}^{\text{in}}$ . Considering a single output  
547 variable  $b$  at a time, the gradient with respect to the input  
548 patch  $T_{N(P)}^{\text{in}}$  is given by

$$549 \quad \frac{\partial \vec{B}^s(P)[b]}{\partial T_{N(P)}^{\text{in}}} = \sum_k \frac{\partial \vec{B}^s(P)[b]}{\partial b_k(P; \Theta_k)} \cdot \frac{\partial b_k(P; \Theta_k)}{\partial T_{N(P)}^{\text{in}}} \quad (17)$$

550 with a similar expression for the gradient with respect to  $\Theta$ .  
551 For a fixed index  $k$ , the derivative  $(\partial \vec{B}^s(P)[b]) / (\partial b_k(P; \Theta_k))$   
552 is given by

$$\begin{aligned} 553 \quad &\frac{\partial \vec{B}^s(P)[b]}{\partial b_k(P; \Theta_k)} \\ 554 \quad &= \prod_{j \neq k} q(s(b, j) b_j(P; \Theta_j); t) \frac{\partial q(s(b, j) b_j(P; \Theta_j); t)}{\partial b_k(P; \Theta_k)} \\ 555 \quad &= \frac{\vec{B}^s(P)[b] \cdot t^{-1} \cdot 1_{|b_k| < t} \cdot s(b, k)}{2q(s(b, k) \cdot b_k(P; \Theta_k))}. \end{aligned} \quad (18)$$

556 The derivative is nonzero when the word  $b$  is active (i.e.,  
557  $\vec{B}^s(P)[b] > 0$ ), and the bit function  $b_k(P)$  is in the dynamic  
558 range  $-t < b_k(P) < t$ .

559 To further compute the derivatives of  $b_k(P)$  with respect  
560 to the input and parameters  $[((\partial b_k(P; \Theta_k)) / dT_{N(P)}^{\text{in}})$   
561 and  $((\partial b_k(P; \Theta_k)) / d\Theta)$ , respectively), we first note that the  
562 forward computation of  $b_k(P)$  estimates the input tensor at  
563 fractional spatial coordinates. This is because the offset parameters  
564  $\Delta x_1, \Delta y_1, \Delta x_2, \Delta y_2$  are trained with gradient descent, which  
565 requires them to be continuous. We use bilinear interpolation  
566 to compute image values at fractional coordinates in  
567 the forward inference. Hence, such bilinear interpolation (of  
568 spatial gradient maps) is required for gradient computation.  
569 For instance,  $((\partial b_k(P)) / \partial \Delta x_1^k)$  is given by

$$\begin{aligned} 570 \quad \frac{\partial b_k(P)}{\partial \Delta x_1^k} &= \frac{\partial T^{\text{in}}(p_x + \Delta x_1^k, p_y + \Delta y_1^k, c)}{\partial \Delta x_1^k} \\ 571 \quad &= \frac{\partial T^{\text{in}}(:, :, c)}{\partial x} [(p_x + \Delta x_1^k, p_y + \Delta y_1^k)] \end{aligned} \quad (19)$$

572 where  $((\partial T^{\text{in}}(:, :, c)) / \partial x)$  is the partial  $x$ -derivative of the  
573 image channel  $T^{\text{in}}(:, :, c)$ , which is estimated numerically and  
574 sampled at  $(x, y) = (p_x + \Delta x_1^k, p_y + \Delta y_1^k)$  to compute (19).

576 The derivatives with respect to other parameters are computed  
 577 in a similar way. Note that the channel index parameters  $c^k$  are  
 578 not learned—these are fixed during network construction for  
 579 each bit function. As for  $((\partial b_k(P; \Theta_k))/dT_{N(P)}^{\text{in}})$ , by consider-  
 580 ing (1), it follows that  $b_k(P; \Theta_k)$  only relates to two fractional  
 581 pixel locations:  $(p_x + \Delta x_1^k, p_y + \Delta y_1^k)$  and  $(p_x + \Delta x_2^k, p_y +$   
 582  $\Delta y_2^k)$ . This implies that a derivative with respect to eight pixels  
 583 is required, as each value of a fractional coordinate's pixel  
 584 is bilinearly interpolated using its four neighboring pixels in  
 585 integer coordinates.

586 The sparsity of the soft-word calculator (both forward and  
 587 backward) is governed by the parameter  $t$  of the sigmoid  
 588 in (14), which acts as a threshold. When  $t$  is large, most of the  
 589 bit functions are ambiguous and soft, i.e., not strictly 0 or 1,  
 590 and the output will be dense. We can control the output's  
 591 sparsity level by adjusting  $t$ . In particular, as  $t \rightarrow 0$ , the  
 592 bit functions become hard, and  $B^s(P)$  converges toward a  
 593 hard fern with a single active output word. While a dense  
 594 flow of information and gradients is required during training,  
 595 fast inference in test time requires a sparse output. Therefore,  
 596 we introduce an annealing scheduling scheme, such that  $t$  is  
 597 initiated as  $t \gg 0$ , set to allow a fraction  $f$  of the values of  
 598 the bit functions to be in the “soft zone”  $[-t, t]$ . The value  
 599 of  $t$  is then gradually lowered to achieve a sharp and sparse  
 600 classifier toward the end of the training phase.

## 601 VII. EXPERIMENTAL RESULTS

602 We experimentally examined and verified the accuracy and  
 603 computational complexity of the proposed CT scheme, applying  
 604 it to multiple small and medium image datasets, consisting  
 605 of up to 250-K images. Such datasets are applicable to IoT and  
 606 non-GPU low-compute devices. We explored the sensitivity of  
 607 CT networks to their main hyperparameters in Section VII-A.  
 608 To this end, we studied the number of ferns  $M$  and the number  
 609 of bit functions per fern  $K$ , and the benefits of distillation  
 610 as a CT training technique. We then compared our CT-based  
 611 architectures to similar CNN classes. Such a comparison is  
 612 nontrivial due to the inherent architectural differences between  
 613 the network classes. Thus, in Section VII-B, we compared  
 614 Deep CT to CNNs with similar depth and width, that is, the  
 615 number of layers and maps per layer. The comparison focuses  
 616 on efficient methods suggested for network quantization and  
 617 binarization, since the CT framework uses similar notions.  
 618 In Section VII-C, the CT is compared with CNN formulations  
 619 with comparable computational complexity, i.e., similar num-  
 620 ber of multiply–accumulate (MAC) operations. We compared  
 621 CT with CNN architectures specifically designed for inference  
 622 efficiency, such as MobileNet [35] and ShuffleNet [29]. The  
 623 comparison was applied by studying the speed:error and  
 624 speed:error:memory tradeoffs of the different frameworks.  
 625 Finally, to exemplify the applicability of Deep CT to low-  
 626 power IoT applications and larger datasets (250-K images),  
 627 it was applied to IoT-based face recognition in Section VII-D.

