

A PROGRAMMABLE APPROACH TO NEURAL NETWORK COMPRESSION

Vinu Joseph¹ Saurav Muralidharan² Animesh Garg³ Michael Garland² Ganesh Gopalakrishnan¹

ABSTRACT

Deep neural networks (DNNs) frequently contain far more weights, represented at a higher precision, than are required for the specific task which they are trained to perform. Consequently, they can often be compressed using techniques such as weight pruning and quantization that reduce both the model size and inference time without appreciable loss in accuracy. However, finding the best compression strategy and corresponding target sparsity for a given DNN, hardware platform, and optimization objective currently requires expensive, frequently manual, trial-and-error experimentation. In this paper, we introduce a programmable system for model compression called **Condensa**. Users programmatically compose simple operators, in Python, to build more complex and practically interesting compression strategies. Given a strategy and user-provided objective (such as minimization of running time), Condensa uses a novel Bayesian optimization-based algorithm to automatically infer desirable sparsities. Our experiments on four real-world DNNs demonstrate memory footprint and hardware runtime throughput improvements of 188x and 2.59x, respectively, using at most ten samples per search. We have released a reference implementation of Condensa at <https://github.com/NVlabs/condensa>.

用于减少推理延迟的推荐压缩策略可能与减少总内存占用所需的压缩策略不同

1 INTRODUCTION

Modern deep neural networks (DNNs) are complex, and often contain millions of parameters spanning dozens or even hundreds of layers (He et al., 2016; Huang et al., 2017). This complexity translates into substantial memory and runtime costs on hardware platforms at all scales. Recent work has demonstrated that DNNs are often over-provisioned and can be compressed without appreciable loss of accuracy. Model compression can be used to reduce both model memory footprint and inference latency using techniques such as weight pruning (Han et al., 2015; Luo et al., 2017), quantization (Gupta et al., 2015), and low-rank factorization (Jaderberg et al., 2014; Denton et al., 2014). Unfortunately, the requirements of different *compression contexts*—DNN structure, target hardware platform, and the user’s optimization objective—are often in conflict. The recommended compression strategy for reducing inference latency may be different from that required to reduce total memory footprint. For example, in a Convolutional Neural Network (CNN), reducing inference latency may require pruning filters to realize speedups on real hardware (Jitzi et al., 2016),

在真实的硬件上实现加速推理并且减少推理延迟，需要修剪CNN的filter

¹University of Utah ²NVIDIA ³University of Toronto. Correspondence to: Vinu Joseph <vinu@cs.utah.edu>.

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). The views, opinions and/or findings expressed are those of the author and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government. Distribution Statement "A" (Approved for Public Release, Distribution Unlimited).

while reducing memory footprint may be accomplished by zeroing out individual weights. Similarly, even for the *same optimization objective*, say reducing inference latency, one may employ filter pruning for a CNN, while pruning 2-D blocks of non-zero weights (Gray et al., 2017) for a language modeling network such as Transformer (Vaswani et al., 2017), since the latter has no convolutional layers. Thus, it is crucial to enable convenient expression of alternative compression schemes, yet none of today’s model compression approaches help the designer tailor compression schemes to their needs.

目前的方法需要手动明确压缩的超参数

Current approaches to model compression also require manual specification of compression hyperparameters, such as **target sparsity**—*the proportion of zero-valued parameters in the compressed model vs. the original*. However, with current approaches, finding the best sparsity often becomes a trial-and-error search, with each such trial having a huge cost (often multiple days for large models such as BERT) and involving training the compressed model to convergence, only to find (in most cases) that the compression objectives are not met. The main difficulty faced by such unguided approaches is that sparsities vary unpredictably with changes in the compression context, making it very difficult to provide users with any guidelines, whatsoever. Therefore, automatic and *sample-efficient* approaches that minimize the number of trials are crucial to support the needs of designers who must adapt a variety of neural networks to a broad spectrum of platforms targeting a wide range of tasks.

To address the above-mentioned problems of flexible ex-

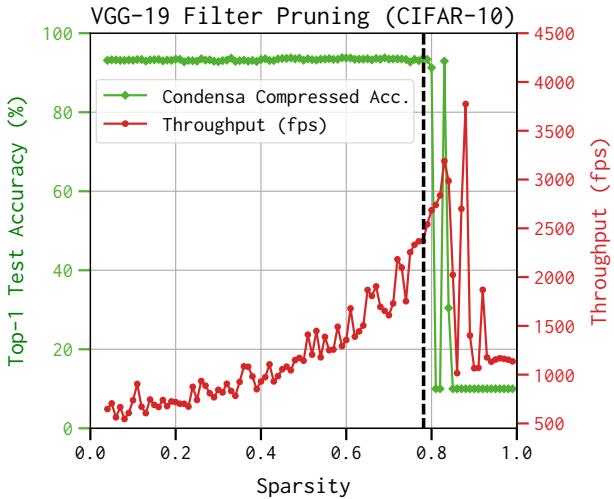


Figure 1. Top-1 test accuracy (green) and throughput (red) vs. sparsity for VGG-19 on CIFAR-10. CONDENA is designed to solve constrained optimization problems of the form “maximize throughput, with a lower bound on accuracy”. In this case, CONDENA automatically discovers a sparsity (vertical dashed line) and compresses the model to this sparsity, improving throughput by $2.59\times$ and accuracy by 0.36%.

pression of compression strategies, automated compression hyperparameter inference, and sample efficiency, we introduce CONDENA, a new framework for programmable model compression. As an illustration of the level of automation provided by CONDENA, consider the problem of improving the inference throughput of VGG-19 (Simonyan & Zisserman, 2014) on the CIFAR-10 image classification task (Krizhevsky et al., 2014). Since VGG-19 is a convolutional neural network, one way to improve its inference performance on modern hardware such as GPUs is by pruning away individual convolutional filters (He et al., 2018a).

由于 VGG-19 是一个卷积神经网络，提高其在现代硬件（如 GPU 上）上的推理性能的一种方法是修剪掉单个卷积过滤器。Figure 1 shows the accuracy and throughput obtained by CONDENA on this task. Here, we plot the compressed model’s top-1 test accuracy and throughput as a function of the sparsity (green and red lines, respectively).¹ CONDENA’s solution corresponds to a sparsity of 0.79 and is depicted as the vertical dashed line. This result is significant for two reasons: (1) using the CONDENA library, the filter pruning strategy employed for this experiment was expressed in less than 10 lines of Python code, and (2) the optimal sparsity of 0.79 that achieves throughput and top-1 accuracy improvements of $2.59\times$ and 0.36%, respectively, was obtained automatically by CONDENA using a sample-efficient constrained Bayesian optimization algorithm. Here, the user didn’t have to specify any sparsities manually, and instead only had to define a domain-specific objective func-

¹Note that these curves are not known a priori and are often extremely expensive to sample; they are only plotted here to better place the obtained solution in context.

tion to maximize (inference throughput, in this case).

