

Designing Energy-Efficient Convolutional Neural Networks using Energy-Aware Pruning

Tien-Ju Yang, Yu-Hsin Chen, Vivienne Sze
Massachusetts Institute of Technology
{tjy, yhchen, sze}@mit.edu

Abstract

Deep convolutional neural networks (CNNs) are indispensable to state-of-the-art computer vision algorithms. However, they are still rarely deployed on battery-powered mobile devices, such as smartphones and wearable gadgets, where vision algorithms can enable many revolutionary real-world applications. The key limiting factor is the high energy consumption of CNN processing due to its high computational complexity. While there are many previous efforts that try to reduce the CNN model size or amount of computation, we find that they do not necessarily result in lower energy consumption, and therefore do not serve as a good metric for energy cost estimation.

To close the gap between CNN design and energy consumption optimization, we propose an energy-aware pruning algorithm for CNNs that directly uses energy consumption estimation of a CNN to guide the pruning process. The energy estimation methodology uses parameters extrapolated from actual hardware measurements that target realistic battery-powered system setups. The proposed layer-by-layer pruning algorithm also prunes more aggressively than previously proposed pruning methods by minimizing the error in output feature maps instead of filter weights. For each layer, the weights are first pruned and then locally fine-tuned with a closed-form least-square solution to quickly restore the accuracy. After all layers are pruned, the entire network is further globally fine-tuned using back-propagation. With the proposed pruning method, the energy consumption of AlexNet and GoogLeNet are reduced by $3.7\times$ and $1.6\times$, respectively, with less than 1% top-5 accuracy loss. Finally, we show that pruning the AlexNet with a reduced number of target classes can greatly decrease the number of weights but the energy reduction is limited.

1. Introduction

In recent years, deep convolutional neural networks (CNNs) have become the state-of-the-art solution to many computer vision applications and are ripe for real-world deployment [1]. However, CNN processing incurs high en-

ergy consumption due to its high computational complexity [2]. As a result, battery-powered devices still cannot afford to run state-of-the-art CNNs due to their limited energy budget. For example, smartphones nowadays cannot even run object classification with AlexNet [3] in real-time for more than an hour. Hence, energy consumption has become the primary issue of bridging CNNs into practical computer vision applications.

In addition to accuracy, the design of modern CNNs also start to incorporate new metrics to make it more favorable in real-world environments. For example, the trend is to simultaneously reduce overall CNN model size and/or simplify the computation while going deeper. This is achieved either by pruning the weights of existing CNNs, i.e., making the filters sparse by setting some of the weights to zero [4–13], or by designing new CNNs with (1) highly bitwidth-reduced weights and arithmetic (e.g., XNOR-Net and BWN [14]) or (2) compact layers with limited number of weights (e.g., Network-in-Network [15], GoogLeNet [16], SqueezeNet [17], and ResNet [18]).

However, counting neither the number of weights nor the amount of computation, i.e., operations, in a CNN reveals its actual energy consumption. A CNN with smaller model size or less number of operations can still have higher overall energy consumption. This is because the sources of energy consumption in a CNN consist of not only computation but also memory accesses. In fact, fetching data from DRAM for an operation consumes *orders of magnitude higher energy* than the computation itself [19], and the energy consumption of a CNN is dominated by memory accesses for filter weights and feature maps. The total number of memory accesses is a function of the CNN shape configuration [20] (i.e., filter kernel size, feature map resolution, number of channels, and number of filters); different shape configurations can lead to different amounts of memory accesses, and thus energy consumption, even under the same number of weights or operations. Therefore, there is still no evidence showing that the aforementioned approaches can directly optimize the energy consumption of a CNN. In addition, there is currently no way for researchers to estimate

the energy consumption of a CNN at design time.

The key to closing the gap between CNN design and energy efficiency optimization is to directly use energy, instead of number of weights or operations, as a metric to guide the design. In order to obtain realistic energy consumption estimation at design time of the CNN, we use a framework proposed in [20] that models the two sources of energy consumption in a CNN: computation and memory accesses, and uses energy numbers extrapolated from actual hardware measurements. We then extend it to further model the impact of data sparsity and bitwidth reduction. The setup targets battery-powered platforms, such as smartphones and wearable devices, where hardware resources (i.e., computation and memory) are limited and energy efficiency is of utmost importance.

We further propose a new CNN pruning algorithm with the goal to minimize overall energy consumption with marginal accuracy degradation. Unlike the previous pruning methods, it directly minimizes the output feature map changes instead of the filter changes and achieves a higher compression ratio (i.e. the number of removed weights divided by the number of total weights). With the ability to directly estimate energy consumption of a CNN, the proposed pruning method identifies the parts of a CNN where pruning can maximally reduce the energy cost, and prunes the weights more aggressively than previously proposed methods to maximize the energy reduction.

