# Distill-Net: Application-Specific Distillation of Deep Convolutional Neural Networks for Resource-Constrained IoT Platforms

Mohammad Motamedi
University of California, Davis
One Shields Avenue
Davis, CA 95616, USA
mmotamedi@ucdavis.edu

Felix A. Portillo
University of California, Davis
One Shields Avenue
Davis, CA 95616, USA
faportillo@ucdavis.edu

Daniel Fong
University of California, Davis
One Shields Avenue
Davis, CA 95616, USA
dfong@ucdavis.edu

Soheil Ghiasi
University of California, Davis
One Shields Avenue
Davis, CA 95616, USA
ghiasi@ucdavis.edu

## ABSTRACT

Many Internet-of-Things (IoT) applications demand fast and accurate understanding of a few key events in their surrounding environment. Deep Convolutional Neural Networks (CNNs) have emerged as an effective approach to understand speech, images, and similar high dimensional data types. Algorithmic performance of modern CNNs, however, fundamentally relies on learning class-agnostic hierarchical features that only exist in comprehensive training datasets with many classes. As a result, fast inference using CNNs trained on such datasets is prohibitive for most resource-constrained IoT platforms. To bridge this gap, we present a principled and practical methodology for distilling a complex modern CNN that is trained to effectively recognize many different classes of input data into an application-dependent essential core that not only recognizes the few classes of interest to the application accurately, but also runs efficiently on platforms with limited resources. Experimental results confirm that our approach strikes a favorable balance between classification accuracy (application constraint), inference efficiency (platform constraint), and productive development of new applications (business constraint).

## 1 INTRODUCTION

Deep Convolutional Neural Networks (CNN) have demonstrated remarkable performance in pattern recognition and classification tasks on high dimensional data, such as visual recognition and speech understanding. State of the art CNNs need to be trained on comprehensive labeled datasets, which include a large number of classes and many examples per class. Such CNNs typically involve billions of operations on millions of parameters, hindering their utilization in IoT applications that demand fast inference on resource-constrained execution platforms.

We observe that a significant subset of IoT applications [5] require recognition of a fairly small number of classes, compared to the number of classes that exist in commonly used labeled training datasets. For example, disaster response drones or security cameras tend to require detection of a few key classes, such as person, tree, animal, road, obstacle and such, while the ILSVRC dataset [30] that is used to train many recent CNNs includes more than 1000 classes, such as various breeds of dogs, types of bugs, and many exotic animals and uncommon objects. Figure 1 shows the diversity of image classes in ILSVRC dataset [30]. The first row includes samples from classes that are beneficial for an outdoor hazard detection application for children. Other rows show images from just a few, out of hundreds, of classes that are not useful for this application. IoT systems are mission-oriented in nature, and thus, require recognition of classes that matter to their application. For many typical applications, such as those arising in precision agriculture, smart environment, navigation, and security [4], the number of classes of interest is limited.

Given that large labeled datasets are expensive to generate, it is unrealistic to expect that for every new application, new labeled training data is available, and a new model can be constructed and trained from scratch. A key open question, thus, has to do with a practical approach to build application-specific CNNs that offer accurate, from a classification viewpoint, and efficient, from an execution perspective, mission-specific inference on resource-constrained IoT platforms. In this paper, we propose a technique towards addressing this very question.

On the surface, the problem may seem misleadingly trivial. It may appear that only the training samples relevant to the target application could be utilized to train a simpler custom CNN. This approach fails to produce CNNs that meet classification accuracy constraints: CNNs are effective in complex tasks, in part because they manage to learn class-agnostic features that exists only in comprehensive datasets. Training a network using a small subset of application-specific examples deprives the network of such features, which in turn, harms its accuracy. In other words, negative examples

**Figure 1: Samples from ILSVRC dataset [30]. The first row shows images from classes that are essential for an outdoor hazard detection application, whereas other rows present samples from just a few classes, out of hundreds of classes, that are not useful for this task. The code of the class to which each image belongs is shown.**

contain significant information, and removing them from the training dataset is detrimental to classification accuracy.

Another naive approach to the problem would be to start with a state of the art trained CNN, and replace its final-stage classifier with an application-specific classifier that detects only a few classes of interest to the application. This approach certainly meets the accuracy requirement. However, it results in a CNN that is computationally very similar to the original model, and thus, does not yield fast inference on IoT platforms. Note that the overwhelming majority of computational workload of a CNN does not occur in its classifier stage [32]. Such a process demands a follow up fine-tuning which includes multiple epochs. The required time for **each** epoch is five times larger than the required time for the entire distillation [20].

Starting with a state of the art trained CNN and the comprehensive publicly-available dataset that was used for its training, we aim to systematically develop a custom CNN that meets the two inference efficiency and classification accuracy constraints, albeit on a small subset of classes. In our proposed methodology, IoT designers specify the classes of interest to the application. Subsequently, we analyze the input CNN, and distill it into its essential core with respect to the target application. Experimental results show that Distill-Nets at most drop the accuracy by 1% for classes of interest to the application, while running significantly faster.

## 2 LANDSCAPE AND CONTRIBUTIONS

In theory, one could design, train and optimize a CNN to meet the requirements of a given IoT application. In practice, this approach faces several challenges and shortcomings:
(1) Time-consuming nature of training deep networks: Designing and training a new industrial-strength CNN for every

new IoT application increases time to market. Such a process requires many rounds of trial-and-error, training, and optimization. Each round of training can take from a day to a week depending on the available hardware and the CNN architecture. In contrast, the proposed distillation approach can finish the task in a matter of minutes.
(2) Inherent difficulty of developing competitive models: While embedded system developers are experts in various aspects of designing time-predictable, and energy-efficient applications for resource-constrained IoT platforms, they are not necessarily proficient in developing and training competitive deep neural networks.
(3) Considerable cost of computational infrastructure to support the training: The process of designing and training a new CNN is computationally intensive and demands significant investment in computing infrastructure.
(4) Data privacy: The proposed distillation algorithm can be performed on a mediocre IoT platform (e.g., Google Nexus 6P) making it feasible to tailor trained CNNs to users' private data. However, retraining a CNN in such a case requires the data to be uploaded to a cloud which may compromise users' privacy.
(5) Dynamic mission assignment: In the training-based approaches, re-purposing a CNN to identify a new set of classes requires it to be retrained. However, since the distillation process we are proposing preserves the knowledge (i.e., weights) of all classes, redeploying a CNN for a new set of image objectives becomes quick and feasible. Distillation merely creates guidance bitmaps for the thread launcher to specify neurons whose computation can be skipped. That is, it preserves the CNN structure, and the corresponding weights (i.e., the knowledge of the original CNN is maintained). As a result, reassigning the purpose of an embedded application with a distilled CNN is simply the matter of discarding or updating the bitmaps.

Our study is performed under two constraints: First, in all of the experiments, the distillation algorithm is targeted to lose at most 1% of the classification accuracy. In practice, in more than 60% of our experiments we achieved less than 0.1% loss in the classification accuracy.