This paper makes the following contributions:

1. It presents CONDENA, a new framework for programmable neural network compression. CONDENA supports the expression of the overall compression strategy in Python using operators provided by its compression library. Since each strategy is a Python function, users are able to programmatically compose elementary schemes to build much more complex and practically interesting schemes.
2. It presents a novel sample-efficient algorithm based on Bayesian optimization (B.O.) in CONDENA for automatically inferring optimal sparsities based on a user-provided objective function. Given CONDENA’s ability to support the expression of meaningful high-level objective functions—for example, the throughput (images/sec) of a convolutional neural network—users are freed from the burden of having to specify compression hyperparameters manually.
3. It demonstrates the effectiveness of CONDENA on three image classification and language modeling tasks, resulting in memory footprint reductions of up to $188\times$ and runtime throughput improvements of up to $2.59\times$ using at most 10 samples per search.

2 BACKGROUND

For a given task such as image classification, assume we have trained a large *reference* model $\bar{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} L(\mathbf{w})$, where $L()$ denotes a *loss function* (e.g., cross-entropy on a given training set), and $\mathbf{w} \in \mathbb{R}^P$. *Model compression* refers to finding a smaller model Θ that can be applied to the same task and ideally achieves the same accuracy as $\bar{\mathbf{w}}$. Model compression can be performed in various ways, and CONDENA currently supports two commonly used techniques: pruning and quantization. In pruning, non-zero values from $\bar{\mathbf{w}}$ are eliminated or “pruned” to obtain Θ . Pruning is usually performed using some kind of thresholding (for eg., magnitude-based) and can be unstructured (prune any non-zero value) or structured (prune only *blocks* of non-zeros). On the other hand, quantization retains the number of parameters in Θ but assigns parameters in $\bar{\mathbf{w}}$ one of K codebook values, where the codebook may be fixed or adaptive. CONDENA supports low-precision approximation, which refers to assigning each parameter in $\bar{\mathbf{w}}$ a corresponding lower-precision representation (for example, converting from 32-bit to 16-bit floating-point) and is equivalent to quantization using a fixed codebook.

DNN Compression Techniques There is considerable prior work on accelerating neural networks using structured weight pruning (Wang et al., 2019; McCarley et al., 2020;

Frankle & Carbin, 2018; Han et al., 2015; Luo et al., 2017; Han et al., 2017; Dong et al., 2017; Han et al., 2016; Polyak & Wolf, 2015; Hu et al., 2016; Anwar & Sung, 2016; Molchanov et al., 2016), quantization (Zhu et al., 2016; Gong et al., 2014) and low-rank tensor factorization (Kossaifi et al., 2020; Lebedev et al., 2014; Xue et al., 2013; Denton et al., 2014; Girshick, 2015). Most of these individual compression schemes for pruning and quantization and their combinations can be expressed in CONDENSA. Two common problems with these existing methods are: (1) determining optimal sparsity at a global (network) level, and (2) distributing global sparsity into per-layer sparsities. We tackle these problems efficiently and systematically using our Bayesian and L-C optimizers, respectively, as described in Section 3.

Automated Model Compression Automating model compression involves finding both an optimal compression strategy for a given \bar{w} , along with its corresponding compression hyperparameters such as target sparsity with minimal manual intervention. Current state-of-the-art frameworks in this domain include AMC (He et al., 2018b) and AutoCompress (Liu et al., 2019), which use reinforcement learning and simulated annealing, respectively, to automatically find desirable target sparsities for a fixed pruning strategy. CONDENSA, in contrast, supports the programmable expression of a wide variety of compression strategies (not just pruning). Also, in the context of automated model compression, each sample corresponds to training the compressed model to convergence, and can be extremely expensive to compute; unfortunately, techniques such as reinforcement learning, which is used in AMC (He et al., 2018b), can be highly sample-inefficient (Mnih et al., 2013). To minimize the number of samples drawn, CONDENSA uses a novel and sample-efficient Bayesian optimization-based algorithm for automatically arriving at desirable target sparsities. While Bayesian optimization has previously been demonstrated to work well for general hyperparameter optimization in machine learning and neural architecture search (Snoek et al., 2012; Dai et al., 2019), to the best of our knowledge, we are the first to use sample-efficient search via Bayesian optimization for obtaining compression hyperparameters.

General Compression Algorithms and Tools General accuracy recovery algorithms capable of handling a wide variety of compression techniques provide the foundation for systems like CONDENSA. Apart from the L-C algorithm (Carreira-Perpinán, 2017) which CONDENSA uses, other recent accuracy recovery algorithms have been proposed. ADAM-ADMM (Zhang et al., 2018) proposes a unified framework for structured weight pruning based on ADMM that performs dynamic regularization in which the regularization target is updated in each iteration. DCP (Zhuang et al., 2018) introduces additional losses into the network to increase the discriminative power of inter-

mediate layers and select the most discriminative channels for each layer by considering the additional loss and the reconstruction error. CONDENSA can readily support such algorithms as additional optimizers as described in Section 3. Neural network distiller (Zmora et al., 2018), TensorFlow model optimization toolkit (Google, 2019) and NNCF (Kozlov et al., 2020) are three recent open-source model compression frameworks that support multiple compression schemes. While these projects share a number of common goals with CONDENSA, they differ in two important ways: first, they do not support the expression of schemes as imperative programs containing control-flow, iteration, recursion, etc. (Distiller requires a declarative compression specification in YAML, while the TensorFlow model optimization toolkit operates by modifying the DNN computation graph directly); second, these frameworks do not support automatic compression hyperparameter optimization for black-box objective functions.

3 CONDENSA FRAMEWORK

Figure 2 provides a high-level overview of the CONDENSA framework. As shown on the left-hand side of the figure, a user compresses a pre-trained model \bar{w} by specifying a compression scheme and an objective function f . Both the scheme and objective are specified in Python using operators from the CONDENSA library; alternatively, users may choose from a selection of commonly used built-in schemes and objectives. The CONDENSA library is described in more detail in Section 3.1. Apart from the operator library, the core framework, shown in the middle of the figure, consists primarily of two components: (1) the constrained Bayesian optimizer for inferring optimal sparsities, and (2) the L-C optimizer for accuracy recovery. These components interact with each other as follows: at each iteration, the Bayesian optimizer samples a sparsity s , which is fed into the L-C optimizer. The L-C optimizer distributes this global sparsity across all the layers of the network and performs accuracy recovery (this process is described in more detail in Section 3.3), passing the final obtained accuracy $A(s)$ back to the Bayesian optimizer. The compressed model w obtained by the L-C optimizer is also used to evaluate the user-provided objective function f , the result of which is fed into the Bayesian optimizer. Based on these inputs ($A(s)$ and $f(w)$), the Bayesian optimizer decides the next point to sample. The sparsity that satisfies both the accuracy and objective constraints (s^*) is used to obtain the final compressed model (denoted by Θ in the figure). The Bayesian and L-C optimizers are described in more detail in Sections 3.2 and 3.3, respectively.