In summary, the key contributions of this work include:

- **Energy Estimation Methodology:** Since the number of weights or operations does not necessarily serve as a good metric to guide the CNN design towards higher energy efficiency, we directly use energy consumption estimation of a CNN to guide its design. This energy estimation methodology is based on a framework proposed in [20] for realistic battery-powered systems, *e.g.*, smartphones, wearable devices, etc. We then further extend it to model the impact of data sparsity and bitwidth reduction. We will release it as a tool to facilitate future research on CNN design for real-world applications.
- **Energy-Aware Pruning:** We propose a new layer-by-layer pruning method that can aggressively reduce the number of non-zero weights by minimizing changes in feature maps as opposed to changes in filters. To maximize the energy reduction, the algorithm starts pruning the layers with the most energy consumption instead of with the largest number of weights, since pruning becomes more difficult as more layers are pruned. Each layer is first pruned and the preserved weights are locally fine-tuned with a closed-form least-square solution to quickly restore the accuracy and increase the compression ratio. After all the layers are pruned, the entire network is further globally fine-tuned by back-propagation. As a result, for AlexNet, we can reduce energy con-

sumption by $3.7\times$ after pruning, which is $1.7\times$ lower than pruning with the popular network pruning method proposed by [6]. Even for a compact CNN, such as GoogLeNet, the proposed pruning method can still reduce energy consumption by $1.6\times$. As many embedded applications only require a limited set of classes, we also show the impact of pruning AlexNet for a reduced number of target classes.

- **Energy Consumption Analysis of CNNs:** We evaluate the energy versus accuracy trade-off of widely-used or pruned CNN models. Our key insights are that (1) maximally reducing weights or the number of MACs in a CNN does not necessarily result in optimized energy consumption, (2) deeper CNNs with fewer weights, *e.g.*, GoogLeNet and SqueezeNet, do not necessarily consume less energy than shallower CNNs with more weights, *e.g.*, AlexNet, (3) convolutional (CONV) layers, instead of fully-connected (FC) layers, dominate the overall energy consumption in a CNN, and (4) sparsifying the filters can provide equal or more energy reduction than reducing the bitwidth (even to binary) of weights.

2. Energy Estimation Methodology

2.1. Background and Motivation

Multiply-and-accumulate (MAC) operations in the CONV and FC layers account for over 99% of total operations in state-of-the-art CNNs [3, 16, 18, 21], and therefore dominate both processing runtime and energy consumption. The energy consumption of MACs comes from computation and memory accesses for the required data, including both filter weights and feature maps. While the amount of computation increases linearly with the number of MACs, the amount of required data does not necessarily scale accordingly due to data reuse, i.e., the same data value is used for the computation of multiple MACs. This implies that some data have a higher impact on energy than others, since they are accessed more often. In other words, removing the data that is reused more has the potential to yield higher energy reduction.

Data reuse in a CNN arises in many ways, and is determined by the shape configurations of different layers. In CONV layers, due to its weight sharing property, each weight and input activation are reused many times according to the resolution of output feature maps and the kernel size of filters, respectively. In both CONV and FC layers, each input activation is also reused across all filters for different output channels within the same layer. When input batching is applied, each weight is further reused across all input feature maps in both types of layers. Overall, CONV layers usually present much more data reuse than FC layers. Therefore, as a general rule of thumb, each weight and activation in CONV layers have a higher impact on energy

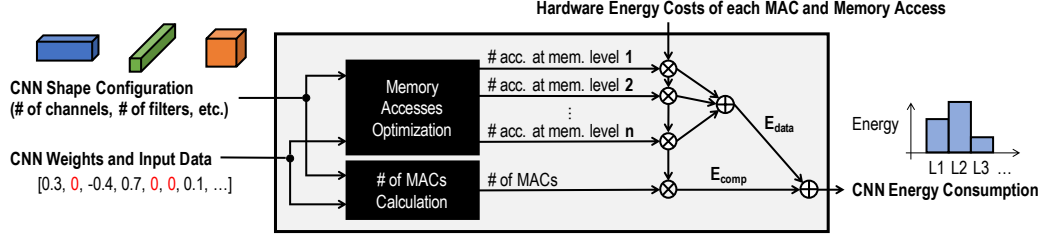


Figure 1. The energy estimation methodology is based on a framework proposed by [20], which can find out the optimized memory accesses at each level of the memory hierarchy that results in the lowest energy consumption. We then further account for the impact of data sparsity and bitwidth reduction, and use energy numbers extrapolated from actual hardware measurements to calculate the energy for both computation and data movement.

than in FC layers.

While data reuse serves as a good metric for comparing relative energy impact of data, it does not directly translate to the actual energy consumption. This is because modern hardware processors implement multiple levels of memory hierarchy, *e.g.*, DRAM and multi-level caches, to amortize the energy cost of memory accesses. The goal is to access data more from the less energy-consuming memory levels, which usually have less storage capacity, and thus minimize data accesses to the more energy-consuming memory levels. Therefore, the total energy cost to access a single piece of data with many reuses can vary a lot depending on how the accesses spread across different memory levels, and minimizing overall energy consumption using the memory hierarchy is the key to energy-efficient processing of CNNs.