Second, IoT devices are required to upload their private or proprietary data to a computation cloud, if they need a CNN to be tailored to their specific application. We believe that users' privacy protection must be considered an integral part of any IoT solution. The proposed distillation is designed under this constraint. In section 8.2.2, we show that a mediocre IoT device can run the algorithm on users' private data in a reasonable time budget.

The main contributions of this paper are: First, for the first time, we propose an approach to tailor a CNN for an IoT application whose objects of interest are a fraction of categories that the original CNN was designed to classify. Second, the proposed approach requires no additional **data**, training or fine-tuning, thus, addressing practical concerns, such as resource-aware inference, dynamic updates, and productive development.

Note that the term "distillation" is coincidentally used in some other articles such as [7, 13] to refer to a CNN training approach that utilizes the knowledge of an ensemble of trained CNNs in addition to the training set. It is essential to emphasize that while this paper uses the same term, it targets a substantially different problem. In the distillation algorithm proposed in this article, we offer a solution for tailoring a CNN to an IoT application whose objects of interest are a subset of categories that the original CNN was designed to classify.

## 3 CONVOLUTIONAL NEURAL NETWORKS

State-of-the-art CNNs are capable of making mostly correct classifications about images with invariance to the shape, size, rotation, or position of the desired object in a scene [23, 32]. The key feature of CNNs is the convolutional layers, which consist of sets of trainable 3D filter-banks. These filter-banks extract features from Input Feature Maps (IFMs) to produce sets of Output Feature Maps (OFMs). By stacking multiple layers, a CNN can extract more abstract information from simple edges and colors to entire objects and patterns as the network gets deeper. However, having too many layers increases the capacity of the neural network (i.e., larger number of parameters), which makes it prone to over-fitting.

In essence, each 3D filter-bank is a neuron with a restricted receptive field. In the rest of this paper, the terms filter-bank and neuron are used interchangeably to convey the message clearly.

## 4 PROBLEM STATEMENT

Let us assume that a given CNN, $\Psi$, is trained using a large dataset whose data belongs to $\alpha$ distinct classes which are shown by set $A$ in Equation (1). Let us also assume that for a particular IoT application, only $\beta$ classes that are shown in set $B$ are required, such that, $B \subset A$. The question is, can we use $\Psi$ to create a smaller CNN that is adequate for classifying members of set $B$?

$$A = \{a_1, a_2, \cdots, a_\alpha\}, \ B = \{b_1, b_2, \cdots, b_\beta\} \qquad (1)$$

The original CNN can serve the purpose. However, for each inference, it extracts all the features that are required for understanding members of set $A$, some of which will never be used to classify images of set $B$. In order to improve the execution time and energy consumption, we are interested in finding neurons that solely react to members of set $\overline{B}$ and removing them. $\overline{B}$ is the complement of set $B$ where set $A$ is the universe.

## 5 CNN-BASED FEATURE EXTRACTION

One of the most important characteristics of CNNs that distinguishes them from other image recognition approaches is their ability in automatic feature extraction. Given a large labeled dataset and sufficient computing resources, a CNN with an appropriate capacity can automatically extract adequate feature-sets that can be used for image classification. Features are learned in a layer-based fashion from the shallowest layer to the deepest one. In the very first layer, basic
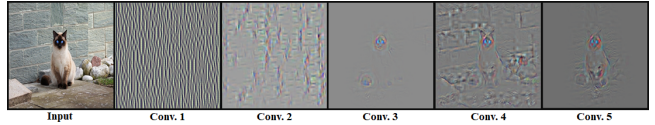


**Figure 2: Strongest feature map for each layer of AlexNet [23]. As we move towards deeper layers, features are more complex, exaggerate discriminative parts, and have a greater invariance.**

feature extractors such as line, curve, and dot detectors form. Even though such basic features do not directly contribute to the classification step, they have a very critical role in improving the final accuracy as they are the building blocks based on which feature extractors in subsequent layers form.

In shallow layers, features are simple and structurally related to the input. In contrast, features that are obtained in deep layers are more complicated and semantically related to the input. The functionality of neurons that generate such features (semantically related) is location-agnostic. That is, the stimulus that makes them fire is the existence of a pattern not its location.

We used Deconvolutional Neural Networks [6, 34] to visualize feature maps for different layers of AlexNet [23]. Figure 2 illustrates the strongest feature map per each layer of the CNN. For this specific example, rich features start forming in convolutional layer #3. In subsequent layers, CNN improves the quality of the features to make them even more distinctive for different classes in the classification task. Features which are extracted in deep layers intensify discriminative parts of objects. In this example, AlexNet mainly uses the face, eyes, and nose, for distinguishing a cat.

In this paper, we use Class-Agnostic Features (CAFs) and Class-Dependent Features (CDFs) terms to refer to features that are extracted in shallow layers and deep layers, respectively. The former exist in all images and the latter exist in images that belong to one or a few specific classes. Likewise, Class-Agnostic Neurons (CANs) and Class-Dependent Neurons (CDNs) are terms used for neurons that generate CAFs and CDFs, respectively. CNNs utilize all of the labeled data (from all classes) to train CANs. However, they are limited to the data of each class for training its CDNs.

When the labeled data that is available for training a particular CNN is inadequate, it is possible to use another dataset, which includes the same type of data, to train CANs. Subsequently, the original dataset must be used to fine-tune CANs and train CDNs.

## 6 NEUROACTIVITY ANALYSIS IN CNNS

Since Output Feature Maps (OFMs) store neurons' outputs, their sparsity is a good measure of neuron activities in each layer. Neuron activity is the reaction of a neuron to a particular input, i.e., whether it fires and if so, how strongly it does. In the first part of this section, we propose an approach for quantifying the sparsity of an OFM. We hypothesize that OFMs which hold outputs of CANs are less sparse compared to those that store outputs of CDNs since CANs fire more

frequently. In the second part of this section, we check the validity of this hypothesis in GoogLeNet [32].

## 6.1 Dormant Neurons Density

Computation that a neuron performs on a tile $t$ of IFMs is shown in Equation (2), in which $f^t$, $K$, $N$, and $b$ are neuron's output for the tile $t$, kernel size, number of IFMs, and the bias value, respectively. In this equation, $X^t$ is a matrix including a tile of the IFM, i.e., the neuron activities from the previous layer, $W$ is a matrix that contains parameters of a neuron, and the "+" superscript shows the positive part of the function.

If Equation (3) holds (linear flattening), Equation (2) can be rewritten in a simplified form which is shown in Equation (4).