据我们所知，我们是第一个通过贝叶斯优化使用样本高效搜索来获得压缩超参数的人。

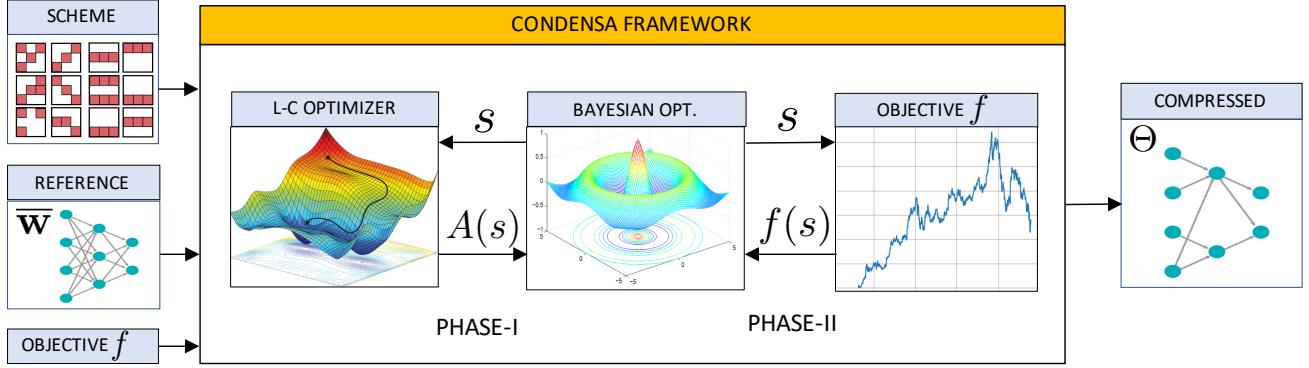


Figure 2. CONDENSA framework overview. The user provides the pre-trained model ($\bar{\mathbf{w}}$), a compression scheme, and an objective function f . CONDENSA uses the Bayesian and L-C optimizers to infer an optimal target sparsity s^* and corresponding compressed model Θ .

3.1 Condensa Library

The CONDENSA Library provides a set of operators for constructing complex compression schemes programmatically in Python. Three sets of operators are currently supported: (1) the `quantize` and `dequantize` operators for converting network parameters from a 32-bit floating-point representation to a lower-precision one such as 16-bit floating-point, and in the opposite direction, respectively; (2) the `prune` operator for unstructured magnitude-based pruning, and (3) the `filter_prune`, `neuron_prune`, and `blockprune` operators for pruning blocks of nonzeros (structure pruning). Each operator can be applied on a per-layer basis.

CONDENSA’s tight integration with the Python ecosystem makes the expression of common compression patterns more natural. For example, operators can be combined with conditional statements to selectively compress layers based on properties of the input DNN and/or target hardware platform, as shown below:

```
# Prune only non-projection layers in ResNets
if not layer.is_projection: prune(layer)
# Quantize only if FP16 hardware is available
if platform_has_fast_fp16(): quantize(layer)
```

Similarly, the use of iteration statements obviates the need for applying compression operators individually for each layer, resulting in more concise and readable schemes. This is in contrast to frameworks such as Distiller (Zmora et al., 2018) which require a per-layer declarative compression specification.

Pre-built Schemes In addition to the layer-wise operators described above, the CONDENSA Library also includes a set of pre-built compression *schemes* that operate on the full model. CONDENSA includes schemes for unstructured and structured pruning, quantization, and composition of individual schemes. These schemes handle a number of low-level details such as magnitude threshold computation from

a sparsity, filter/neuron/block aggregation, etc., enabling non-expert users to quickly get started with CONDENSA without knowledge of low-level implementation details. The current set of pre-built schemes is listed in Table 1.

Listing 1 provides a concrete example of invoking CONDENSA to compress a model. Here, we first train the reference models (lines 2-3) and instantiate the pre-built `FilterPrune` scheme for structured pruning (line 6). We also define our objective function to be throughput (line 8) and specify that it must be maximized (line 10); note that while users may define their own objective functions, CONDENSA also comes bundled with some common objective functions such as model memory footprint and throughput. Next, we instantiate the L-C optimizer (line 12) and the model compressor (Compressor class in Listing) automatically samples and evaluates global sparsities as described in Section 3.2 and returns the final compressed model.

3.2 Sample-Efficient Bayesian Optimization

It is intuitive to split the problem of finding optimal target sparsities into two stages: (1) find the highest target sparsity that loses at most ϵ accuracy w.r.t the original uncompressed model $\bar{\mathbf{w}}$, and (2) in a constrained sparsity regime obtained from stage (1), optimize a user-provided objective function f (e.g., throughput, or memory footprint) and return the solution as the final sparsity. For both stages, CONDENSA utilizes Bayesian optimization as shown in Figure 2.

Bayesian Optimization (B.O.) is an optimization framework based on continually updating a *probabilistic model* with measurements of a function to be optimized (Jones et al., 1998). Given a set of parameters to be optimized, B.O. makes black-box calls to the objective, updates the probabilistic model with the new information, and selects the next point to evaluate using an *acquisition function* that combines information about the expectation and uncertainty of

```

1 # Construct pre-trained model
2 criterion = torch.nn.CrossEntropyLoss()
3 train(model, num_epochs, trainloader, criterion)
4
5 # Instantiate compression scheme
6 prune = condensa.schemes.FilterPrune()
7 # Define objective function
8 tput = condensa.objectives.throughput
9 # Specify optimization operator
10 obj = condensa.searchops.Maximize(tput)
11 # Instantiate L-C optimizer
12 lc = condensa.optimizers.LC(steps=30, lr=0.01)
13 # Build model compressor instance
14 compressor = condensa.Compressor(
15     model=model, # Trained model
16     objective=obj, # Objective
17     eps=0.02, # Accuracy threshold
18     optimizer=lc, # Accuracy recovery
19     scheme=prune, # Compression scheme
20     trainloader=trainloader, # Train dataloader
21     testloader=testloader, # Test dataloader
22     valloader=valloader, # Val dataloader
23     criterion=criterion # Loss criterion
24 )
25 # Obtain compressed model
26 wc = compressor.run()

```

Listing 1. Example usage of the CONDENSA library.

a function value under the probabilistic model. CONDENSA employs a Gaussian Process (G.P.) model for B.O. due to its favorable statistical and computational characteristics (Srinivas et al., 2009). It is worth highlighting that B.O. leverages principled Bayesian inference to trade off exploration and exploitation, and is sample-efficient for non-convex black-box functions such as the ones optimized by CONDENSA (Jones et al., 1998).

In CONDENSA’s two-stage optimization pipeline, we first find a sparsity s_{acc} that constrains the model accuracy function A to the provided ϵ . We then constrain the *sparsity search space* to $(0, s_{acc})$ while optimizing the user-provided objective function f . Note that we assume that A decreases monotonically w.r.t. sparsity in the region $(0, s_{acc})$. For each stage, CONDENSA uses a distinct acquisition function to guide the next best point for function evaluation.

Stage 1: Solving Accuracy Constraints Recall that in the first stage of the sparsity inference process, we aim to find the highest sparsity s_{acc} that loses at most ϵ accuracy w.r.t. the original reference model \bar{w} . To this end, we first define a *Level-Set* L that represents $Acc(\bar{w}) - \epsilon$ and aim to find the point on the accuracy curve of the compressed model that intersects with L ; the sparsity corresponding to this solution will be s_{acc} . We propose a novel acquisition function to find s_{acc} named Domain-Restricted Upper Confidence Bound (DR-UCB).