2.2. Methodology

With the idea of exploiting data reuse in a multi-level memory hierarchy, Chen et al. [20] have presented a framework that can estimate the energy consumption of a CNN for inference. As shown in Fig 1, for each CNN layer, the framework calculates energy consumption by dividing it into two parts: computation energy consumption, E_{comp} , and data movement energy consumption, E_{data} . E_{comp} is calculated by counting the number of MACs in the layer and weighs it with the energy consumed by running each MAC operation in the computation core. E_{data} is calculated by counting the number of memory accesses at each level of the memory hierarchy presented in the hardware, and weighs it with the energy consumed by each access of that memory level. To obtain the number of memory accesses, [20] proposes an optimization procedure to search for the optimal number of accesses for all data types (feature maps and weights) at all levels of memory hierarchy that results in the lowest energy consumption. For energy numbers of each MAC operation and memory accesses, we use numbers extrapolated from actual hardware measurements targeting battery-powered platform setups [22].

Based on the aforementioned framework, we have created a methodology that further accounts for the impact of data sparsity and bitwidth reduction on energy consump-

tion. For example, we assume that the computation of a MAC and its associated memory accesses can be skipped completely when either of its input activation or weight is zero. Lossless data compression is also applied on the sparse data to save the cost of both on-chip and off-chip data movement. The impact of bitwidth is quantified by scaling the energy cost of different hardware components accordingly. For instance, the multiplier energy consumption scales with bitwidth quadratically, while memory only scales its energy linearly.

2.3. Potential Impact

With this methodology, we are able to quantify the difference in energy cost between various previously proposed CNN models and methods, such as increasing data sparsity or aggressive bitwidth reduction. More importantly, it provides a gateway for researchers to assess the energy of a CNN at design time, which can be used as a feedback that leads to CNN designs with significantly reduced energy consumption. In Sec. 4, we will describe an energy-aware pruning method that uses the proposed energy estimation methodology for deciding the pruning priority.

3. CNN Pruning: Related Work

Weight pruning. There is a large body of work that aims to reduce CNN model size by pruning while maintaining accuracy. LeCun et al. [4] and Hassibi et al. [7] remove the weights based on the sensitivity of the final objective function to that weight (*i.e.*, remove the weights of least sensitivity first). Han et al. [5, 6] ignores the output sensitivity and uses a magnitude-based method, which removes the small-magnitude weights first. Jin et al. [8] and Gao et al. [9] extend the magnitude-based method to allow the restoration of the pruned weights in the previous iterations, with tightly coupled pruning and retraining stages, for greater model compression. However, all the above methods evaluate whether to prune each weight independently and do not account for correlation between weights [10]. When the compression ratio is large, the aggregate impact of many weights can have a large impact on the output; thus,

failing to considering the combined influence of the weights on the output limits the achievable compression ratio.

Filter pruning. Rather than investigating the removal of each individual weight (fine-grained pruning), there is also work that investigates removing entire filters (coarse-grained pruning). Hu et al. [11] proposed removing filters that frequently generate zero outputs after the ReLU layer in the validation set. Srinivas et al. [12] proposed merging similar filters into one. Mariet et al. [13] proposed merging filters with similar output activations into one, but the method only works on fully-connected layers. Unfortunately, these coarse-grained pruning approaches tend to have lower compression ratios than fine-grained pruning for the same accuracy.

Previous work directly targets reducing model size rather than energy consumption. However, as discussed in Sec. 1, the number of weights alone does not dictate the energy consumption. Hence, the energy consumption of their pruned CNNs in the previous work is not minimized.

To address issues highlighted above, we propose a new fine-grained pruning algorithm that specifically targets energy-efficiency. It utilizes the estimated energy provided by the methodology described in Sec. 2 to guide the proposed pruning algorithm to aggressively prune the layers with the highest energy consumption with marginal impact on accuracy. Moreover, the pruning algorithm considers the joint influence of weights on the final output feature map, thus enabling both a higher compression ratio and a larger energy reduction. The combination of these two approaches result in CNNs that are more energy-efficient and compact than previously proposed approaches.

The proposed energy-efficient pruning algorithm can be combined with other techniques to further reduce the energy consumption, such as bitwidth reduction of weights or feature maps [14, 23, 24], weight sharing and Huffman coding [6], student-teacher learning [25], filter decomposition [26, 27] and pruning feature maps [28].

4. Energy-Aware Pruning

Our goal is to reduce the energy consumption of a given CNN by sparsifying the filters without significant impact on the network accuracy. The key steps in the proposed energy-aware pruning are shown in Fig. 2, where the input is a pre-trained model and the output is a sparser model with lower energy consumption.

In **Step 1**, the pruning order of the layers is determined based on the energy as described in Sec. 2. Step 2, 3 and 4 removes, restores and fine-tunes weights, respectively, for each layer and this inner loop is repeated for each layer in the network. *Pruning* and *restoring* weights involve choosing weights, while *fine-tuning* weights involves changing the value of the weights, all while minimizing output feature map error. In **Step 2**, a simple magnitude-based prun-

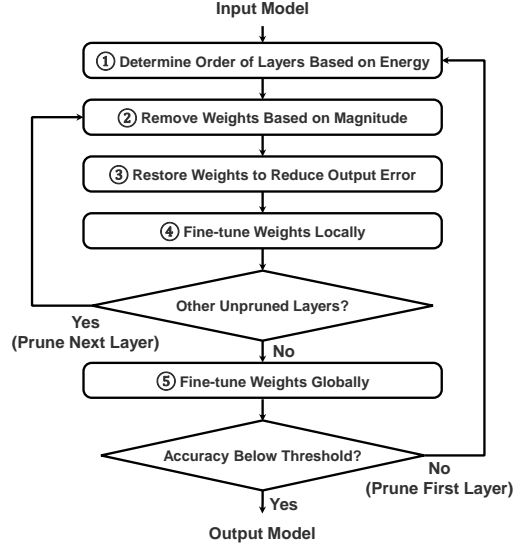


Figure 2. Flow of energy-aware pruning.