628 The proposed novel CT layer does not use standard CNN  
 629 components such as activations, convolutions, or FCs. Thus,  
 630 the Deep CT word calculator layer and the sparse voting  
 631 table layer were implemented from scratch, in unoptimized C,

632 compiled as MEX files, and applied within a MATLAB  
 633 implementation. Our experiments were carried out on the same  
 634 computer using CPU-only implementations for all schemes.

635 Similar to other inherently novel layers, such as Spatial  
 636 Transformer Networks, it is challenging to implement a con-  
 637 verging end-to-end training process. To facilitate the conver-  
 638 gence of the proposed scheme, a three-phase training scheme  
 639 was applied. First, the lower half of the CNN (layers #1–#6  
 640 in Table VIII) was trained, by adding the average pooling  
 641 (AP), Softmax (SM), and cross-entropy (CE) layers on top  
 642 of layer #6. Then, the full architecture was trained with  
 643 layers #1–#6 frozen. Finally, following Hinton et al. [15],  
 644 we unfroze all layers and refined the entire network using  
 645 a distillation-based approach. The CT was trained to optimize  
 646 a convex combination of CE loss and Kullback–Leibler diver-  
 647 gence with respect to the probabilities induced by a teacher  
 648 CNN. The temperature parameter was initially set at  $t = 4$  and  
 649 gradually annealed to  $t = 1$  toward the end of the training.  
 650 All Deep CT networks were trained from scratch using random  
 651 Gaussian initialization of the parameters.

### 652 A. Hyperparameters Sensitivity

653 The expressive power of the Deep CT transformation is  
 654 related to the number of CTs per layer  $M$ , and the number  
 655 of bit functions  $K$  used in each CT. We tested the sen-  
 656 sitivity of the Deep CT network to these parameters and  
 657 the suggested annealing schedule. These experiments applied  
 658 CT networks with four layers to the SVHN, CIFAR-10, and  
 659 CIFAR-100 datasets, and two layers to the MNIST dataset.  
 660 The networks are similar to those shown in the Appendix,  
 661 in Tables VI and VII. We used a fixed patch size  $l = 5$  and  
 662 input and output dimensions of  $D_i = D_o = 32$ , respectively.  
 663 It is best to apply CT networks with a patch size larger than  
 664 the kernel size of CNN-based networks, since CT patches  
 665 are sparsely represented. Comparatively, CNN complexity is  
 666 quadratic with respect to patch size, whereas CT networks' 667 computational complexity is independent of patch size.

668 The classification errors with respect to the bit functions  $K$   
 669 (for a fixed  $M$ ), and the number of CTs  $M$  (for a fixed  $K$ ), are  
 670 shown in Fig. 2(a) and (b), respectively. For most datasets,  
 671 there is a range in which the error exponentially decreases  
 672 as the number of network parameters increases, whereas the  
 673 accuracy improvement saturates for  $M > 8$  and  $K > 8$  for the  
 674 deeper networks (datasets other than MNIST).

675 Fig. 2(c) shows the average number of active words in a CT  
 676 as a function of the training time for CIFAR-10. The fraction  $f$   
 677 of ambiguous bits (bits whose value is not set to either 0 or 1)  
 678 was started at 0.2 and was gradually reduced exponentially,  
 679 whereas the sigmoid threshold  $t$  was adjusted periodically  
 680 according to the required  $f$ . In our experiments, the sigmoid  
 681 threshold was set to  $t = 0.5$ . The number of active words  
 682 per CT is thus gradually reduced and converged to a single  
 683 active word in the final classifier. The CNN with four layers  
 684 obtained a test error of 12.55%, compared with 12.49% with  
 685 the same architecture without annealing. Therefore, annealing-  
 686 based training led to a sparser classifier without significantly  
 687 degrading the classification accuracy.

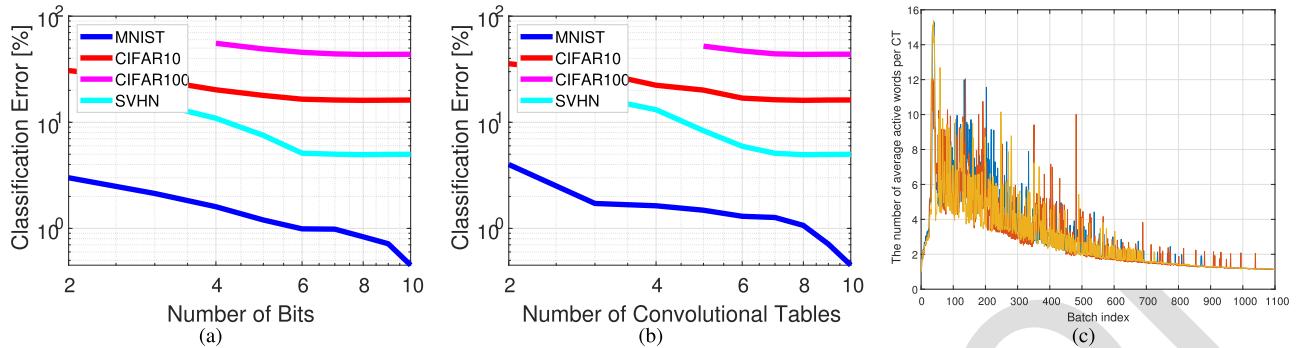


Fig. 2. Hyperparameters study on MNIST, CIFAR10, CIFAR100, and SVHN datasets. (a) Deep CT network error versus the number of bits  $K$  in a log–log plot. In all layers, the number of CTs  $M$  was fixed to 10. (b) Deep CT network error versus the number of CTs  $M$  in a log–log plot, for a fixed number of bits  $K = 10$ . (c) Average number of active words per CT was evaluated using the CIFAR-10 dataset and plotted against the training batch index. Each convolution table is represented by a different color: yellow, blue, and red. As expected, the number of active words decreases during the training process due to the annealing mechanism and eventually converges to a single active word.