Scheme	Description
Quantize(dtype)	Quantizes network weights to given datatype <code>dtype</code> .
Prune()	Performs unstructured pruning of network weights.
NeuronPrune(criteria)	Aggregates and prunes neurons (1D blocks) according to <code>criteria</code> .
FilterPrune(criteria)	Aggregates and prunes filters (3D blocks) according to <code>criteria</code> .
StructurePrune(criteria)	Combines neuron and filter pruning.
BlockPrune(criteria, bs)	Aggregates and prunes n-D blocks of size <code>bs</code> according to <code>criteria</code> .
Compose(slist)	Composes together all schemes in <code>slist</code> .

Table 1. List of pre-built compression schemes in CONDENSA.

DR-UCB builds upon an existing level-set black-box optimization technique named ILS-UCB (Garg et al., 2016), which is characterized by two properties: (1) it prioritizes searching in the region where the level set intersects the accuracy curve, (2) it does not seek to precisely learn the shape of the entire accuracy curve. However, in CONDENSA, since accuracy values can be safely assumed to decrease monotonically with increasing sparsity, we notice that it is also possible to progressively restrict the search domain of sparsities based on whether the currently sampled point meets the level-set constraints. In DR-UCB, we exploit this property to greatly improve sample efficiency over ILS-UCB. Mathematically, we define s_t , the sparsity value sampled at iteration t using DR-UCB, as follows:

$$\begin{aligned} s_t = \operatorname{argmax}_s & (1 - \gamma)\sigma(s) - \gamma|\mu(s) - L| \\ \text{s.t. } & s_t > s_i \quad \forall i \in [0, t-1], \quad \mathcal{B}_f(s_t) \geq L \end{aligned} \quad (1)$$

Here, \mathcal{B}_f represents the L-C accuracy function, and s_t is (1) greater than all the previous sparsities s_i , and (2) satisfies the level set constraint $\mathcal{B}_f(s_t) \geq L$. We achieve this by minimizing the difference between the GP’s mean curve $\mu(s)$ and the level set using the term $|\mu(s) - L|$ in (1); the parameter γ controls the trade-off between exploitation and exploration. Algorithm 1 illustrates how DR-UCB is employed to efficiently find s_{acc} .

Stage 2: Optimizing the User-Defined Objective Once we find a sparsity s_{acc} that satisfies the user-provided accuracy constraints in stage 1, our next objective is to find the

Algorithm 1 Bayesian Sparsity Inference with Domain Restriction

```

1: procedure BODR-UCB( $\mathcal{B}_f$ ,  $L$ ,  $T$ )
   ▷  $\mathcal{B}_f$ : Function to optimize
   ▷  $L$ : Level set
   ▷  $T$ : # Iterations
2:   GP  $\leftarrow$  GP-Regressor.initialize()
3:    $s_0 \leftarrow 0$ ;  $D \leftarrow (0, 1)$ ;  $\mathbf{X} \leftarrow \emptyset$ 
4:   for  $t \leftarrow 1, 2, \dots, T - 1$  do
5:      $s_t \leftarrow \text{argmax}_D \text{DR-UCB}(s | \mathbf{X}_{0:t-1})$ 
6:      $y_t \leftarrow \mathcal{B}_f(s_t)$ 
7:     if  $s_t > s_{t-1}$  and  $y_t \geq L$  then
8:        $D \leftarrow (s_t, 1)$ 
9:     end if
10:     $\mathbf{X}_{0:t} \leftarrow \{\mathbf{X}_{0:t-1}, (s_t, y_t)\}$ 
11:    GP.Update( $\mathbf{X}_{0:t}$ )
12:   return  $s_{T-1}$ 

```

final sparsity s^* that optimizes the user-defined objective function f in the constrained sparsity domain $(0, s_{acc})$. For this, we employ the Upper and Lower Confidence Bound (UCB/LCB) acquisition functions for function maximization and minimization, respectively (Srinivas et al., 2009).

3.3 Accuracy Recovery using L-C

As described earlier in this section, given a reference model, compression scheme, and compression hyperparameter values (obtained automatically by the Bayesian hyperparameter optimization subsystem described in Section 3.2), CONDENSA tries to recover any accuracy lost due to compression. While the compressed model, denoted as Θ , can be obtained by directly zeroing out lower-magnitude parameters from the reference model $\bar{\mathbf{w}}$ (a technique referred to as *direct compression*), the resulting model Θ is generally sub-optimal w.r.t. the loss since the latter is ignored in learning Θ . Instead, we desire an *accuracy recovery algorithm* that obtains an *optimally compressed model* with locally optimal loss. An effective accuracy recovery mechanism for CONDENSA must ideally have three important attributes: (1) able to handle all the compression operators supported by CONDENSA, (2) be efficient with relatively low overheads, and (3) provide optimality guarantees whenever possible. In this paper, we use the recently proposed L-C algorithm (Carreira-Perpinán & Idelbayev, 2018), since it satisfies all three of the above requirements. In L-C, model compression is formulated as a constrained optimization problem:

$$\min_{\mathbf{w}, \Theta} L(\mathbf{w}) \quad \text{s.t.} \quad \mathbf{w} = \mathcal{D}(\Theta) \quad (2)$$

Here, the *decompression mapping* $\mathcal{D} : \Theta \in \mathbb{R}^Q \rightarrow \mathbf{w} \in \mathbb{R}^P$ maps a low-dimensional parameterization to uncompressed model weights, and the *compression mapping* $\mathcal{C}(\mathbf{w}) = \text{argmin}_{\Theta} \|\mathbf{w} - \mathcal{D}(\Theta)\|^2$ behaves similar to the inverse of \mathcal{D} .

This formulation naturally supports a number of well-known compression techniques. In particular, pruning is defined as $\mathbf{w} = \mathcal{D}(\Theta) = \Theta$ where \mathbf{w} is real and Θ is constrained to have fewer nonzero values by removing (zeroing out) lower magnitude weights; low-precision approximation defines a constraint $w_i = \theta_i$ per parameter where w_i is in a higher-precision representation and θ_i is in a lower-precision one.

Eq. 2 is non-convex due to two reasons: (1) the original problem of training the reference model is already non-convex for models such as DNNs, making the objective function of Eq 2 non-convex, and (2) the decompression mapping $\mathcal{D}(\Theta)$ typically adds another layer of non-convexity caused by an underlying combinatorial problem. While a number of non-convex algorithms may be used to solve Eq 2, we focus on the augmented Lagrangian (AL) method (Wright & Nocedal, 1999) implemented in the L-C algorithm (Carreira-Perpinán & Idelbayev, 2018) in this paper, since it is relatively efficient and easy to implement. As its name indicates, the L-C algorithm alternates between two steps: a learning (L) step which trains the uncompressed model but with a quadratic regularization term, and a compression (C) step, which finds the best compression of \mathbf{w} (the current uncompressed model) and corresponds to the definition of compression mapping \mathcal{C} . We refer the reader to (Carreira-Perpinán & Idelbayev, 2018) for a more detailed description of the L-C algorithm. Other recent AL-based algorithms that could potentially be used include ADMM (Zhang et al., 2018) and DCP (Zhuang et al., 2018).