ing method is used to aggressively and quickly remove the weights above the target compression ratio (e.g., if the target compression ratio is 30%, 35% of the weights are removed in this step). The number of extra weights removed is determined empirically. In **Step 3**, the correlated weights that have the greatest impact on reducing output error are restored to their original non-zero value to reach the target compression ratio (e.g., restore 5% of weights). In **Step 4**, the preserved weights are locally fine-tuned with a closed-form least-square solution to further decrease the output feature map error. Each of these steps are described in detail in Sec. 4.2, Sec. 4.3 and, Sec. 4.4.

Once each individual layer has been pruned using Step 2 to 4, **Step 5** performs a global fine-tuning of weights across the entire network using back-propagation as described in Sec. 4.5. All these steps are iteratively performed until the final network can no longer maintain a given accuracy, e.g., 1% accuracy loss.

Compared to the previous magnitude-based approaches [5, 6, 8, 9], the main difference of this work is the introduction of Steps 1, 3, and 4. Step 1 enables pruning to minimize the energy consumption. Step 3 and 4 increase the compression ratio to further decrease the energy consumption.

4.1. Determine Order of Layers based on Energy

As more layers are pruned, it becomes increasingly difficult to remove weights because the accuracy approaches the given accuracy threshold. Accordingly, layers that are pruned early on tend to have a higher compression ratio than the layers that follow. Thus, in order to maximize the overall energy reduction, we prune the layers that consume the most energy first. Specifically, we used the energy estima-

tion from Sec. 2 and determine the pruning order of layers based on their energy consumption. As a result, the layers that consume the most energy achieve higher compression ratios and energy reduction. In the beginning of each outer loop iteration in Fig. 2, the new pruning order is redetermined according to the new energy estimation of each layer.

4.2. Prune Weights based on Magnitude

For a fully-connected layer, $Y_i \in R^{k \times 1}$ is the i_{th} output feature map across k images and is computed from

$$Y_i = X_i A_i + B_i \mathbf{1}, \quad (1)$$

where $A_i \in R^{m \times 1}$ is the i_{th} filter among all n filters with m weights, and $X_i \in R^{k \times m}$ denotes the corresponding k input feature maps, $B_i \in R$ is the bias, and $\mathbf{1} \in R^{k \times 1}$ is a vector where all entries are one. For a convolutional layer, we can convert the convolutional operation into matrix multiplication operations, by converting the input feature maps into a Toeplitz matrix, and compute the output feature maps with a similar equation as Eq. 1.

To sparsify the filters without impacting the accuracy, the simplest method is pruning weights with magnitudes smaller than a threshold, which is referred to as magnitude-based pruning [5, 6]. The advantage of this approach is that it is fast, and works well when a small amount of weights is removed and thus the correlation between weights only have a minor impact on the output. However, as more weights are pruned, this method introduces large output errors as the correlation between weights become more critical. For example, if most of the small-magnitude weights are negative, the output feature map error will become large once lots of these small negative weights are removed using the magnitude-based pruning. It would be desirable to remove a large positive weight to compensate for the introduced error instead of removing more smaller negative weights. Thus, we only use magnitude-based pruning for fast initial pruning of each layer. We then introduce additional steps that account for the correlation between weights to reduce the output error due to magnitude-based pruning.

4.3. Restore Weights to Reduce Output Error

It is the error in the output feature map, and not the filter, that affects the overall network accuracy. Therefore, we focus on minimizing the error of the output feature maps instead of that of the filters. To achieve this, we model the problem as the following ℓ_0 -minimization problem:

$$\begin{aligned} \bar{A}_i &= \arg \min_{\bar{A}_i} \left\| \hat{Y}_i - X_i \bar{A}_i \right\|_F^2, \\ \text{subject to } \left\| \bar{A}_i \right\|_0 &\leq q, \quad i = 1, \dots, n, \end{aligned} \quad (2)$$

where \hat{Y}_i denotes $Y_i - B_i \mathbf{1}$ and q is the number of non-zero weights we want to retain in all filters. Unfortunately, solving this ℓ_0 -minimization problem is NP-hard. Therefore, a greedy algorithm is proposed to approximate it.

The algorithm starts from pruned filters $\check{A} \in R^{m \times n}$, obtained from the magnitude-based pruning in Step 2. These filters are overpruned, meaning they are pruned at a higher compression ratio than the target compression ratio. These filters have the support S , where S is a set of the indices of non-zero weights in the filters. It then iteratively restores weights until the number of non-zero weights is equal to q , which reflects the target compression ratio.

The residual of each filter, which indicates the current output feature map difference we need to minimize, is initialized as $Y_i - B_i \mathbf{1} - \check{A}_i$. In each iteration, out of the weights not in the support of a given filter S_i , we select the weight that reduces the ℓ_1 -norm of the corresponding residual the most, and adds it to the support S_i . The residual then is updated by taking this new weight into account.