We define $X''$ to be a sorted equivalent of $X'$ whose values are in a descending order. Let us assume that function $M$, which is shown in Equation (5), is a mapping between indexes of corresponding elements in $X'$ and $X''$. By creating $W''$ from $W'$ using mapping $M$, it is possible to compute $f^t$ using Equation (6).

$$f^t = \left( \sum_{i=1}^{N} \sum_{j=1}^{K} \sum_{k=1}^{K} W_{i,j,k} \times X_{i,j,k}^t + b \right)^+ \tag{2}$$

$$\forall\, l \in \mathbb{N}_{NK^2}:$$

$$if\ l = K \times i \times K + j + k \Rightarrow W_l' = W_{i,j,k}\ and\ X_l' = X_{i,j,k}^t \tag{3}$$

$$f^t = \left( \sum_{l=1}^{NK^2} W_l' \times X_l' + b \right)^+ \tag{4}$$

$$if\ j = Mi \Rightarrow X_j'' = X_i' \tag{5}$$

$$f^t = \left( \sum_{l=1}^{NK^2} W_l'' \times X_l'' + b \right)^+ \tag{6}$$

Since a considerable number of neurons are dormant in each layer, X is expected to be a sparse vector. Assuming a sparsity ratio of $\tau_t$ for $X$, one can compute $f^t$ using Equation (7).

$$f^t = \left( \sum_{l=1}^{NK^2 - \tau_t NK^2} W_l'' \times X_l'' + b \right)^+ \tag{7}$$

Neural networks are designed to be noise resilient. Hence, it is possible to approximate small elements in $X$ with zeros to increase the sparsity ratio from $\tau_t$ to $\tau_t'$. This can be performed using a cutoff threshold: $TH_t$.

The value of $\tau_t'$, which is an indicator of the number of neurons that did not fire for a given input, is called dormant neurons' density in the rest of the paper. The difference between $\tau_t$ (neurons which are actually dormant) and $\tau_t' - \tau_t$ (neurons that are considered dormant) is important. The former simply reflects the number of zeros in OFMs. However, the latter is the number of small elements that were approximated by zero without changing the classification results. The largest value of $TH_t$ for **each tile** of a given input can be found by performing a binary search on the range of possible values and observing the effect on the classification accuracy.

For ease and efficiency of implementation, we compute a single sparsity ratio for each layer of a CNN (instead of each tile) whose upper bound can be calculated using Equation (8) where $T$ is the total number of tiles in the IFMs of layer
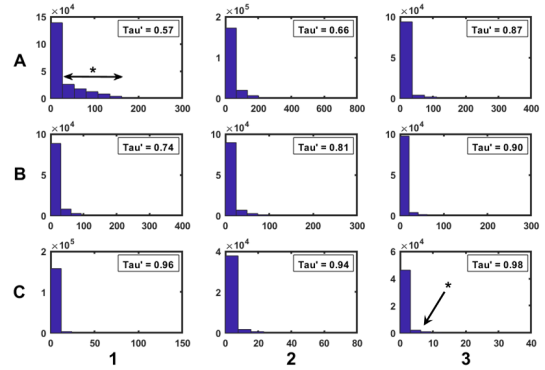


**Figure 3:** Histogram of OFMs of main layers of GoogLeNet for a sample input. By moving towards deeper layers, the value of $\tau'$, i.e., the number of neurons that have negligible impact on classification result, increases. A1, A2, A3, B1, B2, B3, C1, C2, and C3 are outputs of layers: CL2, I3A, I4A, I4B, I4C, I4D, I4E, I5A, and I5B, respectively.

$l$. That is, for each layer of a CNN finding only one cutoff value ($TH_l$) suffices.

$$\tau_l' \leq max\{\tau_1', \tau_2', \cdots, \tau_T'\} \tag{8}$$

The value of dormant neurons' density is expected to be smaller for shallow layers which mainly consist of CANs. However, it ought to be larger for deep layers which include a large number of CDNs.

## 6.2 Case Study: Dormant Neurons' Density in Different Layers of GoogLeNet

Figure 3 illustrates the histogram of OFM values for nine major layers of GoogLeNet. In the shallowest layer (A1), a considerable number of neurons are active (shown by an asterisk). This phenomenon was expected since most of the CANs which reside in the first layer collaborate to extract features that will be used by subsequent layers. As we move towards the deeper layers, the value of $\tau'$ gradually increases from 0.57 to 0.98. That is, in the very last Inception layer (C3), which mainly includes CDNs, 98% of neurons do not contribute to classifying the given input. Comparing the asterisk in (A1) to the one in (C3) shows how activity patterns change as we approach the deepest layers of a CNN.

We visualized nine random OFMs per each of the major layers of GoogLeNet and the results are depicted in Figure 4 (the order of layers is congruent with that of Figure 3). There are two important patterns to observe. First, the variation of the dormant neurons' density from the very first layer (A1) to the very last one (C3), and how an increasing area of heatmaps gets cooler. Second, in the deep layers, some of the active neurons that fire have very high output values. Those seem to be mainly CDNs that have a high sensitivity for the class to which the given input belongs.
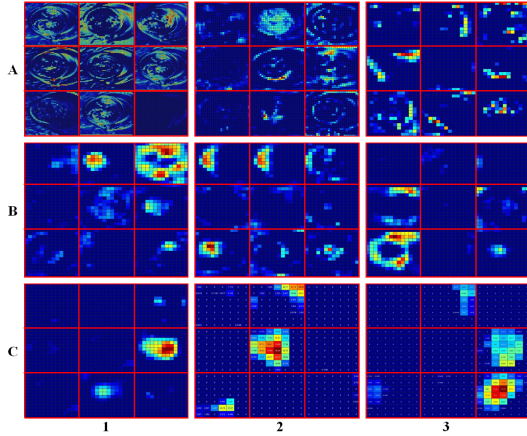
4

**Figure 4: Nine random OFMs per each major layer of GoogLeNet. The order of layers is congruent with that of Figure 3. As we go deeper in the network, more neurons become dormant. However, those that fire have higher output values. This indicates that the CNN is forming a strong opinion about the label of the input.**



**Figure 5: Dormant and active neurons. OFMs A and B in layer #$l$ contain the output of neurons A and B, respectively. A red cross in location $i, j$ of an OFM indicates that the corresponding neuron has fired in that particular position.**

## 7  CNN DISTILLATION

Running server-grade CNNs on resource-constrained IoT devices can be inefficient from various facets including:

**(1) Energy Consumption:** It is possible to efficiently accelerate CNNs on various IoT devices [8, 11, 27, 29, 35]. However, even after such an implementation, the required energy budget can be large for most battery-operated platforms. For example, Motamedi et al. showed that an efficient acceleration of GoogLeNet on mobile System-on-Chips (SoCs), requires 3.24J energy per inference [26].

**(2) Applicability in IoT Domains:** Many IoT applications demand fast and accurate understanding of a few key events in their surrounding environment. Recognizing unimportant events, which imposes extra stress on the energy source, is not the best use of our constrained resources in an IoT platform. As an example, in developing an IoT application for pet monitoring, recognition of presence or absence of the animal in a camera feed can be sufficient. In such cases, we do not need the CNN to be able to distinguish the breed of the dog. Instead, it is expected to work reliably in terms of execution time, energy consumption, and heat dissipation. Nonetheless, using GoogLeNet for such tasks, executes all of its neurons, including those whose outputs are unimportant (e.x., breed detection) or not applicable (e.x., whale detector).