3.4 Implementation

The CONDENSA library and L-C optimizer are implemented in Python and are designed to interoperate seamlessly with the PyTorch framework (Paszke et al., 2017). While we chose PyTorch for its widespread use in the machine learning community, it is worth noting that CONDENSA's design is general and that its features can be implemented in other similar frameworks such as TensorFlow (Abadi et al., 2016) and MXNET (Chen et al., 2015). We currently use a publicly available Python library for Bayesian global optimization with Gaussian Processes (fmfn, 2019).

Network Thinning Condensa comes pre-built with three *structure pruning* schemes: filter, neuron, and block pruning, as shown in Table 1. The application of these schemes may yield *zero structures*, which refer to blocks of zeros within a DNN's parameters. *Network thinning* refers to the process of identifying and removing such zero structures and consequently reducing the number of floating-point operations executed by the target hardware platform. Condensa employs a three-phase network thinning algorithm for structure pruning: in the first phase, we construct an in-memory graph representation of the target DNN. PyTorch makes this non-trivial, as its eager execution semantics preclude it from

ever building a full graph-based representation of the DNN. To overcome this, we trace a forward execution path of the DNN and use it to construct an in-memory representation based on the ONNX format. In the next phase, we create a *thinning strategy* by analyzing the dependencies between the nodes of the graph constructed in the first phase. This step primarily involves keeping track of tensor dimension changes in a node due to thinning and ensuring that the corresponding tensor dimensions of the node’s successors are appropriately adjusted. Due to the possibility of complex dependence patterns such as skip nodes in real-world DNNs (for example, deep residual networks (He et al., 2016)), this step is the most challenging to implement. In the final phase, we apply the thinning strategy obtained in phase 2 and physically alter tensor shapes to obtain the final thinned network. The Condensa Library provides the `thin` method which can be used to thin a given compressed model.

4 EVALUATION

We conduct extensive experiments and fully analyze CONDENSA on three real-world tasks:

(1) Image Classification on CIFAR-10 The CIFAR-10 dataset (Krizhevsky et al., 2014) consists of $50k$ training and $10k$ testing 32×32 images in 10 classes. We train the VGG-19 (Simonyan & Zisserman, 2014) and ResNet56 (He et al., 2016) neural networks on this dataset for 160 epochs with batch normalization, weight decay (10^{-4}), decreasing learning rate schedules (starting from 0.1) and augmented training data.

(2) Image Classification on ImageNet Here, we use the VGG-16 neural network (Simonyan & Zisserman, 2014) trained on the challenging ImageNet task (Deng et al., 2009), specifically the ILSVRC2012 version. We use PyTorch (Paszke et al., 2017) and default pretrained models as a starting point.

(3) Language Modeling on WikiText-2 We trained a 2-layer LSTM model to perform a language modeling task on the WikiText-2 dataset (Merity et al., 2016). We used a hidden state size of 650 and included a dropout layer between the two RNN layers with a dropout probability of 0.5. The LSTM received word embeddings of size 650. For training, we used truncated Backpropagation Through Time (truncated BPTT) with a sequence length of 50. The training batch size was set to 30, and models were optimized using SGD with a learning rate of 20. This setup is similar to the one used by Yu et al. (Yu et al., 2019).

We optimize the networks in each task for two distinct objectives described below:

Objective 1: Minimize Memory Footprint The memory footprint of a model is defined as the number of bytes con-

sumed by the model’s *non-zero* parameters. Reducing the footprint below a threshold value is desirable, especially for memory-constrained devices such as mobile phones, and can be accomplished through either pruning or quantization, or both. For reducing footprint, we define a compression scheme that performs unstructured pruning of each learnable layer (except batch normalization layers), and then quantizes it to half-precision floating-point, yielding an additional 2x reduction. We denote this scheme by P+Q and implement it using the CONDENSA library as follows (see Table 1 for the full list of schemes):

```
from schemes import Compose, Prune, Quantize
scheme = Compose([Prune(), Quantize(float16)])
```

Objective 2: Maximize Throughput Inference throughput is defined as the number of input samples processed by a model per second, and is commonly used for measuring real-world performance. For CIFAR-10 and ImageNet, we measure hardware inference throughput of the compressed model in the objective function. We use an NVIDIA Titan V GPU with the TensorRT 5 framework to obtain throughput data. For WikiText-2, due to the lack of optimized block-sparse kernels for PyTorch, we measure the floating-point operations (FLOPs) of the compressed model instead as a proxy for inference performance. To improve throughput, we focus on removing entire blocks of non-zeros, such as convolutional filters, since they have been proven to improve performance on real-world hardware (He et al., 2018a; Gray et al., 2017). For CIFAR-10 and ImageNet, we use filter pruning, since all the networks we consider are CNNs. In WikiText-2, we employ block pruning with a block size of 5.

Bayesian Optimizer Settings We use a Gaussian Processes prior with the Matern kernel ($\nu = 2.5$), length scale of 1.0 and α value of 0.1 with normalization of the predictions. For the GP regressor, the noise level in the covariance matrix is governed by another parameter, which we set to a very low value of 10^{-6} . For the DR-UCB acquisition function, we use a γ value of 0.95 for all our experiments with a bias towards sampling more in the area of level set, with the intention that the Bayesian optimizer results in a favorable sparsity level in as few samples as possible. We implemented DR-UCB using the `fmpn/B0` package (fmpn, 2019).

L-C Optimizer Settings The L-C optimizer was configured as follows: for all experiments, we use $\mu_j = \mu_0 a^j$, with $\mu_0 = 10^{-3}$ and $a = 1.1$ where j is the L-C iteration. For CIFAR-10 and ImageNet, we use the SGD optimizer in the learning (L) step with a momentum value of 0.9, with the learning rate decayed from 0.1 to 10^{-5} over each mini-batch iteration. We use the Adam optimizer in the L-step of WikiText-2 with a fixed learning rate of 10^{-4} . We ran between 4000-5000 mini-batch iterations in each L-step,

Table 2. CONDENSA performance results on CIFAR-10, ImageNet, and WikiText-2. Here, s^* represents the target sparsity obtained by CONDENSA, r_c is the memory footprint reduction, and s_F the FLOP reduction. The level-set, represented by ϵ , is set to 2% below baseline in all experiments.

METHOD	DATASET	NETWORK	s^*	ACCURACY	r_c	THROUGHPUT
BASELINE	CIFAR-10	VGG19-BN		92.98%	1×	1×
CONDENSA P+Q	CIFAR-10	VGG19-BN	0.99	93.26%	188.23×	N/A
CONDENSA FILTER	CIFAR-10	VGG19-BN	0.79	93.34%	1.35×	2.59×
BASELINE	CIFAR-10	RESNET56		92.75%	1×	1×
AMC (HE ET AL., 2018B)	CIFAR-10	RESNET56	N/A	90.1%	N/A	$s_F = 2\times$
CONDENSA P+Q	CIFAR-10	RESNET56	0.95	91.42%	31.14×	N/A
CONDENSA FILTER	CIFAR-10	RESNET56	0.63	93.18%	1.14×	1.17×
BASELINE	IMAGENET	VGG16-BN		91.50%	1×	1×
FILTER PRUNING (HE ET AL., 2017)	IMAGENET	VGG16-BN		89.80%	$\approx 4\times$	N/A
AUTOCOMPRESS (LIU ET AL., 2019)	IMAGENET	VGG16-BN	N/A	90.90%	6.4×	N/A
AMC (HE ET AL., 2018B)	IMAGENET	VGG16-BN	N/A	90.1%	N/A	$s_F = 1.25\times$
CONDENSA P+Q	IMAGENET	VGG16-BN	0.93	89.89%	29.29	N/A
CONDENSA FILTER	IMAGENET	VGG16-BN	0.12	90.25%	1×	1.16×
BASELINE	WIKITEXT-2	LSTM		LOG-PERPLEXITY: 4.70	1×	1×
LOTTERY TICKET(YU ET AL., 2019)	WIKITEXT-2	LSTM	N/A	LOG-PERPLEXITY: 4.70	$\approx 10\times$	N/A
CONDENSA P+Q	WIKITEXT-2	LSTM	0.92	LOG-PERPLEXITY: 4.75	4.2×	N/A
CONDENSA BLOCK	WIKITEXT-2	LSTM	0.60	LOG-PERPLEXITY: 4.62	1.1×	$s_F = 2.14$