We restore weights from the filter with the largest residual in each iteration. This prevents the algorithm from restoring weights in filters with small residual, which will likely have less effect on the overall output feature map error. This could occur if the weights were selected based solely on the largest ℓ_1 -norm improvement for any filter.

To speed up this restoration process, we restore multiple weights within a given filter in each iteration. The p weights with the top- p maximum ℓ_1 -norm improvement are chosen. As a result, we reduce the frequency of computing residual improvement for each weight, which takes a significant amount of time. We adopt p equal to 2 in our experiments, but a higher p can be used.

4.4. Fine-tune Weights Locally

The previous two steps select a subset of weights to preserve, but do not change the values of the weights. In this step, we perform the least-square optimization on each filter to change the values of these weights to further reduce the output error and restore network accuracy:

$$\bar{A}_{i,S_i} = \arg \min_{\bar{A}_{i,S_i}} \left\| \hat{Y}_i - X_{i,S_i} \bar{A}_{i,S_i} \right\|_F^2, \quad \bar{A}_{i,S_i^c} = 0, \quad (3)$$

where the subscript S_i means choosing the non-pruned weights from the i_{th} filter and the corresponding columns from X_i . The least-square problem has a closed-form solution, which can be solved efficiently.

4.5. Fine-tune Weights Globally

After all the layers are pruned, we fine-tune the whole network using back-propagation with the pruned weights fixed at zero. This step can be used to globally fine-tune the weights to achieve higher accuracy. Fine-tuning the whole network is time-consuming and requires careful tuning of several hyper-parameters. In addition, back-propagation can only restore accuracy within certain accuracy loss. However, since we first fine-tune weights locally, part of the accuracy has already been restored, which enables more weights to be pruned under a given accuracy

Table 2. Compression ratio¹ of each layer in AlexNet.

# of Classes	[6]	This Work			
	1000	1000	100	10 (Random)	10 (Dog)
CONV1	16%	83%	86%	89%	89%
CONV2	62%	92%	97%	97%	96%
CONV3	65%	91%	97%	98%	97%
CONV4	63%	81%	88%	97%	95%
CONV5	63%	74%	79%	98%	98%
FC1	91%	92%	93%	~100%	~100%
FC2	91%	91%	94%	~100%	~100%
FC3	74%	78%	78%	~100%	~100%

¹ The number of removed weights divided by the number of total weights. The higher, the better.

loss tolerance. As a result, we increase the compression ratio in each iteration, reducing the total number of retraining iterations and overall retraining time.

5. Experiment Results

5.1. Pruning Method Evaluation

We evaluate the effectiveness of our energy-aware pruning on AlexNet [3], GoogLeNet v1 [16] and SqueezeNet v1 [17] and compare it with the state-of-the-art magnitude-based pruning method with the publicly available models [6].¹ The accuracy and energy consumption are measured on the ImageNet ILSVRC 2014 dataset [29]. Since the energy-aware pruning method relies on the output feature maps, we use the training images for both pruning and retraining. All accuracy numbers are measured on the validation images. To estimate the energy consumption with the proposed methodology in Sec. 2, we assume all values are represented as 16-bit numbers, except where otherwise specified, to fairly compare the energy consumption of networks. The hardware parameters used are similar to [22].

Table 1 summarizes the results.² The batch size is 44 for AlexNet and 48 for other two networks. All the energy-aware pruned networks have less than 1% accuracy loss with respect to the other corresponding networks. For AlexNet and SqueezeNet, our method achieves better results in all metrics (i.e., number of weights, number of MACs, and energy consumption) than magnitude-based pruning [6]. For example, the number of MACs is reduced by another 3.2 \times and the estimated energy is reduced by another 1.7 \times with a 15% smaller model size on AlexNet. Table 2 shows a comparison of the energy-aware pruning and the magnitude-based pruning across each layer; our method gives a higher compression ratio for all layers, especially for

¹The proposed energy-aware pruning can be easily combined with other techniques in [6], such as weight sharing and Huffman coding.

²We use the models provided by MatConvNet [30] or converted from Caffe [31] or Torch [32], so the accuracies may be slightly different from that reported by other works.

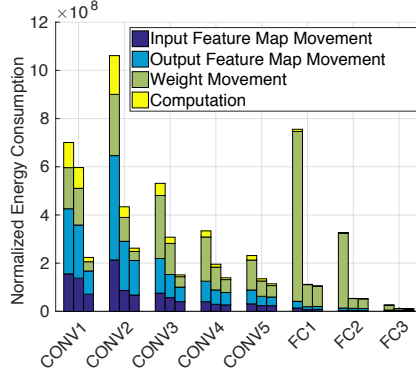


Figure 3. Energy consumption breakdown of different AlexNets in terms of the computation and the data movement of input feature maps, output feature maps and filter weights. From left to right: original AlexNet, AlexNet pruned by [6], AlexNet pruned by energy-aware pruning.

CONV1 to CONV3, which consume most of the energy.

Our approach is also effective on compact models. For example, on GoogLeNet, the achieved reduction factor is 2.9 \times for the model size, 3.4 \times for the number of MACs and 1.6 \times for the estimated energy consumption.