On one hand, using very high-end CNNs for IoT-class applications can waste the available computing power. On the other hand, since training CANs benefits from data of all classes, retraining a new CNN with a limited number of training examples (only from those classes that belong to the scope of interest of a particular task) can lead to an inferior performance. In this section, we propose a principled methodology for identifying CANs and CDNs whose outputs are essential for understanding events that are in the interest of a particula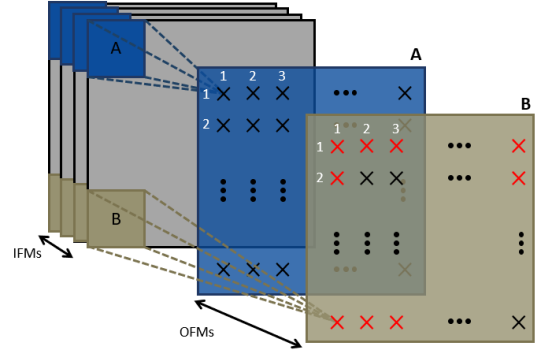r IoT application. The detected neurons form a smaller CNN which is capable of fulfilling the required tasks while avoiding unnecessary computations.

### 7.1  Neuron Removal

*7.1.1  Complete Neuron Removal.* Let us assume that a neuron always remains dormant for a particular application. That is, it either never fires or, if it fires, the output value is negligible. This is illustrated in Figure 5, where neuron A refuses to fire on all of the different spatial locations of IFMs. If this phenomenon deterministically repeats for a specific learning task, it is computationally beneficial to extract such a neuron. Therefore, neuron A in Figure 5 should be removed which leads to elimination of OFM A and saves all of the computations that should have been performed to create it. Neuron removal in a layer $l$ results in computation reduction in the same layer and computation reduction in the subsequent layer (since it decreases the number of OFMs in layer $l$. OFMs in layer $l$ are IFMs to layer $l + 1$). Equation (9) and Equation (10) show the number of Multiply-Accumulation (MAC) operations which are skipped in layers $l$ and $l + 1$, respectively. In these equations, $W_{out}$, $H_{out}$, $N$, $K$, and $S$ stand for width of OFMs, height of OFMs, number of IFMs, kernel size, and convolution stride, respectively. Subscript $l$ emphasizes that values of parameters vary across different layers of the neural network.

$$\#\text{skipped MACs}_l = \frac{Wout_l \times Hout_l \times N_l \times K_l^2}{S_l^2} \quad (9)$$

$$\#\text{skipped MACs}_{l+1} = \frac{Wout_{l+1} \times Hout_{l+1} \times M_{l+1} \times K_{l+1}^2}{S_{l+1}^2} \quad (10)$$

*7.1.2  Partial Neuron Removal.* As illustrated in Figure 5, there are cases where a neuron fires on some regions of IFMs and stays dormant on other regions (Neuron B). If a neuron remains dormant in a certain spatial location for a particular task, it is computationally efficient to skip calculating its value in that specific location.

The number of dormant neurons that can be removed varies based on the original CNN and the type of the IoT application that the CNN is being distilled for. This number also varies among different layers of a CNN. In what follows we continue our discussion by offering a principled approach for recognition and partial or complete removal of dormant neurons.

## 7.2 Neuron Removal vs. Weight Pruning

Weight pruning is concerned with manipulating neuron interconnections by approximating small weight values with zeros [12]. Each pruning changes the topology of CNN architecture and attenuates the relationship of two neurons. Pruning increases the sparsity of weight matrix; however, that does not significantly drop the execution time since the current computing infrastructures are very inefficient when it comes to numerical calculation on non-uniform sparse matrices [12, 21, 32].

**Fallacy**: If we realize $W_{i,j,k}$ in Equation (2) equals zero during runtime, avoiding the computation of $W_{i,j,k} \times X_{i,j,k}^t$ drastically improves the power consumption.

The aforementioned statement is not quite accurate since under 45 nm CMOS technology, a 32-bit operation on average requires 1 pJ energy, while a 32-bit cache access demands 6 pJ and a 32-bit DRAM access demands 640 pJ. Given a hit rate of $h$, Equation (11) compares the required energy for performing multiplication with the required energy for loading the operands. For instance, if cache's hit rate equals to 0.95, avoiding the multiplication improves the energy consumption by 2.58%. That is, 97.42% of the energy budget is dedicated to loading the weight value. An analogous analysis can be performed to show that the effect of such optimization is also limited on the execution time. In addition, using a branch to check the value of $W$ decreases the efficiency of speculative instruction execution which can further decrease the performance.

$$\text{energy saving} = \frac{1}{1 + h \times 6 + 1 - h \times 640} \times 100 \qquad (11)$$

In contrast, neuron removal targets a group of weights which are semantically and structurally related (a neuron) and aims to eliminate all of them together when possible. A complete neuron removal requires no further overhead. That is, during runtime it is not required to check a flag to determine if a neuron had been removed. Hence, the energy consumption and execution time proportionally decrease by the sum of the values that Equation (9) and Equation (10) compute. For the partial neuron removals, however, a bitmap is required for bookkeeping purposes to determine if a neuron fires in a spatial location. For each location, if the bitmap's value is one, all of the computations of the corresponding neuron in that specific place has to be performed. Otherwise, all the instructions which includes $2NK^2$ memory and $NK^2$ MAC operations can be skipped.

In addition, modern computing devices, especially GPUs, do not handle branching effectively [31]. The key difference between partial neuron removal, rather than pruning, is that removing a neuron impacts a large chunk of computations. Therefore, the outcome easily compensates the overhead of decision-making processes. In addition, complete neuron removal does not require branching.

After pruning a CNN, it must preserve its ability to classify different images. Hence, a pruned network can generate distinctive features that are used to discriminate between different classes. That is, in a pruned CNN there are kernels which extract features that are exclusively used for identification of images of a particular class. This is the sufficient condition for distillation to work. Hence it is possible to distill a pruned CNN or prune a distilled CNN.

It is worth noting that pruning removes weights which are not useful for **any** class, whereas distillation eliminates those neurons that are not beneficial for recognizing classes that belong to the **scope of interest of a particular IoT application**. Hence, distillation and pruning target different problems.

## 7.3 CNN Distillation Algorithm

*7.3.1 Neuroactivity Measurement.* For each member of training dataset that belongs to the classes of set $B$ (from Equation (1)), the neural network has to be run in inference mode. The output of each layer, which is an indicator of neurons' reaction to each particular input, has to be preserved. We use notation $OFM_{b_i,l}^j$ to show the responses of neurons that reside in layer $l$ of the neural network to $j^{th}$ element of set $Train_{b_i}$. The set $Train_{b_i}$ includes all of the samples of the training section of ILSVRC dataset [30] that belong to the class $b_i$. The cardinality of this set is 1300 for most classes.