with a higher number of iterations in the first L-step (30k for CIFAR-10 and ImageNet, and 7k for WikiText-2) as recommended by (Carreira-Perpinán & Idelbayev, 2018). We ran 5, 30, and 50 L-C iterations for WikiText-2, ImageNet, and CIFAR-10, respectively; compared to CIFAR-10, we ran relatively fewer iterations for ImageNet due to its significantly higher computational cost, and ran an extra 5 fine-tuning iterations instead. We use the same mini-batch sizes as during training for all experiments, and use validation datasets to select the best model during compression (we perform a 9:1 training:validation split for CIFAR-10 since it doesn’t include a validation dataset).

4.1 Results

We present the memory footprint reductions and inference throughput improvements obtained by CONDENSA for each of the three tasks we evaluate in Table 2. For each task, we list the sparsity obtained by the CONDENSA Bayesian optimizer (s^* in the table), its corresponding accuracy/perplexity (top-1 accuracy, top-5 accuracy, and log perplexity for CIFAR-10, ImageNet, and WikiText-2, respectively), memory footprint reductions using pruning and quantization (column labeled r_c), and inference throughput/FLOP improvements using filter/block pruning. We also compare our approach with recent work on automated model compression. For CIFAR-10 and ImageNet, we compare our results with AMC (He et al., 2018b) and AutoSlim (Liu et al., 2019), and for WikiText-2, we compare with (Yu et al., 2019). Since AMC (He et al., 2018b) and (Yu et al., 2019) do not report actual runtime numbers on hardware, we

report the corresponding FLOP improvements instead (values marked s_F). We also use FLOP reduction as a metric for LSTM block pruning, as described above.

Using the P+Q scheme designed to minimize memory footprint, CONDENSA is able to obtain compression ratios up to 188×, which surpasses those of frameworks such as AutoCompress. While AMC and AutoCompress only report theoretical FLOP improvements on CIFAR-10 and ImageNet, the filter pruning strategy implemented using CONDENSA yields real-world runtime improvements of up to 2.59× on an NVIDIA Titan V GPU. Since AMC and AutoCompress do not report the number of samples evaluated to arrive at their solutions, we are unable to directly compare sample efficiencies with these frameworks; however, we notice that CONDENSA obtains desirable model sparsities using a fixed 10 iterations per search in all experiments. Finally, while we set the level set to be 2% below the accuracy of the reference model in all our experiments, we notice that CONDENSA-compressed models often exceed baseline accuracy.

4.2 Sparsity Profile Analysis

Figures 3 and 4 illustrate how a compressed model’s accuracy, inference performance, and memory footprint vary w.r.t. sparsities for the CIFAR-10 and WikiText-2 tasks. All three of these functions are *assumed to be unknown* in our problem formulation, but we compute them explicitly here to better understand the quality of solutions produced by CONDENSA. For each figure, compression accuracies (shown in green) are obtained by running the L-C algorithm to convergence for 100 sparsity values ranging from 0.9 to

--- Direct Compression Acc. — Condensa Compression Acc. — Objective Fn.

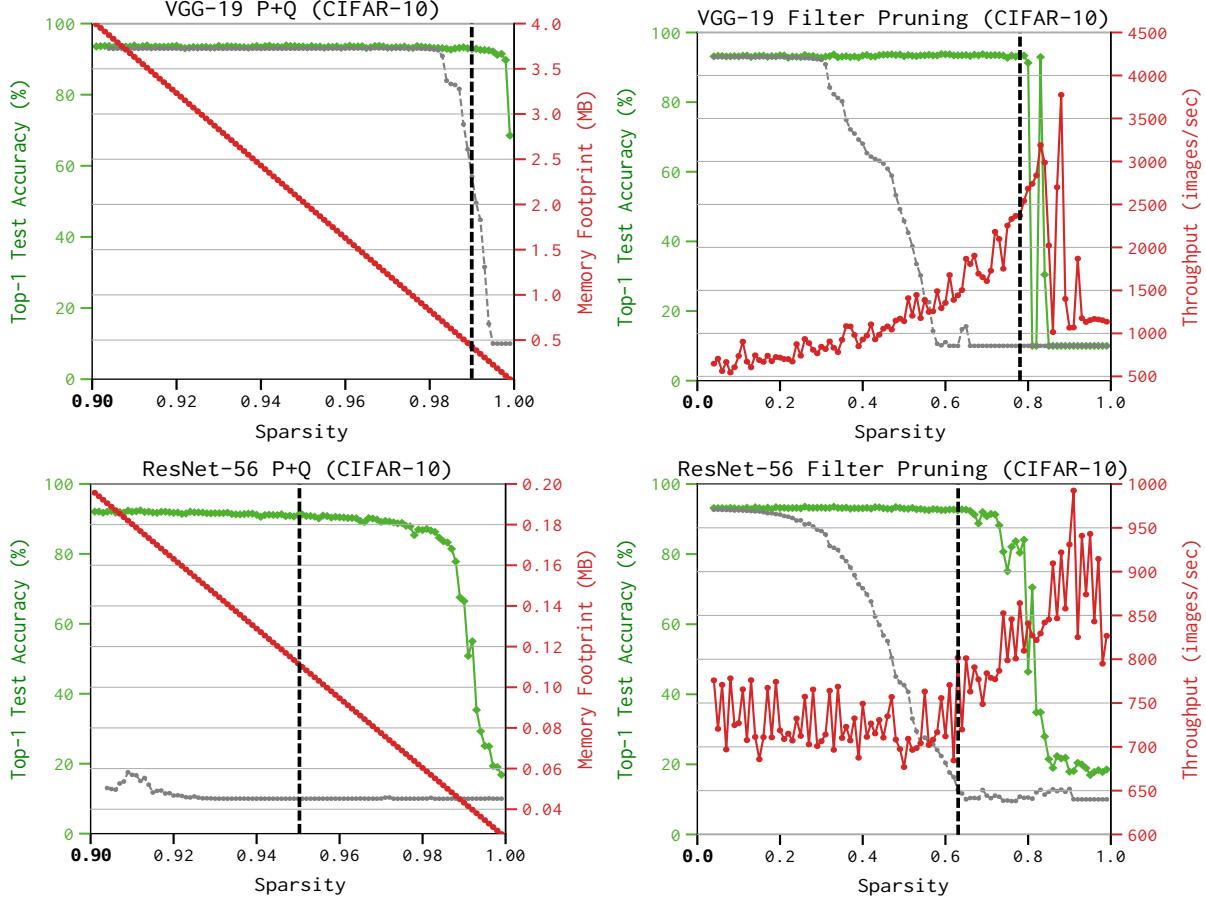


Figure 3. CONDENSA sparsity profiles for VGG19-BN and ResNet56 for CIFAR-10. Column 1 shows the problem of the form “minimize *memory footprint* with a lower bound on accuracy”, while Column 2 illustrates “maximize *throughput* with a lower bound on accuracy”. The DC line (gray) shows accuracy values if no accuracy recovery with L-C is performed. Note that the x-axis ranges are different: the plots on the left have sparsities ranging from 0.9 to 1.0 while those on the right have values ranging from 0 to 1.