5.2. Energy Consumption Analysis

We also evaluate the energy consumption of popular CNNs. In Fig. 4, we summarize the estimated energy consumption of CNNs relative to their top-5 accuracy. The result reveals the following key observations:

- **Convolutional layers consume more energy than fully-connected layers.** Fig. 3 shows the energy breakdown of the original AlexNet and two pruned AlexNet models. Although most of the weights are in the fully-connected (FC) layers, convolutional (CONV) layers account for most of the energy consumption. Take the original AlexNet as an example. All CONV layers contain 3.8% of the total weights, but consume 72.6% of the total energy. (1) In CONV layers, the energy consumption of the input and output feature maps is much higher than that of FC layers. Compared to FC layers, CONV layers require a larger number of MACs, which involve loading inputs from the storage and writing the outputs to the storage. Accordingly, a large number of MACs lead to a large amount of weight and feature map movement and hence high energy consumption. (2) The weight-related energy consumption of all CONV layers is similar to that of all FC layers. Each weight in CONV layers are used more frequently than that in FC layers, which causes the small weight-related energy consumption difference. Obviously, the number of weights is not an accurate estimator for energy consumption and pruning a weight from CONV layers contributes more to energy reduction than pruning a weight from FC layers. As a network goes deeper, *e.g.*, ResNet [18], the CONV layers

Table 1. Performance metrics of various dense and pruned models.

Model		Top-5 Accuracy	# of Non-zero Weights ($\times 10^6$)	# of Non-skipped MACs ($\times 10^8$) ¹	Normalized Energy ($\times 10^9$) ¹
AlexNet	(Original)	80.43%	60.90 (100%)	3.71 (100%)	3.97 (100%)
AlexNet	([6])	80.37%	6.79 (11%)	1.79 (48%)	1.85 (47%)
AlexNet	(Energy-Aware Pruning)	79.56%	5.73 (9%)	0.56 (15%)	1.06 (27%)
GoogLeNet	(Original)	88.26%	7.00 (100%)	7.41 (100%)	7.63 (100%)
GoogLeNet	(Energy-Aware Pruning)	87.28%	2.37 (34%)	2.16 (29%)	4.76 (62%)
SqueezeNet	(Original)	80.61%	1.24 (100%)	4.51 (100%)	5.28 (100%)
SqueezeNet	([6])	81.47%	0.42 (33%)	3.30 (73%)	4.61 (87%)
SqueezeNet	(Energy-Aware Pruning)	80.47%	0.35 (28%)	1.93 (43%)	3.99 (76%)

¹ The estimated energy is measured per image and normalized to the energy of a MAC operation with two 16-bit inputs.

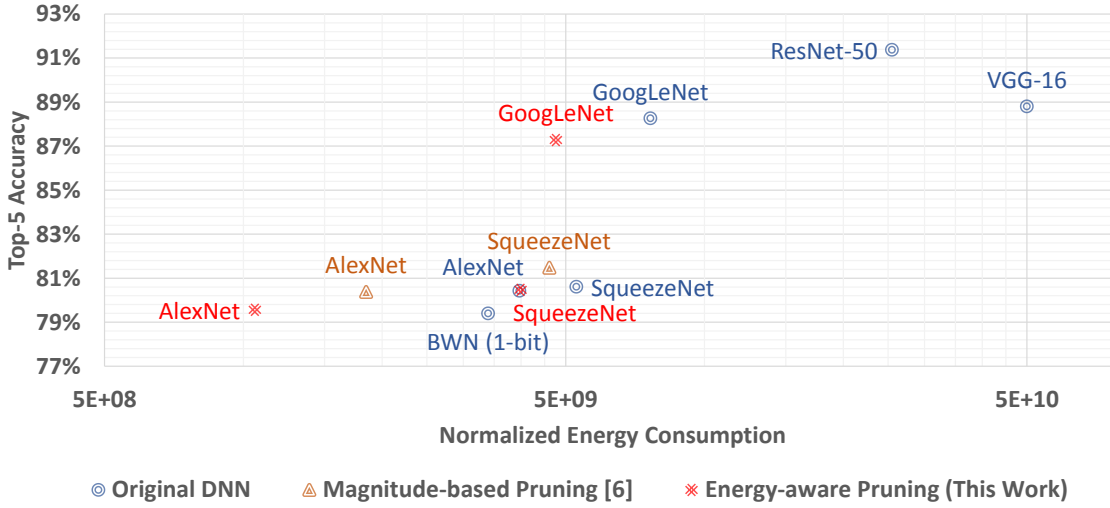


Figure 4. Accuracy versus energy trade-off of popular CNN models. Models pruned with the energy-aware pruning provide a better accuracy vs. energy trade-off (steeper slope).

dominate both energy consumption and the model size. By pruning, the energy consumption of all the three data types are reduced because of skipping computation. The energy-aware pruning prunes CONV layers effectively, so significant energy reduction is achieved.