*7.3.2 Heatmap Generation.* To detect all neurons which play a considerable role in recognizing members of a class, for every layer of the neural network, we need to compute the average neuron activities for each class $b_i$ using Equation (12). In this equation, $OFM_{b_i,l}$ indicates the average responses of neurons of layer $l$ to all training data of class $b_i$, and $OFM_{b_i,l}^j$ is computed in the last step.

$$OFM_{b_i,l} = \frac{\sum_{j=1}^{|Train_{b_i}|} OFM_{b_i,l}^j}{|Train_{b_i}|} \qquad (12)$$

*7.3.3 Bitmap Generation.* The heatmap which is generated in the previous step ($OFM_{b_i,l}$) shows the average reaction of each neuron to elements of $Train_{b_i}$. This value is zero for most of CDNs and can be negligible for some CANs. Using a threshold $TH_l$, it is possible to sort out dormant neurons and save their location in a bitmap as it is shown in Equation (13) where, $M$ is the number of OFMs, $N$ is the number of IFMs, $Win$ is the width of IFMs, and $Hin$ is the height of IFMs. In addition, bitmap is a 3D data structure that keeps track of neurons whose computations can be skipped.

$\forall\, m \in \mathbb{N}_M, \forall\, h \in \mathbb{N}_{Hin}, \forall\, w \in \mathbb{N}_{Win} :$

$$bitmap_{b_i,l}mhw = \begin{cases} 1 & OFM_{b_i,l}mhw \geq TH_l \\ 0 & \text{otherwise} \end{cases}$$
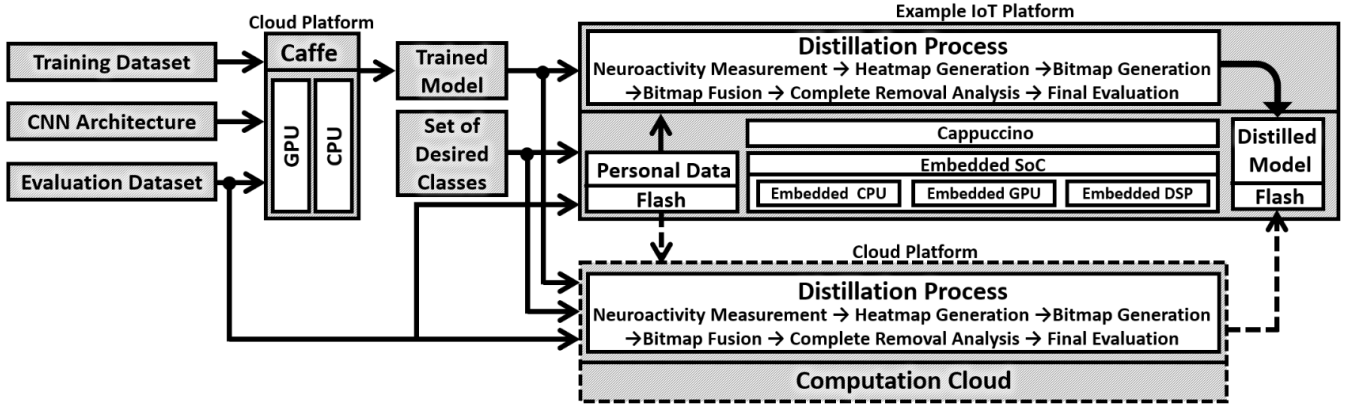
$$(13)$$

**Figure 6:** The distillation process takes a trained CNN and eliminates neurons which are responsible for understanding subjects that are out of an IoT application's scope of interest. The distillation process can be performed on a mobile device or a computation cloud.

*7.3.4 Bitmap Fusion.* The bitmaps which are computed in Equation (13) are class specific. That is, applying $bitmap_{b_i,l}$ to a CNN eliminates all of its sections that do not contribute in classifying members of class $b_i$. In order to successfully fulfill our purpose, i.e., be able to classify members of set $B$, all of the corresponding bitmaps must be fused. This can be achieved using the Hadamard Product as it is shown in Equation (14) in which the "$\sim$" symbol shows an element-wise inversion.

$$bitmap_l = \sim \sim bitmap_{b_1,l} \circ \sim bitmap_{b_2,l} \circ \cdots \circ \sim bitmap_{b_{|B|},l} \tag{14}$$

In Equation (14), $bitmap_l$ is a 3D data structure that shows neurons whose outputs are essential for understanding images that belong to classes of set $B$. Therefore, if $bitmap_l mhw = 0$, neuron $m$ of layer $l$ in location $h, w$ can be partially removed.

*7.3.5 Identifying Candidates for Complete Removal.* The structure of bitmaps is analogous to that of OFMs. That is, each layer of it includes the outputs of a single neuron. Therefore, if all elements of layer $m$ in a bitmap equal to zero, the corresponding neuron ($m$) can be completely removed. In other words, neuron $m$ can be completely removed if the result of Equation (15) equals to zero.

$$\sum_{h=1}^{Hout_l} \sum_{w=1}^{Wout_l} bitmap_l mhw \tag{15}$$

The distillation process is briefly illustrated in the block diagram of Figure 6.

## 7.4 The Required Time for Distilling a CNN

The CNN distillation algorithm is designed with users' privacy protection in mind. Therefore, we developed it to be reasonably fast and energy efficient in order to be used on edge devices. This enables users to distill a CNN using their personal data without being forced to upload the data to an external computation cloud. In general, three important factors should be considered when it comes to algorithm design for IoT platforms:

*7.4.1 Energy Consumption.* The distillation process has to be performed once. The results (i.e., neuron removal bitmaps) will be saved and can be used indefinitely. Given the offline nature of the process, it should be performed when the mobile platform is connected to an electrical grid. Hence, for this specific application, the impact of energy consumption is not as critical as other factors.

*7.4.2 Memory Footprint.* During the distillation, a mobile platform has to keep $l$ different OFMs in the memory. Given the average OFM size of $28 \times 28 \times 522$, and average length of 60 convolutional layers for deep CNNs, the required memory budget is 94 MB, which is affordable for most IoT devices that are qualified to host artificial intelligence algorithms.

*7.4.3 Execution Time.* Motamedi et al. offered Cappuccino, which is a platform for inference software synthesis for mobile system-on-chips [27]. Their approach requires 61.80 ms for each inference using AlexNet on a mediocre mobile platform (Google Nexus 6P). Since our distillation algorithm only requires the inference information, running it for a dataset of 10K private images on the aforementioned device requires 10.3 minutes of computation.

When privacy is not an issue, the neural network distillation can be performed on a powerful workstation using platforms such as Caffe [20] or TensorFlow [1]. Subsequently, the computed bitmaps which include guidance for neuron removal will be transferred to the target IoT device.

## 8 RESULTS AND DISCUSSIONS

## 8.1 Experimental Setup

In this section, we use GoogLeNet as a deep CNN to perform different experiments. GoogLeNet, which is trained on ILSVRC dataset [30], has the state of the art performance in classification and is currently being used by Google in their computation cloud [21]. In our experiments, we use the ILSVRC 2012 train dataset for network distillation process and take advantage of the validation part of the same dataset to measure the effect of neural network distillation on the

classification accuracy. The ILSVRC dataset contains images of both natural and human objects with labels that indicate the presence or absence of an object in a scene. The dataset serves as one of the standard benchmarks for which novel CNN architectures can evaluate their performance.