TABLE I  
JOINT HYPERPARAMETERS TUNING FOR THE FOUR- AND SIX-LAYER CT NETWORKS DESCRIBED IN TABLES VII AND II.  $M$  AND  $K$  DENOTE THE CONVOLUTION TABLES' SIZE AND THE NUMBER OF BITS, RESPECTIVELY

| K        | M        | Layers# | Distil | Error [%]    |
|----------|----------|---------|--------|--------------|
| 9        | 8        | 6       | No     | 11.31        |
| 8        | 9        | 6       | No     | 12.31        |
| <b>8</b> | <b>8</b> | 6       | No     | <b>11.05</b> |
| 7        | 8        | 6       | No     | 13.93        |
| 8        | 7        | 6       | No     | 13.34        |
| <b>8</b> | <b>8</b> | 6       | Yes    | <b>10.25</b> |
| 9        | 8        | 4       | No     | 12.88        |
| 8        | 9        | 4       | No     | 13.05        |
| <b>8</b> | <b>8</b> | 4       | No     | <b>12.49</b> |
| 7        | 8        | 4       | No     | 13.25        |
| 8        | 7        | 4       | No     | 12.90        |
| <b>8</b> | <b>8</b> | 4       | Yes    | <b>10.91</b> |

To establish the best hyperparameters, additional experiments were carried out for multiple  $M$  and  $K$  combinations around  $M = 8$ ,  $K = 8$ . The experiments were conducted using the CIFAR-10 dataset, and the CT networks having four and six layers, detailed in Tables II and VII, respectively. In addition, we examined training with and without distillation from a teacher. The  $M = 8$ ,  $K = 8$  configuration was found to be the most accurate. Distillation was found to be beneficial, providing accuracy improvements of 1% – 2%.

### B. Accuracy Comparison of Low-Complexity Architectures

Although deep CT and CNN networks use very different computational elements, both use a series of intermediate representations to perform classification tasks. Hence, we consider CT and CNN backbone networks similar if they have the same number of intermediate representations (depth) and the same number of maps in each representation (width). In this section, we report experiments comparing CNN and CT networks that have identical configurations in these respects. Although the intermediate spatial tensors were chosen to be similar, we allowed for different top inference layers. Those were chosen to optimize the classification accuracy for each network class. Thus, the CNN-based networks used

max-pooling and an FC layer, whereas the CT only used an average pooling layer.

We compared the classification accuracy of CT networks with standard CNNs and contemporary low-complexity CNN architectures based on CNN binarization schemes. The methods we compared against include the BinaryConnect (BC) [7],<sup>1</sup> XNOR-Net (XN)<sup>2</sup>, XNOR-Net++ (XN++) [4], Bi-Real Net (Bi-Real) [28], SLB [42]<sup>3</sup>, and TABCNN [27],<sup>4</sup> using their publicly available implementations. These methods, [34] and [27] in particular, also encode the input activity using a set of binary variables. In contrast to our CT table-based approach, they use dense encodings and apply a binary dot product to their encoded input. Finally, we applied the state-of-the-art DSQ scheme [12]<sup>5</sup> to the conventional CNN to quantize it to 2, 4, and 6 bits.

All networks were applied to the MNIST [25], CIFAR-10 [21], CIFAR-100 [22], and SVHN [33] datasets. For the MNIST dataset, we applied CT networks and CNNs consisting of two layers, as this relatively small dataset did not require additional depth. For the other datasets, networks consisting of four and six layers were applied.

The architectures of the six-layer networks are detailed in Table II, and the architectures of the two and four layers used are detailed in the Appendix in Tables VI and VII, respectively. The classification accuracy results are reported in Table III. On the left-hand side of the table, all methods are compared using the cross-entropy loss. The proposed CT networks outperformed the binarization and discretization schemes in all cases and were outperformed by conventional CNNs. In the right-hand column of the table, we show the results of a CT network using a distillation loss and a VGG16 [46] teacher network. These CT networks improve the accuracy significantly further, and the margin of CNNs over them is also significantly smaller. Several key findings can be drawn from these results: First, the good accuracy of CT networks indicates that they can be trained using standard SGD

<sup>1</sup><https://github.com/MatthieuCourbariaux/BinaryConnect>

<sup>2</sup><https://github.com/rarilurelo/XNOR-Net>

<sup>3</sup><https://github.com/zhaohui-yang/Binary-Neural-Networks/tree/main/SLB>

<sup>4</sup><https://github.com/layog/Accurate-Binary-Convolution-Network>

<sup>5</sup><https://github.com/ricky40403/DSQ>

TABLE II

SIX-LAYER DEEP CT AND CNN NETWORKS APPLIED TO THE CIFAR10, CIFAR100, AND SVHN DATASETS. THE CNN BINARIZATION SCHEMES WERE APPLIED TO THE CNN ARCHITECTURE, WHERE  $l$  AND  $K$  DENOTE THE PATCH SIZE AND THE NUMBER OF BITS,  $M$  IS THE NUMBER OF CTS AND  $D_o$  IS THE NUMBER OF OUTPUT MAPS

| Deep CT Layer      | Parameters<br>$l, K, M, D_o$ | Output                   | Convolution Layers | Parameters          | Output                   |
|--------------------|------------------------------|--------------------------|--------------------|---------------------|--------------------------|
| 1 CT Layer 1       | 7, 8, 8, 32                  | $26 \times 26 \times 32$ | Conv,ReLU          | $l = 3, D_o = 32$   | $30 \times 30 \times 32$ |
| 2 Average Pooling  | $l = 3$                      | $24 \times 24 \times 32$ | Max Pooling        | $l = 3, Stride = 2$ | $14 \times 14 \times 32$ |
| 3 CT Layer 2       | 5, 8, 8, 32                  | $20 \times 20 \times 32$ | Conv,ReLU          | $l = 3, D_o = 32$   | $12 \times 12 \times 32$ |
| 4 Average Pooling  | $l = 3$                      | $18 \times 18 \times 32$ | Max Pooling        | $l = 2$             | $11 \times 11 \times 32$ |
| 5 CT Layer 3       | 5, 8, 8, 64                  | $14 \times 14 \times 64$ | Conv,ReLU          | $l = 3, D_o = 64$   | $9 \times 9 \times 64$   |
| 6 Average Pooling  | $l = 3$                      | $12 \times 12 \times 64$ | Max Pooling        | $l = 2$             | $8 \times 8 \times 64$   |
| 7 CT Layer 4       | 5, 8, 8, 128                 | $8 \times 8 \times 128$  | Conv,ReLU          | $l = 3, D_o = 128$  | $6 \times 6 \times 128$  |
| 8 Average Pooling  | $l = 2$                      | $7 \times 7 \times 128$  | Max Pooling        | $l = 2$             | $5 \times 5 \times 128$  |
| 7 CT Layer 5       | 3, 8, 8, 64                  | $5 \times 5 \times 64$   | Conv,ReLU          | $l = 3, D_o = 64$   | $4 \times 4 \times 64$   |
| 8 Average Pooling  | $l = 2$                      | $4 \times 4 \times 64$   | Max Pooling        | $l = 2$             | $3 \times 3 \times 64$   |
| 9 CT Layer 6       | 3, 8, 8, 10                  | $2 \times 2 \times 10$   | Conv,ReLU          | $l = 3, D_o = 10$   | $1 \times 1 \times 10$   |
| 10 Average Pooling | $l = 2$                      | $1 \times 1 \times 10$   | FC                 | $D_o = 10$          | $1 \times 1 \times 10$   |
| 11 SoftMax         |                              |                          | SoftMax            |                     | 1                        |