- **Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights.** One network design strategy for reducing the size of a network without sacrificing the accuracy is to make a network thinner but deeper. However, does this mean the energy consumption is also reduced? From Table 1, a network architecture having a smaller model size does not necessarily have lower energy consumption. For instance, SqueezeNet is a compact model and a good fit for memory-limited applications. By making a network thinner but deeper than AlexNet, it achieves a similar accuracy but with $49\times$ size reduction. However, as the network goes deeper, more CONV layers are used. To preserve the accuracy, the size of the feature maps cannot be reduced significantly until the final few layers. Hence, the newly added CONV layers involve a large amount of data movement, resulting in higher energy consumption.

- **Reducing the number of weights can provide lower energy consumption than reducing the bitwidth of weights.** From Fig. 4, the AlexNet pruned by the proposed method consumes less energy than BWN. BWN uses an AlexNet-like architecture with binarized weights. Through binarizing the weights, only the weight-related and computation-related energy consumption is reduced. However, pruning can achieve extra energy reduction on moving feature maps. Since the weights in CONV1 and FC3 of BWN are not binarized to preserve the accuracy, BWN does not reduce their energy consumption. Moreover, to compensate for the accuracy loss of binarizing weights, CONV2, CONV4 and CONV5 layers in BWN use $2\times$ the number of weights in the corresponding layers of the pruned AlexNet, which leads to higher energy consumption for feature map movement.
- **A lower number of MACs does not necessarily lead to lower energy consumption.** For example, the pruned GoogleNet has a fewer MACs but consumes more energy than the SqueezeNet pruned by [6]. That is because they have different data reuse, which is determined by the shape configurations, as discussed in Sec. 2.1.

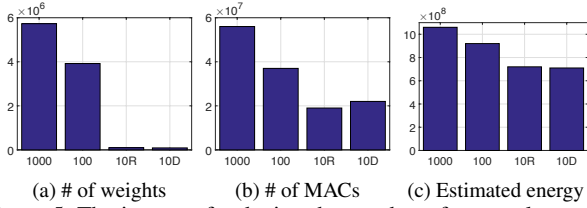


Figure 5. The impact of reducing the number of target classes on the three metrics. The x-axis is the number of target classes. 10R and 10D denote the 10-random-class model and the 10-dog-class model, respectively.

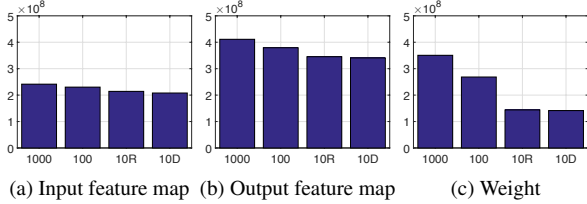


Figure 6. The energy breakdown of models with different numbers of target classes.

From Fig. 4, we also observe that the energy consumption scales exponentially with linear increase in accuracy. For instance, GoogLeNet consumes $2\times$ energy of AlexNet for 8% accuracy improvement, and ResNet-50 consumes $3.3\times$ energy of GoogLeNet for 3% accuracy improvement.

In summary, the model size (i.e., the number of weights \times bitwidth) and the number of MACs do not directly reflect the energy consumption of a layer or a network. There are other factors like the data movement of the feature maps, which are often overlooked. Therefore, with the proposed energy estimation methodology, researchers can have a clearer view of CNNs and more effectively design low-energy-consumption networks.

5.3. Number of Target Class Reduction

In many applications, the number of classes can be significantly less than 1000. We study the influence of reducing the number of target classes by pruning on model size, number of MACs, and energy consumption. AlexNet is used as the starting point. The number of target classes is reduced from 1000 to 100 to 10. The target classes of the 100-class model and one of the 10-class model are randomly picked, and that of another 10-class model are different dog breeds. These models are pruned with less than 1% top-5 accuracy loss for the 100-class model and less than 1% top-1 accuracy loss for the two 10-class models.

From Fig. 5, as the number of target classes reduces, the number of weights and MACs and the estimated energy consumption decrease. However, they reduce at different rates with the model size dropping the fastest, followed by the number of the MACs the second, and the estimated energy reduces the slowest.

According to Table 2, for the 10-class models, almost

all the weights in the FC layers are pruned, which leads to a very small model size. Because the FC layers work as classifiers, most of the weights that are responsible for classifying the removed classes are pruned. The higher-level layers, such as CONV4 and CONV5, which contain filters for extracting more specialized features of objects, are also significantly pruned. CONV1 is pruned less since it extracts basic features that are shared among all classes. As a result, the number of MACs and energy consumption do not reduce as rapidly as the number of weights. Thus, we hypothesize that the layers closer to the output of the networks shrink more rapidly with the number of classes.

As the number of classes reduces, the energy consumption becomes less sensitive to weight sparsity. From the energy breakdown (Fig. 6), the energy consumption of feature maps gradually saturates due to data reuse and the memory hierarchy. For example, each time one input activation is loaded from DRAM into the chip, it is used multiple times by several weights. If any one of these weights is not pruned, the activation still needs to be fetched from DRAM. Moreover, we observed that sometimes the sparsity of feature maps decreases after we reduce the number of target classes, which causes higher energy consumption for moving the feature maps.

From Table 2 and Fig. 5 and 6, the compression ratio and the performance of the two 10-class models are similar. Hence, we hypothesize that the pruning performance mainly depends on the number of target classes and the type of the preserved classes is less influential.