Currently, there are two major approaches for measuring the classification accuracy: top-1 accuracy rate and top-5 accuracy rate [30]. The former compares the ground truth against the prediction with the highest probability; however, the latter compares the ground truth against the top 5 predictions. The classification result is accepted if the ground truth exists among the top five predictions. It is difficult to achieve a high top-1 accuracy. Nonetheless, in most real-world situations, a CNN is expected to have one correct prediction. Hence, in all of our experiments, we use the top-1 accuracy rate. The classification accuracy is measured on the preserved classes.

Even though techniques such as neural network ensembling and excessive cropping can be used to improve the classification accuracy in competitions, in resource-constrained IoT devices performing such computations is infeasible or inefficient. Therefore, in all of our experiments we use a single crop per input frame and only query a single CNN. Distilling a CNN for set $B$ (Equation 1) requires $\sum_{i=1}^{\beta} |Train_{b_i}|$ training examples where $Train_{b_i}$ is the set of training data that belongs to class $b_i$. Since the cardinality of $Train_{b_i}$ is approximately identical for different classes, the required data for distillation is $\beta\alpha$ of the required data for training the CNN, where $\alpha$ is the number of all classes available in the dataset.

## 8.2  Case Study: Emergency Vehicle Recognition

*8.2.1  Application-Specific CNN Distillation.* GoogLeNet is trained to recognize classes that include emergency vehicles as well as many other classes which are not in the scope of interest of this particular task. The goal is to distill this CNN to remove those of its neurons that do not contribute to recognizing emergency vehicles.

If desired, it is always possible to lose some accuracy for achieving higher neuron removal rates. However, all distillation examples that are mentioned in this paper are performed without sacrificing accuracy. We distilled the neural network according to the algorithm which is described in subsection 7.3 to compute the heatmaps. The normalized histograms of heatmap values for the major layers of GoogLeNet are shown in Figure 7. The order of CNN layers is congruent with the one in Figure 3. As we expected, in the second convolutional layer, where most CANs that are structurally related to the input reside, the number of dormant neurons is very small. However, as we move towards the end of the neural network, this number starts to increase.

As it is highlighted in the last Inception layer (C3 in Figure 7), a relatively small number of neurons have very large output values in deep layers. Even though those neurons are small in count, they are important in determining the outcome of the classification.
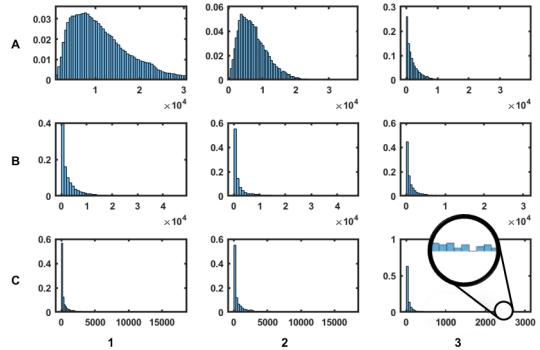


**Figure 7:** **Normalized histogram of heatmap values for the most important layers of GoogLeNet. The neural network is being distilled to mainly focus on recognizing emergency vehicles. The order of CNN layers is congruent with that of Figure 3.**

*8.2.2  Experimental Results.* Since a complete neuron removal in layer $l$ eliminates an entire OFM, the effect of such neuron removal on the required computation is identical to the impact of $Wout_l \times Hout_l$ partial neuron removal. Therefore, in the interest of simplicity, we count each complete neural removal as $Wout_l \times Hout_l$ partial neuron removal for the rest of the paper. Let us use Compression Ratio (CR) variable, which is defined in Equation (16), to demonstrate the amount of computations that could be skipped in distilling a CNN for a particular task. In Equation (16), $pr$ and $cr$ are the number of neurons that are partially removed and the number of neurons that are completely removed, respectively.

$$\mathrm{CR}_l = \frac{W_{out_l} \times H_{out_l} \times M_l}{W_{out_l} \times H_{out_l} \times M_l - cr_l \times W_{out_l} \times H_{out_l} + pr_l} \times 100 \tag{16}$$

It is worth noting that a decision on partial or complete removal of a neuron impacts a large chunk of computations. That is, for each and every single weight of that specific neuron, we skip the following operations: first, loading the weight data, second, computing the MAC result, and third, writing the output back.

Table 1 shows the compression ratio of each major layer of GoogLeNet for distilling it to another neural network which is capable of recognizing emergency vehicles only.

There are four notable points regarding the compression rates that are presented in Table 1:

(1) Starting from Inception layer 3B, the compression rates become considerable. A compression rate of 1.65X in that layer indicates that the distillation process eliminates 40% of unnecessary computations. This number increases semi-monotonically as we move towards deeper layer of the neural network. In the very last Inception layer, a compression rate of 4.71X means 78.78% of the unnecessary computations of that layer are removed.

(2) In Inception layer 3A, the compression rate is negligible. This is the phenomenon that we expected. That is, CANs which have a higher density in shallow layers cannot be removed.

8

**Table 1:** Compression ratios of different layers of GoogLeNet for distilling it to a smaller CNN that is capable of recognizing emergency vehicles. Numbers that are reported for each Inception layer are averaged over all convolutional layers that exist in that layer.

| I3A | I3B | I4A | I4B | I4C | I4D | I4E | I5A | I5B | Total MFLOPs [32] | MFLOPs Skipped | Speedup | Accuracy Loss |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-------------------|----------------|---------|---------------|
| 1.05X | 1.65X | 1.71X | 1.68X | 2.05X | 2.86X | 3.49X | 4.26X | 4.71X | **1107** | **539** | 1.95X | **0.0%** |

(3) The distillation is performed with classification accuracy preservation mindset. Therefore, the classification accuracy of the CNN on the classes that include emergency vehicles is not reduced after distillation. However, one can achieve higher compression rates, when small reduction of accuracy is tolerable.

(4) The distillation process, which only needs to be performed once per application, takes 39 minutes on a Qualcomm Snapdragon 810 SoC when the algorithm is accelerated using Cappuccino [27]. We also implemented the algorithm using Caffe [20]. The Caffe-based implementation finished the distillation process in 30.09 seconds using a Tesla K40 GPU.

## 8.3  GoogLeNet Distillation for Random Classes

Classes that contain emergency vehicles are semantically related. That is, all emergency vehicles are some sort of cars and share a number of feature-sets, such as wheels and headlights, that a CNN uses to recognize vehicles. This leads to a higher compression rate since classes that belong to our scope of interest do not demand a high variety of CDNs. To study the impact of distillation on subsets which includes classes with no particular semantical relationship, we used the distillation algorithm offered in Section 7.3 to distill GoogLeNet for three different subsets of ILSVRC datasets whose classes are chosen randomly. The first, second, and third subset contains 5, 10, and 15 random classes, respectively.

The results of GoogLeNet distillation for these classes are shown in Table 2. Increasing the diversity of elements of a subset for which a CNN is being distilled can have an adverse effect on the achievable compression ratio. Augmenting the cardinality of the subset of interest and handpicking the elements of that subset to have minimal semantical relationship with each other are both effective in increasing the diversity. In the experiments whose results are shown in Table 2, both approaches are used. Since for a subset with a high diversity a wide variety of neurons have to contribute to make the classification process successful, it is expected to observe smaller compression ratio as diversity increases. In an extreme case, when the size of the subset equals the number of classes that the CNN is originally trained on, we expect the compression rate to become one.