TABLE III

COMPARISON OF THE CLASSIFICATION ERROR PERCENTAGE. WE COMPARE THE RESULTS FOR CNNs CONSISTING OF TWO LAYERS (MNIST) OR FOUR AND SIX LAYERS (OTHER DATASETS). ALL QUANTIZED/BINARIZED CNNs WERE DERIVED BY APPLYING THE CORRESPONDING QUANTIZATION/BINARIZATION SCHEMES TO THE BASELINE CNN HAVING THE SAME NUMBER OF LAYERS (TABLES VI, VII AND II). THE MOST ACCURATE RESULTS ARE MARKED IN **BOLD**

| Dataset         | DeepCT       | BC [7] | XN [34] | XN++ [4] | TABCNN [27] | DSQ [12] 2b | DSQ [12] 4b | DSQ [12] 6b | Bi-Real [28] | SLB [42] | DeepCT + distil. | CNN   |
|-----------------|--------------|--------|---------|----------|-------------|-------------|-------------|-------------|--------------|----------|------------------|-------|
| Two layers [%]  |              |        |         |          |             |             |             |             |              |          |                  |       |
| <b>MNIST</b>    | <b>0.35</b>  | 0.61   | 0.51    | 0.39     | 0.49        | 0.41        | 0.38        | 0.36        | 0.36         | 0.37     | 0.39             | 0.334 |
| Four layers [%] |              |        |         |          |             |             |             |             |              |          |                  |       |
| <b>CIFAR10</b>  | <b>12.49</b> | 14.47  | 14.77   | 12.97    | 13.97       | 12.87       | 12.71       | 12.68       | 12.72        | 12.66    | 10.91            | 10.86 |
| <b>CIFAR100</b> | <b>39.17</b> | 50.27  | 47.32   | 41.03    | 45.32       | 40.14       | 40.05       | 39.96       | 39.97        | 39.55    | 37.12            | 36.81 |
| <b>SVHN</b>     | <b>4.75</b>  | 6.21   | 5.88    | 4.92     | 5.11        | 4.85        | 4.79        | 4.77        | 4.79         | 4.80     | 4.46             | 4.27  |
| Six layers [%]  |              |        |         |          |             |             |             |             |              |          |                  |       |
| <b>CIFAR10</b>  | <b>11.05</b> | 12.56  | 12.91   | 11.88    | 12.13       | 11.35       | 11.27       | 11.14       | 11.19        | 11.16    | 10.25            | 9.88  |
| <b>CIFAR100</b> | <b>37.95</b> | 42.43  | 40.14   | 39       | 39.37       | 38.31       | 38.25       | 38.12       | 38.15        | 38.05    | 36.05            | 35.58 |
| <b>SVHN</b>     | <b>4.12</b>  | 4.95   | 4.54    | 4.25     | 4.28        | 4.21        | 4.18        | 4.15        | 4.17         | 4.25     | 3.96             | 3.77  |

optimization, such as CNNs, without significant overfitting. Second, the similar precision of the CT and CNN networks with a similar intermediate representation structure indicates that the accuracy is more dependent on the representation hierarchy structure than the particular computing elements used.

### C. Tradeoff Comparison: Accuracy, Speed, and Memory

We studied the Error:Speed and Error:Speed:Memory tradeoffs enabled by the CT networks and compared them with those of contemporary state-of-the-art schemes. The experiments were conducted using the CIFAR10 and CIFAR100 datasets. Seven models of the CT network were trained, with operation count budgets  $3.1 \times 12 \times$  lower than those of the baseline CNN in Table VII ( $7.54 \cdot 10^6$ ). Similar to the baseline CNN, the CT models had six layers. The variations were derived by changing the hyperparameters  $M$  and  $K$  and the number of channels at intermediate layers.

We compare the tradeoffs of the CT networks with those of the MobileNet V2 [35] and ShuffleNet V2 [29] efficient CNN architectures. The first is based on inverted bottleneck residual

TABLE IV  
SPEED AND MEMORY TRADEOFFS COMPARING THE PROPOSED CT, MOBILENET2 [35], AND SHUFFLENET2 [29], USING THE CIFAR-10 AND CIFAR-100 DATASETS. EACH ROW RELATES TO A PARTICULAR CLASSIFICATION PERCENTAGE ERROR AND SHOWS THE CORRESPONDING SPEED AND MEMORY TRADEOFFS. SPEEDUP AND MEMORY RATIOS ARE COMPUTED BY  $(R_{\text{CNN}}/R_{\text{Net}})$ , WHERE  $R_{\text{CNN}}$  IS THE RESOURCES REQUIRED BY THE BASELINE CNN, AS IN TABLE III, AND  $R_{\text{Net}}$  IS THE NETWORK'S REQUIREMENTS

| Error [%] | Speedup |      | Memory ratio |      |
|-----------|---------|------|--------------|------|
|           | CT [35] | [29] | CT [35]      | [29] |
| CIFAR-10  |         |      |              |      |
| 10.6      | 4.91    | 0.72 | 1.47         | 0.09 |
| 11.6      | 5.69    | 1.59 | 3.51         | 0.14 |
| 12.6      | 14.29   | 3.42 | 6.74         | 0.45 |
| CIFAR-100 |         |      |              |      |
| 36.6      | 4.91    | 0.72 | 2.80         | 0.09 |
| 38.1      | 5.69    | 1.59 | 3.51         | 0.14 |
| 39.6      | 13.74   | 3.42 | 6.74         | 0.28 |

blocks, where depthwise convolutions are performed on high-depth representations, and the second used blocks in which half of the channels are processed, followed by a channel