1.0 (for P+Q), and from 0 to 1 for the filter and block pruning schemes; collecting each such point requires between 30 minutes to 8 hours of time on a single NVIDIA Tesla V100 GPU. We are unable to show the full profile for ImageNet due to its significantly higher computation cost: collecting each data point for compression accuracy requires over 12 hours of compute time on a node with 8 Tesla V100 GPUs. Inference throughput, FLOPs, and memory footprint data is collected for each compressed model and depicted by red lines in the figures (right-hand-side y-axis). We also show direct compression (DC) accuracies in gray for comparison (DC is described in more detail in Section 3.3). In each figure, the sparsity found by CONDENSA is shown as a black vertical dashed line.

We notice three important trends in Figures 3 and 4: (1) CONDENSA consistently finds solutions near the ‘knee’ of

the L-C accuracy curves, signifying the effectiveness of the DR-UCB acquisition function; (2) local minima/maxima is avoided while optimizing the objective function, demonstrating that the UCB acquisition function for objective function optimization is working as expected, and (3) the knee of the D-C accuracy curves occur at significantly lower sparsity ratios; the L-C optimizer, on the other hand is able to recover accuracies up to much higher sparsities.

4.3 Layerwise Runtime Performance

In this section, we analyze how improving throughput using compression translates to execution time improvements for each layer on actual hardware. For this experiment, we focus on VGG-19 on CIFAR-10, since it has a relatively simple structure and is easy to analyze on a layer-by-layer basis. We use filter pruning with a target sparsity of 0.79

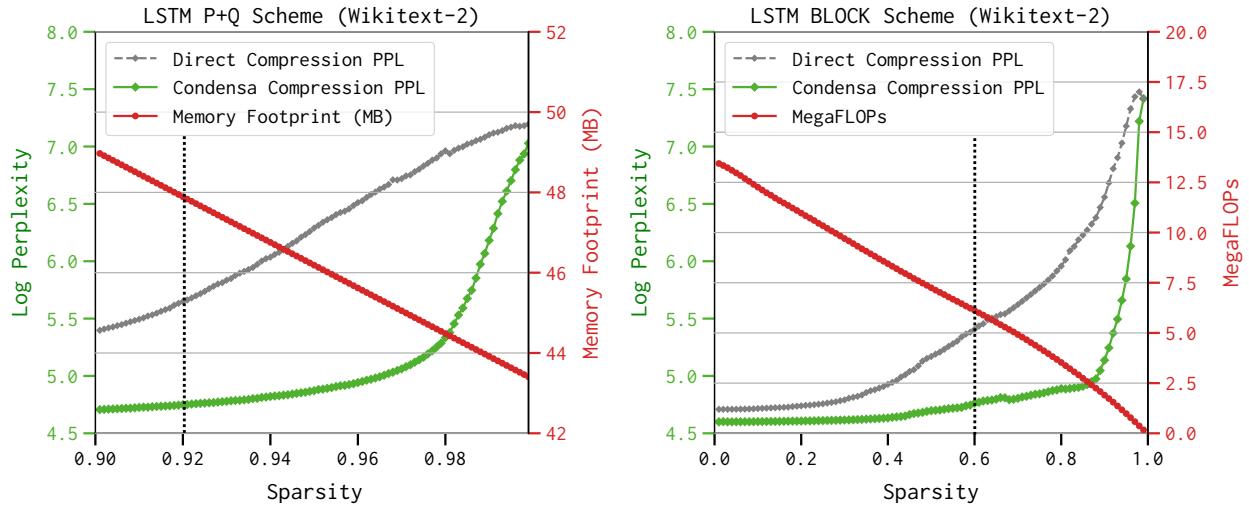


Figure 4. WikiText-2 two-layer LSTM results for pruning + quantization (left) and block pruning with block size of 5 (right). Note that the x-axis ranges are different: the plot on the left has sparsity values ranging from 0.9 to 1.0 while the one on the right has values ranging from 0 to 1.

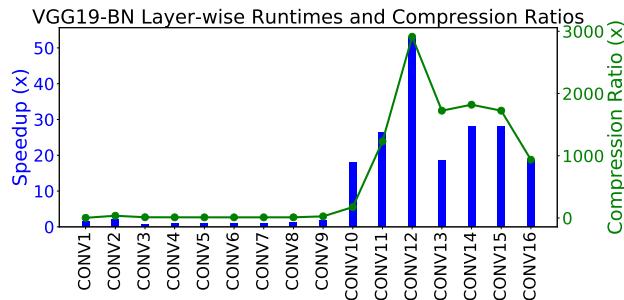


Figure 5. TensorRT runtimes and compression ratios of convolutional layers in VGG19-BN (filter pruning).

(found by the Bayesian optimizer, as shown in Table 2) for this experiment. Figure 5 shows layer-by-layer mean runtimes collected over 100 runs using TensorRT (blue bars, left y-axis), and compression ratios (green line, right y-axis) for filter pruning. We only show data for convolutional layers as they dominate computation time for this network. We make two key observations: (1) runtime speedups on real hardware are largely correlated with compression ratios, but may be affected by hardware and implementation details (e.g., compare conv13 with conv14 in the Figure), and (2) higher compression ratios and corresponding speedups for the later layers of the network, which indicates that distributing a given global sparsity evenly across network layers may not always be optimal, and algorithms such as L-C are essential to automatically finding desirable distributions of sparsity across layers.

5 ACKNOWLEDGMENTS

This material is based upon work supported by DARPA under Contract No. HR0011-18-3-0007, NSF award CCF-1704715 and a CIFAR AI Chair award. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the U.S. Government. Distribution Statement "A" (Approved for Public Release, Distribution Unlimited).

6 CONCLUSIONS

This paper has presented CONDENSA, which is a flexible programming system for model compression and corresponding hyper-parameter optimization. We have demonstrated CONDENSA’s effectiveness and ease-of-use on a range of state-of-the-art DNNs for image classification and language modeling, and achieved memory footprint reductions of up to $188\times$ and runtime throughput improvements of up to $2.59\times$ using at most 10 samples per search. With the initial framework in place, we envision a number of directions to expand on CONDENSA’s capability. For example, we plan to augment automatic sparsity inference with support for additional compression hyperparameters such as block sizes in block-sparsification (Gray et al., 2017), and data types for quantization. Our long-term goal is a framework that makes model compression easier, more flexible, and accessible to a wide range of users.