## 8.4  Distillation on CIFAR and MNIST Datasets

To further inspect the effectiveness of the proposed approach, we distilled two other CNNs: The first CNN is the TensorFlow [1] reference design for classifying CIFAR-10 [22] dataset and the second one is the reference design for classifying MNIST [25]. The CIFAR-10 dataset is a collection of color images which contains 50000 training and 10000 test

samples. The MNIST dataset includes handwritten images, and has a training set of 60000 samples, and a test set of 10000 samples.

In each experiment, the number of target classes is decreased, the CNN is distilled, and the effect on the parameter reduction is studied. In a CNN with 10 classes, one can hope that at most 10% of all parameters are exclusively used for understanding each class. Therefore, in an ideal setting, removing a class and distilling the CNN can at most eliminate 10% of the parameters. We use the ideal case as a reference to gauge the effectiveness of the distillation. Table 3 presents the distillation results on the aforementioned CNNs. The proposed algorithm can achieve near-linear parameter reduction specifically when a smaller number of classes are removed. In all experiments, the loss in classification accuracy is restricted to be at most 1%. The distillation is performed on the training subset of each dataset (50000 samples for CIFAR-10, and 60000 samples for MNIST). Subsequently, the distilled CNNs are tested on the evaluation subset of the datasets.

## 9  RELATED WORK

Deep CNNs, which involve many millions of parameters and billions of operations, are very compute-intensive. Such a huge computational demand is out of the budget for most of contemporary IoT platforms. Therefore, currently the mainstream approach for using artificial intelligence algorithms on such platforms is cloud-based computation.

To address this issue, the research community has put forth a number of ways to improve the computational complexity of such algorithms on IoT platforms. In what follows we briefly review those approaches and indicate where our proposal stands.

(1) **Bit-compression and Quantization:** It has been shown that inference using CNNs do not necessarily require 32 bits floating point operations [3, 10, 12, 18]. The process can be performed using 16 bits or less fixed point arithmetic. This is particularly very beneficial for FPGA-based implementation of CNNs where we have the luxury of tailoring the architecture to the algorithm's requirements.

(2) **Compact Model Design:** Instead of trying to compress cumbersome CNNs, it is more effective to develop a compact model in the first place. The expand-reduce model which is used in GoogLeNet [32] and SqueezeNet [19] is a notable example. The idea behind SqueezeNet [19] is used to develop CNNs that address other vision applications such as point cloud segmentation [33]. MobileNet [14] and SqueezeNext [9] are two other CNNs which are designed with a constrained parameter budget mindset. MobileNet [14] further improves AlexNet [23] in terms of parameter budget while keeping

**Table 2:** Compression of major layers of GoogLeNet for distilling it to smaller CNNs tailored to recognize members of three sets that we mentioned in Subsection 8.3. Numbers that are reported for each Inception Layer are averaged over all convolutional layers that it includes.

| | I3A | I3B | I4A | I4B | I4C | I4D | I4E | I5A | I5B | Total MFLOPs [32] | MFLOPs Skipped | Speedup | Accuracy Loss |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subset 3 | 1.01X | 1.43X | 1.44X | 1.13X | 1.50X | 1.80X | 1.46X | 1.38X | 1.98X | **1107** | **315** | 1.4X | **0.00%** |
| Subset 2 | 1.02X | 1.43X | 1.53X | 1.30X | 1.73X | 2.28X | 1.90X | 1.47X | 2.02X | **1107** | **382** | 1.5X | **0.20%** |
| Subset 1 | 1.03X | 1.46X | 2.09X | 1.34X | 1.86X | 2.83X | 3.25X | 2.24X | 4.35X | **1107** | **485** | 1.8X | **0.13%** |

**Table 3:** Distillation results on the TensorFlow reference designs for MNIST and CIFAR-10 datasets. In all experiments, the accuracy loss is projected to be less than 1.0%. The fine-tuning column shows the results that a fine-tuning-based approach would yield.

| CNN | Removed Classes (%) | Total KFLOPS | Skipped KFLOPS | Parameter Removal (%) | | | Accuracy Loss |
|---|---|---|---|---|---|---|---|
| | | | | Achieved | Fine-tuning | Ideal | |
| TF-MNIST | 10 | 18014 | 1801 | 10 | 0.008 | 10 | 0.02% |
| | 20 | 18014 | 3602 | 20 | 0.016 | 20 | 0.05% |
| | 30 | 18014 | 5403 | 30 | 0.024 | 30 | 0.06% |
| | 40 | 18014 | 7204 | 40 | 0.031 | 40 | 0.27% |
| | 50 | 18014 | 9555 | 49 | 0.039 | 50 | 0.97% |
| | 60 | 18014 | 9555 | 49 | 0.047 | 60 | 0.99% |
| | 70 | 18014 | 10808 | 57 | 0.055 | 70 | 0.98% |
| | 80 | 18014 | 13737 | 71 | 0.063 | 80 | 0.93% |
| TF-CIFAR | 10 | 25112 | 2511 | 10 | 0.022 | 10 | 0.01% |
| | 20 | 25112 | 5022 | 20 | 0.045 | 20 | 0.01% |
| | 30 | 25112 | 7534 | 30 | 0.064 | 30 | 0.0% |
| | 40 | 25112 | 10045 | 40 | 0.09 | 40 | 0.01% |
| | 50 | 25112 | 12556 | 50 | 0.112 | 50 | 0.06% |
| | 60 | 25112 | 15067 | 60 | 0.135 | 60 | 0.13% |
| | 70 | 25112 | 17277 | 68 | 0.157 | 70 | 0.48% |
| | 80 | 25112 | 19525 | 76 | 0.179 | 80 | 0.93% |

the same accuracy. SqueezeNext [9] further improves MobileNet [14] by decreasing the parameter budget by 13%. Some of the proposed compact CNNs use residual layers such as DenseNet [17] and CondenseNet [16]. Huang et al. further improved DenseNet [17] by introducing anytime classification [15].

(3) **FPGA-based Acceleration:** There are numerous proposals for FPGA-based CNN inference [8, 28, 29, 35]. However, there are two major issues constraining widespread usage of FPGA-based CNN acceleration. First, the time-consuming nature of hardware acceleration which increases the engineering costs. Second, the fact that not all IoT devices are equipped with an FPGA fabric.

(4) **Mobile SoC-based Acceleration:** Mobile SoC-based platforms are projected to be a major player in the emerging IoT landscape. Hence, an increasing number of solutions for accelerating CNNs on such platforms has been proposed [2, 24, 26, 27].