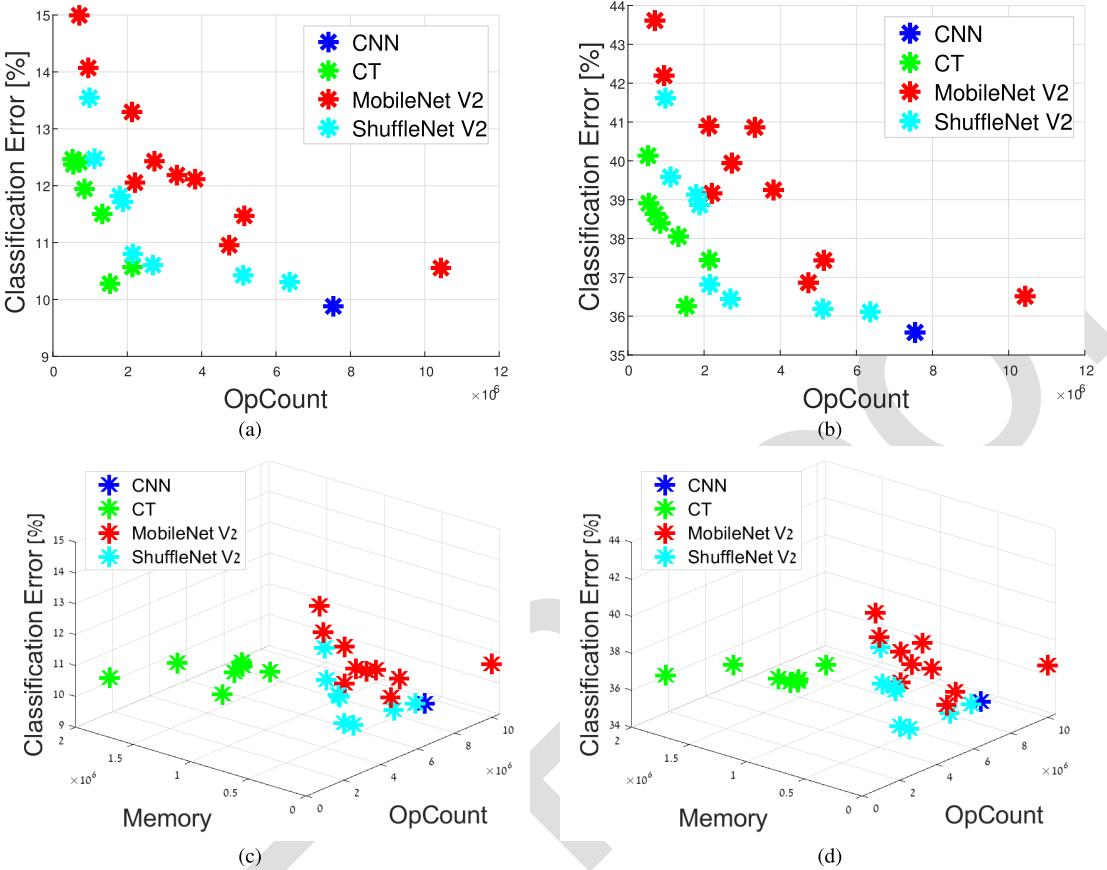


Fig. 3. Performance tradeoffs of the deep CT using an architecture of six layers. (a) Speed:Error tradeoff for several families of models applied to the CIFAR-10 dataset. The models are based on the CNN in Table VII, CT networks, MobileNet V2 [35], and ShuffleNet V2 [29]. We considered architectures using up to  $10^7$  operations. (b) Speed:Error tradeoff on the CIFAR-100 dataset. (c) Speed:Error:Memory tradeoff on CIFAR-10. (d) Speed:Error:Memory tradeoff on CIFAR-100. CT networks provide better Speed:Error tradeoffs in this domain while requiring additional memory. The speed is measured by the MACC operation count. The memory is given by the number of parameters, whereas the error is the classification percentage error.

TABLE V  
ACCURACY PERCENTAGE OF DEEP CT, CNN, MOBILENET V2 [35], AND SHUFFLENET V2 [29] WHEN APPLIED TO SUBSETS OF THE CASIA-WEBFACE DATASET

| Dataset    | Deep CT | BC    | XN     | TABCNN | DSQ 2 bits | DSQ 4 bits | DSQ 6 bits | MobileNet | ShuffleNet | CNN   |
|------------|---------|-------|--------|--------|------------|------------|------------|-----------|------------|-------|
| CASIA 25K  | 87.21   | 85.36 | 86.286 | 86.95  | 87.05      | 87.15      | 87.19      | 87.20     | 87.186     | 89.96 |
| CASIA 250K | 90.66   | 87.95 | 88.35  | 88.89  | 90.06      | 90.27      | 90.35      | 90.63     | 90.54      | 93.25 |

shuffle operation. We used the code provided by the corresponding authors<sup>6,7</sup> and adapted it to the relevant low-compute regime by adjusting their hyperparameters. For MobileNet, the networks were generated by adjusting the expansion parameter  $t$  (Table I in [35]), the number of channels in intermediate representations, and the number of bottleneck modules in a block (between 1 and 2). For ShuffleNet, the networks were generated by changing the number of channels in the intermediate layers and the number of bottleneck modules in a block (between 1 and 4).

Error:Speed results are shown in Fig. 3(a) and (b) for CT networks and competing lean network approaches. These plots show that the proposed CT networks significantly outperform the baseline CNN and competing architectures. To enable

a more explicit assessment, the left-hand side of Table IV shows the speedup achieved by the most efficient networks in each framework, for several error values on the CIFAR-10 and CIFAR-100 datasets. In the high-accuracy end, while losing at most 1% error w.r.t. to the baseline CNN, the CT network enables 4.91 $\times$  acceleration over the CNN and an acceleration of 1.75–6.8 $\times$  over competing methods. When more significant errors are allowed (2.7% for CIFAR-10 and 4% for CIFAR-100), the CT networks provide the accelerations of 13 $\times$ –14 $\times$  over the baseline CNN, significantly outperforming the 6.7 $\times$  acceleration obtained by the best competitor [29]. The baseline CNN uses  $7.54 \cdot 10^7$  MACC operations,  $181 \cdot 10^6$  parameters, and obtains the errors of 9.88 and 35.58 on CIFAR-10 and CIFAR-100, respectively. The tables for CIFAR-10 and CIFAR-100 are similar, as the same architectures are preferable for both datasets in each category of error and network type.