6. Conclusion

This work presents an energy-aware pruning algorithm that directly uses energy estimation of a CNN to guide the pruning process in order to optimize for the best energy-efficiency. The energy of a CNN is estimated by a methodology that models the computation and memory accesses of a CNN and uses energy numbers extrapolated from actual hardware measurements. It enables more accurate energy consumption estimation compared to just using the model size or the number of MACs. With the estimated energy for each layer in a CNN model, the algorithm performs layer-by-layer pruning, starting from the layers with the highest energy consumption to the layers with the lowest energy consumption. For pruning each layer, it removes the weights which have the smallest joint impact on the output feature map. The experiments show that the proposed pruning method reduces the energy consumption of AlexNet and GoogLeNet, by $3.7\times$ and $1.6\times$, respectively, compared to their original dense models. The influence of pruning the AlexNet with the number of target classes reduced is explored and discussed. The results show that by reducing the number of target classes, the model size can be greatly reduced but the energy reduction is limited.

References

- [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, pp. 436–444, May 2015.
- [2] “GPU-Based Deep Learning Inference: A Performance and Power Analysis.” Nvidia Whitepaper, 2015.
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks,” in *NIPS*, pp. 1097–1105, 2012.
- [4] Y. LeCun, J. S. Denker, and S. A. Solla, “Optimal Brain Damage,” in *NIPS*, pp. 598–605, 1990.
- [5] S. Han, J. Pool, J. Tran, and W. J. Dally, “Learning both weights and connections for efficient neural networks,” in *NIPS*, pp. 1135–1143, 2015.
- [6] S. Han, H. Mao, and W. J. Dally, “Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding,” in *ICLR*, 2016.
- [7] B. Hassibi and D. G. Stork, “Second order derivatives for network pruning: Optimal brain surgeon,” in *NIPS*, pp. 164–171, 1993.
- [8] X. Jin, X. Yuan, J. Feng, and S. Yan, “Training Skinny Deep Neural Networks with Iterative Hard Thresholding Methods,” *arXiv preprint arXiv:1607.05423*, 2016.
- [9] Y. Guo, A. Yao, and Y. Chen, “Dynamic Network Surgery for Efficient DNNs,” in *NIPS*, 2016.
- [10] R. Reed, “Pruning algorithms - a survey,” *IEEE Transactions on Neural Networks*, vol. 4, no. 5, pp. 740–747, 1993.
- [11] H. Hu, R. Peng, Y.-W. Tai, and C.-K. Tang, “Network Trimming: A Data-Driven Neuron Pruning Approach towards Efficient Deep Architectures,” *arXiv preprint arXiv:1607.03250*, 2016.
- [12] S. Srinivas and R. V. Babu, “Data-free parameter pruning for Deep Neural Networks,” in *BMVC*, 2015.
- [13] Z. Mariet and S. Sra, “Diversity Networks,” in *ICLR*, pp. 1–11, 2016.
- [14] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks,” in *ECCV*, 2016.
- [15] M. Lin, Q. Chen, and S. Yan, “Network in Network,” in *ICLR*, 2014.
- [16] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going Deeper With Convolutions,” in *CVPR*, pp. 1–9, 2015.
- [17] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, “SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size,” *arXiv preprint arXiv:1602.07360*, 2016.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *CVPR*, 2016.
- [19] M. Horowitz, “Computing’s energy problem (and what we can do about it),” in *ISSCC*, 2014.
- [20] Y. Chen, J. Emer, and V. Sze, “Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks,” in *ISCA*, pp. 367–379, 2016.
- [21] K. Simonyan and A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” in *ICLR*, 2014.
- [22] Y. Chen, T. Krishna, J. Emer, and V. Sze, “Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks,” in *ISSCC*, pp. 262–263, 2016.
- [23] M. Courbariaux, Y. Bengio, and J.-P. David, “Binaryconnect: Training deep neural networks with binary weights during propagations,” in *NIPS*, 2015.
- [24] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, “Quantized Convolutional Neural Networks for Mobile Devices,” in *CVPR*, 2016.
- [25] J. Ba and R. Caruana, “Do deep nets really need to be deep?,” in *NIPS*, 2014.
- [26] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin, “Compression of Deep Convolutional Neural Networks for Fast and Low Power Mobile Applications,” *ICLR*, 2016.
- [27] H. Foroosh, M. Tappen, and M. Pensky, “Sparse Convolutional Neural Networks,” in *CVPR*, 2015.
- [28] B. Reagen, P. Whatmough, R. Adolf, S. Rama, H. Lee, S. Kyu, L. José, G.-y. W. D. Brooks, and W. Power, “Minerva : Enabling Low-Power, Highly-Accurate Deep Neural Network Accelerators,” in *ISCA*, 2016.
- [29] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.
- [30] A. Vedaldi and K. Lenc, “MatConvNet – Convolutional Neural Networks for MATLAB,” in *Proceeding of the ACM Int. Conf. on Multimedia*, 2015.
- [31] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional Architecture for Fast Feature Embedding,” *arXiv preprint arXiv:1408.5093*, 2014.
- [32] R. Collobert, K. Kavukcuoglu, and C. Farabet, “Torch7: A matlab-like environment for machine learning,” in *BigLearn, NIPS Workshop*, 2011.