REFERENCES

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pp. 265–283, 2016.
- Anwar, S. and Sung, W. Compact deep convolutional neural networks with coarse pruning. *arXiv preprint arXiv:1610.09639*, 2016.
- Carreira-Perpinán, M. A. Model compression as constrained optimization, with application to neural nets. part I: General framework. *arXiv preprint arXiv:1707.01209*, 2017.
- Carreira-Perpinán, M. A. and Idelbayev, Y. “learning-compression” algorithms for neural net pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8532–8541, 2018.
- Chen, T., Li, M., Li, Y., Lin, M., Wang, N., Wang, M., Xiao, T., Xu, B., Zhang, C., and Zhang, Z. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- Dai, X., Zhang, P., Wu, B., Yin, H., Sun, F., Wang, Y., Dukhan, M., Hu, Y., Wu, Y., Jia, Y., et al. Chamnet: Towards efficient network design through platform-aware model adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11398–11407, 2019.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Denton, E. L., Zaremba, W., Bruna, J., LeCun, Y., and Fergus, R. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in neural information processing systems*, pp. 1269–1277, 2014.
- Dong, X., Huang, J., Yang, Y., and Yan, S. More is less: A more complicated network with less inference complexity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5840–5848, 2017.
- fmpn. A Python implementation of global optimization with Gaussian processes. <https://github.com/fmpn/BayesianOptimization>, 2019. [Online; accessed 1-September-2019].
- Frankle, J. and Carbin, M. The lottery ticket hypothesis: Training pruned neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Garg, A., Sen, S., Kapadia, R., Jen, Y., McKinley, S., Miller, L., and Goldberg, K. Tumor localization using automated palpation with gaussian process adaptive sampling. In *2016 IEEE International Conference on Automation Science and Engineering (CASE)*, pp. 194–200. IEEE, 2016.
- Girshick, R. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448, 2015.
- Gong, Y., Liu, L., Yang, M., and Bourdev, L. Compressing deep convolutional networks using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014.
- Google. TensorFlow model optimization toolkit. <https://github.com/tensorflow/model-optimization>, 2019. [Online; accessed 1-September-2019].
- Gray, S., Radford, A., and Kingma, D. P. GPU kernels for block-sparse weights. *arXiv preprint arXiv:1711.09224*, 2017.
- Gupta, S., Agrawal, A., Gopalakrishnan, K., and Narayanan, P. Deep learning with limited numerical precision. In *International Conference on Machine Learning*, pp. 1737–1746, 2015.
- Han, S., Pool, J., Tran, J., and Dally, W. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pp. 1135–1143, 2015.
- Han, S., Liu, X., Mao, H., Pu, J., Pedram, A., Horowitz, M. A., and Dally, W. J. Eie: efficient inference engine on compressed deep neural network. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 243–254. IEEE, 2016.
- Han, S., Kang, J., Mao, H., Hu, Y., Li, X., Li, Y., Xie, D., Luo, H., Yao, S., Wang, Y., et al. Ese: Efficient speech recognition engine with sparse LSTM on FPGA. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pp. 75–84. ACM, 2017.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- He, Y., Zhang, X., and Sun, J. Channel pruning for accelerating very deep neural networks. In *International Conference on Computer Vision (ICCV)*, volume 2, 2017.
- He, Y., Dong, X., Kang, G., Fu, Y., and Yang, Y. Progressive deep neural networks acceleration via soft filter pruning. *arXiv preprint arXiv:1808.07471*, 2018a.

- He, Y., Lin, J., Liu, Z., Wang, H., Li, L.-J., and Han, S. Amc: Automl for model compression and acceleration on mobile devices. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 784–800, 2018b.
- Hu, H., Peng, R., Tai, Y.-W., and Tang, C.-K. Network trimming: A data-driven neuron pruning approach towards efficient deep architectures. *arXiv preprint arXiv:1607.03250*, 2016.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269. IEEE, 2017.
- Jaderberg, M., Vedaldi, A., and Zisserman, A. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- Jones, D. R., Schonlau, M., and Welch, W. J. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- Kossaifi, J., Toisoul, A., Bulat, A., Panagakis, Y., Hospedales, T., and Pantic, M. Factorized higher-order cnns with an application to spatio-temporal emotion estimation, 2020.
- Kozlov, A., Lazarevich, I., Shamporov, V., Lyalyushkin, N., and Gorbachev, Y. Neural network compression framework for fast model inference. *arXiv preprint arXiv:2002.08679*, 2020.
- Krizhevsky, A., Nair, V., and Hinton, G. The cifar-10 dataset. *online: http://www.cs.toronto.edu/kriz/cifar.html*, 55, 2014.
- Lebedev, V., Ganin, Y., Rakuba, M., Oseledets, I., and Lempitsky, V. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014.
- Li, H., Kadav, A., Durdanovic, I., Samet, H., and Graf, H. P. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- Liu, N., Ma, X., Xu, Z., Wang, Y., Tang, J., and Ye, J. Autoslim: An automatic dnn structured pruning framework for ultra-high compression rates. *arXiv preprint arXiv:1907.03141*, 2019.
- Luo, J.-H., Wu, J., and Lin, W. Thinet: A filter level pruning method for deep neural network compression. *arXiv preprint arXiv:1707.06342*, 2017.
- McCarley, J. S., Chakravarti, R., and Sil, A. Structured pruning of a bert-based question answering model, 2020.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing Atari with Deep Reinforcement Learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Molchanov, P., Tyree, S., Karras, T., Aila, T., and Kautz, J. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- Polyak, A. and Wolf, L. Channel-level acceleration of deep face representations. *IEEE Access*, 3:2163–2175, 2015.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Snoek, J., Larochelle, H., and Adams, R. P. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pp. 2951–2959, 2012.
- Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Wang, Z., Wohlwend, J., and Lei, T. Structured pruning of large language models, 2019.
- Wright, S. and Nocedal, J. Numerical optimization. *Springer Science*, 35(67-68):7, 1999.
- Xue, J., Li, J., and Gong, Y. Restructuring of deep neural network acoustic models with singular value decomposition. In *Interspeech*, pp. 2365–2369, 2013.
- Yu, H., Edunov, S., Tian, Y., and Morcos, A. S. Playing the lottery with rewards and multiple languages: lottery tickets in rl and nlp. *arXiv preprint arXiv:1906.02768*, 2019.
- Zhang, T., Zhang, K., Ye, S., Li, J., Tang, J., Wen, W., Lin, X., Fardad, M., and Wang, Y. Adam-ADMM: A unified, systematic framework of structured weight pruning for DNNs. *arXiv preprint arXiv:1807.11091*, 2018.

Zhu, C., Han, S., Mao, H., and Dally, W. J. Trained ternary quantization. *arXiv preprint arXiv:1612.01064*, 2016.

Zhuang, Z., Tan, M., Zhuang, B., Liu, J., Guo, Y., Wu, Q., Huang, J., and Zhu, J. Discrimination-aware channel pruning for deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 883–894, 2018.

Zmora, N., Jacob, G., and Novik, G. Neural network distiller, June 2018. URL <https://doi.org/10.5281/zenodo.1297430>.