<sup>6</sup><https://github.com/xiaochus/MobileNetV2>

<sup>7</sup><https://github.com/TropComplique/shufflenet-v2-tensorflow>

TABLE VI

TWO-LAYER DEEP CT AND CNN NETWORKS APPLIED TO THE MNIST DATASET.  $l$  AND  $K$  DENOTE THE SIZE OF THE PATCH AND THE NUMBER OF BITS,  $M$  IS THE NUMBER OF CTs AND  $D_o$  IS THE NUMBER OF OUTPUT MAPS

| Deep CT Layer     | Parameters<br>$l, K, M, D_o$ | Output                    | Convolution Layer | Parameters<br>$l, D_o$ | Output                    |
|-------------------|------------------------------|---------------------------|-------------------|------------------------|---------------------------|
| 1 CT Layer 1      | 9, 10, 10, 100<br>$l = 7$    | $20 \times 20 \times 100$ | Conv,ReLU         | $l = 5, D_o = 100$     | $24 \times 24 \times 100$ |
| 2 Average Pooling |                              | $14 \times 14 \times 100$ | Max Pooling       | $l = 5, Stride = 2$    | $10 \times 10 \times 100$ |
| 3 CT Layer 2      | 9, 10, 10, 10<br>$l = 6$     | $6 \times 6 \times 10$    | Conv,ReLU         | $l = 5, D_o = 10$      | $6 \times 6 \times 10$    |
| 4 Average Pooling |                              | $1 \times 1 \times 10$    | Max Pooling+FC    | $l = 6$                | $1 \times 1 \times 10$    |
| 5 SoftMax         |                              | 1                         | SoftMax           |                        | 1                         |

TABLE VII

FOUR-LAYER DEEP CT AND CNN NETWORKS APPLIED TO THE CIFAR10, CIFAR100, AND SVHN DATASETS. THE CNN BINARIZATION SCHEMES WERE APPLIED TO THE CNN ARCHITECTURE, WHERE  $l$  AND  $K$  DENOTE THE PATCH SIZE AND THE NUMBER OF BITS,  $M$  IS THE NUMBER OF CTs, AND  $D_o$  IS THE NUMBER OF OUTPUT MAPS

| Deep CT Layer     | Parameters<br>$l, K, M, D_o$ | Output                   | Convolution Layers | Parameters          | Output                   |
|-------------------|------------------------------|--------------------------|--------------------|---------------------|--------------------------|
| 1 CT Layer 1      | 7, 8, 8, 32<br>$l = 3$       | $26 \times 26 \times 32$ | Conv,ReLU          | $l = 5, D_o = 32$   | $28 \times 28 \times 32$ |
| 2 Average Pooling |                              | $24 \times 24 \times 32$ | Max Pooling        | $l = 3, Stride = 2$ | $13 \times 13 \times 32$ |
| 3 CT Layer 2      | 7, 8, 8, 32<br>$l = 3$       | $18 \times 18 \times 32$ | Conv,ReLU          | $l = 3, D_o = 32$   | $11 \times 11 \times 32$ |
| 4 Average Pooling |                              | $16 \times 16 \times 32$ | Max Pooling        | $l = 3$             | $9 \times 9 \times 32$   |
| 5 CT Layer 3      | 7, 8, 8, 64<br>$l = 3$       | $10 \times 10 \times 64$ | Conv,ReLU          | $l = 3, D_o = 64$   | $7 \times 7 \times 64$   |
| 6 Average Pooling |                              | $8 \times 8 \times 64$   | Max Pooling        | $l = 3$             | $5 \times 5 \times 64$   |
| 7 CT Layer 4      | 7, 8, 8, 10<br>$l = 2$       | $2 \times 2 \times 10$   | Conv,ReLU          | $l = 3, D_o = 10$   | $3 \times 3 \times 10$   |
| 8 Average Pooling |                              | $1 \times 1 \times 10$   | Max Pooling + FC   | $l = 3$             | $1 \times 1 \times 10$   |
| 9 SoftMax         |                              | 1                        | SoftMax            |                     | 1                        |

TABLE VIII

DEEP CT AND CNN ARCHITECTURES APPLIED TO THE SUBSETS OF THE CASIA-WEBFACE DATASET, WHERE  $l$  AND  $K$  DENOTE THE PATCH SIZE AND THE NUMBER OF BITS,  $M$  IS THE NUMBER OF CTs, AND  $D_o$  IS THE NUMBER OF OUTPUT MAPS

| Deep CT Layer      | Parameters<br>$l, K, M, D_o$ | Output                    | Convolution Layers | Parameters          | Output                   |
|--------------------|------------------------------|---------------------------|--------------------|---------------------|--------------------------|
| 1 CT Layer 1       | 9, 8, 8, 32<br>$l = 5$       | $56 \times 56 \times 32$  | Conv,ReLU          | $l = 3, D_o = 32$   | $62 \times 62 \times 32$ |
| 2 Average Pooling  |                              | $52 \times 52 \times 32$  | Max Pooling        | $l = 3, Stride = 2$ | $30 \times 30 \times 32$ |
| 3 CT Layer 2       | 9, 8, 8, 32<br>$l = 5$       | $44 \times 44 \times 32$  | Conv,ReLU          | $l = 3, D_o = 32$   | $28 \times 28 \times 32$ |
| 4 Average Pooling  |                              | $40 \times 40 \times 32$  | Max Pooling        | $l = 2, Stride = 2$ | $14 \times 14 \times 32$ |
| 5 CT Layer 3       | 9, 8, 8, 64<br>$l = 3$       | $32 \times 32 \times 64$  | Conv,ReLU          | $l = 3, D_o = 64$   | $12 \times 12 \times 64$ |
| 6 Average Pooling  |                              | $30 \times 30 \times 64$  | Max Pooling        | $l = 3$             | $10 \times 10 \times 64$ |
| 7 CT Layer 4       | 9, 8, 8, 128<br>$l = 3$      | $22 \times 22 \times 128$ | Conv,ReLU          | $l = 3, D_o = 128$  | $8 \times 8 \times 128$  |
| 8 Average Pooling  |                              | $20 \times 20 \times 128$ | Max Pooling        | $l = 3$             | $6 \times 6 \times 128$  |
| 9 CT Layer 5       | 9, 8, 8, 256<br>$l = 3$      | $12 \times 12 \times 256$ | Conv,ReLU          | $l = 3, D_o = 256$  | $5 \times 5 \times 256$  |
| 10 Average Pooling |                              | $10 \times 10 \times 256$ | Max Pooling        | $l = 3$             | $3 \times 3 \times 256$  |
| 11 CT Layer 6      | 9, 8, 8, 50<br>$l = 2$       | $2 \times 2 \times 50$    | Conv,ReLU          | $l = 3, D_o = 50$   | $1 \times 1 \times 50$   |
| 12 Average Pooling |                              | $1 \times 1 \times 50$    | Fully Connected    |                     | $1 \times 1 \times 50$   |
| 13 SoftMax         |                              | 1                         | SoftMax            |                     | 1                        |

Fig. 3(c) and (d) depicts the Error:Speed:Memory tradeoff for CT networks and other networks, using the CIFAR-10 and CIFAR-100 datasets, respectively. The memory requirements of the most efficient networks for each accuracy level are listed on the right-hand side of Table IV, as compression ratios with respect to the baseline CNN. Clearly, the better Error:Speed tradeoff of the CT networks comes at the expense of a larger memory footprint (additional parameters). This is expected, as in CT transformations with memory-based retrieval replacing dot-product operations, the compute:memory ratio is much lower. For instance, in the high-accuracy regime, using the CIFAR-10 dataset allows a  $4.9\times$  speedup over the baseline CNN while requiring  $10.7\times$  more parameters. However, the size of excess memory can be controlled by limiting  $K$  (for which the memory demand is exponential) to moderate values. Table IV shows that the fastest CT network enabling the

CIFAR-10 error of 12.6% provides a  $14.3\times$  acceleration while requiring only  $2.2\times$  additional parameters over the baseline CNN.