Our approach is different from the previous work in that we use an application-oriented approach to remove computations whose results will never be used. The CNN distillation algorithm that we proposed in this paper is orthogonal and complementary to the aforementioned approaches. We tested the proposed algorithm on a CNN with an expand-reduce architecture (GoogLeNet [32]) and used a mobile SoC-based inference software synthesizer to accelerate it (Cappuccino [27]). The underlying idea in distillation is drastically different from other approaches. A working CNN must possess kernels that generate distinctive features which make it possible to discriminate between different classes. Distill-Net takes advantage of this characteristic to detect those neurons that are not beneficial for recognizing classes that belong to the scope of interest of a particular IoT application. To the best of our knowledge, an approach for CNN mission reassignment is proposed in this article for the first time.

## 10 CONCLUSION

In this paper, we proposed a systematic approach for recognition and elimination of those sections of a neural network which will never be used in a particular application. First, we studied how different neurons of a CNN react to a stimulus and discussed how class-agnostics neurons and class-dependent neurons contribute to understanding an input. Subsequently, we offered a distillation algorithm that can recognize those neurons which are dormant for all instances that belong to a particular set of classes. Finally, we offered an approach for removing such neurons to avoid performing computations that are not necessary for a specific application. Our experimental results confirms that the approach strikes a favorable balance between classification accuracy and inference efficiency.

## REFERENCES

[1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. https://www.tensorflow.org/ Software available from tensorflow.org.
[2] Moustafa Alzantot, Yingnan Wang, Zhengshuang Ren, and Mani B Srivastava. 2017. Rstensorflow: Gpu enabled tensorflow for deep learning on commodity android devices. In *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*. ACM, 7–12.

**Table 4:** Achieved experimental speedups in running the distilled version of TF-MNIST on Google Nexus 6P. Each experiment has been repeated 10 times and the results are reported in the form of $a \pm s$, where $a$ is the average, and $s$ is the standard deviation.

| | Baseline | Removed Classes (%) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 |
| Execution Time (ms) | $213 \pm 5$ | $199 \pm 5$ | $180 \pm 4$ | $152 \pm 4$ | $132 \pm 1$ | $112 \pm 5$ | $87 \pm 2$ | $69 \pm 3$ | $50 \pm 2$ |
| Speedup | N/A | 1.07X | 1.18X | 1.40X | 1.61X | 1.90X | 2.45X | 3.09X | 4.26X |
| Ideal Speedup | N/A | 1.11X | 1.25X | 1.43X | 1.67X | 2.00X | 2.50X | 3.33X | 5.00X |

[3] Renzo Andri, Lukas Cavigelli, Davide Rossi, and Luca Benini. 2016. YodaNN: An ultra-low power convolutional neural network accelerator based on binary weights. In *VLSI (ISVLSI), 2016 IEEE Computer Society Annual Symposium on*. IEEE, 236–241.

[4] Luigi Atzori, Antonio Iera, and Giacomo Morabito. 2010. The internet of things: A survey. *Computer networks* 54, 15 (2010), 2787–2805.

[5] Debasis Bandyopadhyay and Jaydip Sen. 2011. Internet of things: Applications and challenges in technology and standardization. *Wireless Personal Communications* 58, 1 (2011), 49–69.

[6] Falak Shah Bhagyesh Vikani. 2017. CNN Visualization. https://github.com/InFoCusp/tf_cnnvis/.

[7] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. 2017. Learning Efficient Object Detection Models with Knowledge Distillation. In *Advances in Neural Information Processing Systems*. 742–751.

[8] Yu-Hsin Chen, Tushar Krishna, Joel S Emer, and Vivienne Sze. 2017. Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE Journal of Solid-State Circuits* 52, 1 (2017), 127–138.

[9] Amir Gholami, Kiseok Kwon, Bichen Wu, Zizheng Tai, Xiangyu Yue, Peter Jin, Sicheng Zhao, and Kurt Keutzer. 2018. SqueezeNext: Hardware-Aware Neural Network Design. *arXiv preprint arXiv:1803.10615* (2018).

[10] Philipp Gysel, Mohammad Motamedi, and Soheil Ghiasi. 2016. Hardware-oriented approximation of convolutional neural networks. *arXiv preprint arXiv:1604.03168* (2016).

[11] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. 2016. EIE: efficient inference engine on compressed deep neural network. In *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*. IEEE, 243–254.

[12] Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149* (2015).

[13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).

[14] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).

[15] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q Weinberger. 2017. Multi-scale dense convolutional networks for efficient prediction. *arXiv preprint arXiv:1703.09844* (2017).

[16] Gao Huang, Shichen Liu, Laurens van der Maaten, and Kilian Q Weinberger. 2017. CondenseNet: An Efficient DenseNet using Learned Group Convolutions. *arXiv preprint arXiv:1711.09224* (2017).

[17] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, Vol. 1. 3.

[18] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2016. Quantized neural networks: Training neural networks with low precision weights and activations. *arXiv preprint arXiv:1609.07061* (2016).

[19] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. 2016. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and $< 0.5$ MB model size. *arXiv preprint arXiv:1602.07360* (2016).

[20] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 675–678.

[21] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. 2017. In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*. ACM, 1–12.

[22] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. 2014. The CIFAR-10 dataset. *online: https://www.cs.toronto.edu/ kriz/cifar.html* (2014).

[23] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.

[24] Seyyed Salar Latifi Oskouei, Hossein Golestani, Matin Hashemi, and Soheil Ghiasi. 2016. Cnndroid: Gpu-accelerated execution of trained deep convolutional neural networks on android. In *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 1201–1205.

[25] Yann LeCun. 1998. The MNIST database of handwritten digits. *http://yann. lecun. com/exdb/mnist/* (1998).

[26] Mohammad Motamedi, Daniel Fong, and Soheil Ghiasi. 2017. Machine Intelligence on Resource-Constrained IoT Devices: The Case of Thread Granularity Optimization for CNN Inference. *ACM Transactions on Embedded Computing Systems (TECS)* 16, 5s (2017), 151.

[27] Mohammad Motamedi, Daniel Fong, and Soheil Ghiasi. 2018. Cappuccino: Efficient CNN Inference Software Synthesis for Mobile System-on-Chips. *IEEE Embedded Systems Letters* (2018).

[28] Kalin Ovtcharov, Olatunji Ruwase, Joo-Young Kim, Jeremy Fowers, Karin Strauss, and Eric S Chung. 2015. Accelerating deep convolutional neural networks using specialized hardware. *Microsoft Research Whitepaper* 2, 11 (2015).

[29] Jiantao Qiu, Jie Wang, Song Yao, Kaiyuan Guo, Boxun Li, Erjin Zhou, Jincheng Yu, Tianqi Tang, Ningyi Xu, Sen Song, et al. 2016. Going deeper with embedded fpga platform for convolutional neural network. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 26–35.

[30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* 115, 3 (2015), 211–252. https://doi.org/10.1007/s11263-015-0816-y

[31] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538* (2017).

[32] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich, et al. 2015. Going deeper with convolutions. Cvpr.

[33] Bichen Wu, Alvin Wan, Xiangyu Yue, and Kurt Keutzer. 2017. Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud. *arXiv preprint arXiv:1710.07368* (2017).

[34] Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*. Springer, 818–833.

[35] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. 2015. Optimizing fpga-based accelerator design for deep convolutional neural networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 161–170.