#### D. Face Recognition Test Case

We experimented with CT networks for face recognition using the CASIA-WebFace dataset [43]. Two datasets of  $64 \times 64$  extracted face images were formed. First, we randomly drew a subset of 50 classes (each with 5000 images), 250-K images overall, to evaluate the proposed CT on a large-scale dataset. We also randomly drew 50 identities and a class of negatives (not belonging to the 50 identities), each consisting of 500 images. This exemplifies a typical IoT application, where face recognition is an enabler of personalized IoT services in smart homes, as well as security cameras. For both

datasets, we trained a six-layer CT network and a CNN comparable in terms of the number of layers and intermediate tensor width. The architectures used are detailed in Table VIII in the Appendix. The CT accuracy was compared with the baseline CNN, with the CNN versions produced by the compression techniques used in Section VII-B, and with the efficient CNN frameworks [29], [35] used in Section VII-C. To compare with the six-layer CNN and CT networks, the MobileNet2 and ShuffleNet2 networks were constructed with six blocks: an initial convolution layer and a bottleneck block in the high resolution, followed by two bottleneck blocks in the second and third resolutions. The number of maps at each intermediate representation was identical to the widths used by the CT and CNN networks in Table VIII. The results are presented in Table V. With 25 and 250-K data samples, the CT network outperformed all accelerated methods, with results comparable with MobileNet V2, and was only outperformed by conventional CNN. In particular, the 250-K set results show that the CT can be effectively applied to large datasets.

## VIII. CONCLUSION

In this work, we introduced a novel deep learning framework that utilizes indices computation and table access instead of dot products. The experimental successful validation of our framework implies that both optimization and generalization are not uniquely related to the properties and dynamics of dot-product neuron interaction and that having such neuronal ingredients is an unnecessary condition for successful deep learning. Our analysis and experimental results show that the suggested framework outperforms conventional CNNs in terms of computational complexity:capacity ratio. Empirically, we showed that it enables improved speed:accuracy tradeoff for the low-compute inference regime at the cost of additional processing memory. Such an approach is applicable to a gamut of applications that are either low-compute or are to be deployed on laptops/tablets that are equipped with large memory, but lack dedicated GPU hardware. Future work will extend the applicability of CT networks using an improved GPU-based training procedure. In this context, it should be noted that only basic CT networks were introduced and tested in this work, whereas competing CNN architectures enjoy a decade of intensive empirical research. The fact that CT networks were found competitive in this context is therefore highly encouraging. Combining CT networks with CNN optimization techniques (batch/layer normalization), regularization techniques (dropout), and architectural benefits (Residuals connections) is yet to be tested. Finally, to cope with the memory costs, methods for discretization/binarization of the model's parameters should be developed, similar to CNN discretization algorithms.

## APPENDIX

The architectures used in Section VII-B for networks consisting of two and four layers are described in Tables VI and VII, respectively. The architectures used for CASIA-WebFace data in Section VII-C are detailed in Table VIII.

## REFERENCES

- [1] A. Waissman and A. Bar-Hillel, "Input-dependent feature-map pruning," in *Proc. Int. Conf. Artif. Neural Netw. (ICANN)*, 2018, pp. 706–713. 886  
887
- [2] M. Al-Hami, M. Pietron, R. Casas, S. Hijazi, and P. Kaul, "Towards a stable quantized convolutional neural networks: An embedded perspective," in *Proc. 10th Int. Conf. Agents Artif. Intell.*, 2018, pp. 573–580. 888  
889  
890
- [3] E. Bengio, P.-L. Bacon, J. Pineau, and D. Precup, "Conditional computation in neural networks for faster models," 2015, *arXiv:1511.06297*. 891  
892
- [4] A. Bulat and G. Tzimiropoulos, "XNOR-Net++: Improved binary neural networks," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, 2019, p. 62. 893  
894
- [5] S. R. Bulo and P. Kotschieder, "Neural decision forests for semantic image labelling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2014, pp. 81–88. 895  
896  
897
- [6] W. Chen, T. J. Wilson, S. Tyree, Q. K. Weinberger, and Y. Chen, "Compressing neural networks with the hashing trick," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2015, pp. 2285–2294. 898  
899  
900
- [7] M. Courbariaux, Y. Bengio, and J.-P. David, "BinaryConnect: Training deep neural networks with binary weights during propagations," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Cambridge, MA, USA, 2015, pp. 3123–3131. 901  
902  
903  
904
- [8] S. A. Davis and I. Arel, "Low-rank approximations for conditional feed-forward computation in deep neural networks," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Y. Bengio and Y. LeCun, Eds., 2014, pp. 1–10. 905  
906  
907
- [9] L. Deng and D. Yu, "Deep learning: Methods and applications," *Found. Trends Signal Process.*, vol. 7, nos. 3–4, pp. 197–387, Jun. 2014. 908  
909
- [10] Y. Eidelman, V. D. Milman, and A. Tsolomitis, *Functional Analysis: An Introduction*, vol. 66. 2004. 910
- [11] J. Ephrath, L. Ruthotto, E. Haber, and E. Treister, "LeanResNet: A low-cost yet effective convolutional residual networks," in *Proc. ICML Workshop On-Device Mach. Learn. Compact Deep Neural Netw.*, 2019, pp. 1–11. 912  
913  
914  
915
- [12] R. Gong et al., "Differentiable soft quantization: Bridging full-precision and low-bit neural networks," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 4851–4860. 916  
917  
918
- [13] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2016, pp. 1–14. 919  
920  
921
- [14] M. Hassoun, *Fundamentals of Artificial Neural Networks*. 1995. 922
- [15] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," 2015, *arXiv:1503.02531*. 923  
924
- [16] K. Huang, B. Ni, and X. Yang, "Efficient quantization for neural networks with binary weights and low bitwidth activations," in *Proc. AAAI Conf. Artif. Intell.*, 2019, pp. 3854–3861. 925  
926  
927
- [17] Y. Ioannou et al., "Decision forests, convolutional networks and the models in-between," 2016, *arXiv:1603.01250*. 928  
929
- [18] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," in *Proc. Brit. Mach. Vis. Conf.*, 2014, pp. 1–12. 930  
931  
932
- [19] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Bulo, "Deep neural decision forests," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1467–1475. 933  
934  
935
- [20] R. Krishnamoorthi, "Quantizing deep convolutional networks for efficient inference: A whitepaper," 2018, *arXiv:1806.08342*. 936  
937
- [21] A. Krizhevsky, V. Nair, and G. Hinton, "CIFAR-10 (Canadian Institute for Advanced Research)," Tech. Rep., 938
- [22] A. Krizhevsky, V. Nair, and G. Hinton, "CIFAR-100 (Canadian Institute for Advanced Research)," Tech. Rep., 940  
941
- [23] E. Krupka, A. Vinnikov, B. Klein, A. B. Hillel, D. Freedman, and S. Stachniak, "Discriminative ferns ensemble for hand pose recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2014, pp. 3670–3677. 942  
943  
944  
945
- [24] E. Krupka et al., "Toward realistic hands gesture interface: Keeping it simple for developers and machines," in *Proc. CHI Conf. Human Factors Comput. Syst.*, New York, NY, USA, May 2017, pp. 1887–1898. 946  
947  
948
- [25] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," *ATT Labs Online*, vol. 2, 2010. 949
- [26] Y. LeCun, S. John Denker, and A. Sara Solla, "Optimal brain damage," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, D. S. Touretzky, Ed., 1990, pp. 598–605. 950 AQ:11  
951  
952  
953  
954  
955  
956
- [27] X. Lin, C. Zhao, and W. Pan, "Towards accurate binary convolutional neural network," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 344–352.

- [28] Z. Liu, B. Wu, W. Luo, X. Yang, W. Liu, and K. Cheng, "Bi-Real Net: Enhancing the performance of 1-bit CNNs with improved representational capability and advanced training algorithm," in *Proc. 15th Eur. Conf.*, Sep. 2018, pp. 747–763.
- [29] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "ShuffleNet V2: Practical guidelines for efficient CNN architecture design," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 116–131.
- [30] Y. Mansour, "Pessimistic decision tree pruning based on tree size," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 1997, pp. 195–201.
- [31] B. Martinez, J. Yang, A. Bulat, and G. Tzimiropoulos, "Training binary neural networks with real-to-binary convolutions," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020, pp. 1–11.
- [32] M. Mathieu, M. Hena, and Y. LeCun, "Fast training of convolutional networks through FFTs," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014, pp. 1–11.
- [33] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2011, pp. 1–9.
- [34] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: ImageNet classification using binary convolutional neural networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, B. Leibe, J. Matas, N. Sebe, M. Welling, Eds. Cham, Switzerland: Springer, 2016, pp. 525–542.
- [35] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 4510–4520.
- [36] S. Ren, X. Cao, Y. Wei, and J. Sun, "Global refinement of random forest," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 723–730.
- [37] J. Shotton et al., "Real-time human pose recognition in parts from single depth images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 1297–1304.
- [38] X. Su et al., "Searching for network width with bilaterally coupled network," *IEEE Trans. Pattern Anal. Mach. Intell.*, early access, Dec. 5, 2022, doi: [10.1109/TPAMI.2022.3226777](https://doi.org/10.1109/TPAMI.2022.3226777).
- [39] Y. Tang et al., "Manifold regularized dynamic network pruning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Los Alamitos, CA, USA, Jun. 2021, pp. 5016–5026.
- [40] M. Tu, V. Berisha, M. Woolf, J.-S. Seo, and Y. Cao, "Ranking the parameters of deep neural networks using the Fisher information," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Mar. 2016, pp. 2647–2651.
- [41] B. Wu et al., "Shift: A zero FLOP, zero parameter alternative to spatial convolutions," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 9127–9135.
- [42] Z. Yang et al., "Searching for low-bit weights in quantized neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Red Hook, NY, USA, 2020, pp. 1–12.
- [43] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning face representation from scratch," 2014, *arXiv:1411.7923*.
- [44] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 7130–7138.
- [45] O. T. Yıldız, "VC-dimension of univariate decision trees," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 2, pp. 378–387, Feb. 2015.
- [46] X. Zhang, J. Zou, K. He, and J. Sun, "Accelerating very deep convolutional networks for classification and detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 1943–1955, Oct. 2015.
- [47] Z.-H. Zhou and J. Feng, "Deep forest: Towards an alternative to deep neural networks," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, Aug. 2017, pp. 3553–3559.



**Shay Dekel** received the B.Sc. degree (summa cum laude) in electrical engineering from Bar Ilan University, Ramat-Gan, Israel, in 2005, and the M.Sc. degree (summa cum laude) in electrical engineering from Tel Aviv University, Tel Aviv, Israel, in 2011. He is currently pursuing the Ph.D. degree with the Faculty of Electrical Engineering, Bar Ilan University.

His research interests include neural 3-D reconstruction from sparse views, neuromorphic computing, and cognitive computing.



**Yosi Keller** received the B.Sc. degree in electrical engineering from the Technion-Israel Institute of Technology, Haifa, Israel, in 1994, and the M.Sc. and Ph.D. degrees (summa cum laude) in electrical engineering from Tel Aviv University, Tel Aviv, Israel, in 1998 and 2003, respectively.

He served as a Gibbs Assistant Professor for Yale University, New Haven, CT, USA, from 2003 to 2006. He is currently an Associate Professor with the Faculty of Engineering, Bar Ilan University, Ramat-Gan, Israel. His research interests include computer vision and deep learning.



**Aharon Bar-Hillel** was born in Beer-Sheva, Israel, in 1971. He received the B.A. degree in math and philosophy from Tel-Aviv University, Tel Aviv, Israel, in 1999, and the Ph.D. degree in neural computation from The Hebrew University of Jerusalem, Jerusalem, Israel, in 2006.

He was a Machine Learning and Computer Vision Researcher with Intel Research, from 2006 to 2008, GM Research, from 2009 to 2012, and a Microsoft Research, from 2013 to 2016. From 2016 to 2022, he served as a Senior Lecturer with the Department of Industrial Engineering and Management, Ben-Gurion University, Negev, Israel. He is currently a Researcher with Saips, a Ford Company. His research interests include interpretation and acceleration of deep networks, as well as applications of computer vision and machine learning for ADAS, agricultural phenotyping, and visual tracking.