

---

# DROPBACK: FULL DNN TRAINING ON A PRUNED WEIGHT BUDGET

---

Anonymous Authors<sup>1</sup>

## ABSTRACT

We introduce a technique that reduces deep neural networks both during and after training by constraining the total number of weights updated during backpropagation to those with the highest total gradients. The remaining weights are forgotten and their initial value is regenerated at every access to avoid storing them in memory. This dramatically reduces the number of off-chip memory accesses during both training and inference, a key component of the energy needs of DNN accelerators. By ensuring that the total weight diffusion remains close to that of baseline unpruned SGD, networks pruned using DropBack are able to maintain high accuracy across network architectures — including networks previously identified as difficult to compress, such as Densenet and WRN. Applying DropBack to ResNet18 on ImageNet, we observe an 11.7× weight reduction with no accuracy loss, and up to 24.4× with a small accuracy impact.

## 1 INTRODUCTION

of vendors announcing low-power deep neural network accelerators that reduce computation costs to specifically target mobile and embedded applications (Intel, 2017; Samsung, 2018; Qualcomm, 2017; Kingsley-Hughes, 2017). To address the comparatively smaller memories in mobile devices — an order of magnitude less capacity and two orders of magnitude less bandwidth than a datacentre-class GPU — many researchers have proposed pruning, quantizing, and compressing neural networks, typically making a tradeoff of accuracy versus network size (LeCun et al., 1990; Hassibi et al., 1993; Han et al., 2015b;a; Wu et al., 2015; Choi et al., 2016; Alvarez & Salzmann, 2017; Ge et al., 2017; Zhou et al., 2017; Luo et al., 2017; Yang et al., 2017, etc.). The past year has seen the arrival of on-device inference, with a number of comparably little attention, however, has been paid to the problem of *on-device training*, which has so far been limited to simple models (e.g., Apple, 2017). Training is much more energetically expensive than inference: one iteration on a single sample involves roughly three times as many weight accesses (forward pass, gradient computation, weight update), and typically many iterations on many samples are required. Energy efficiency is limited by off-chip memory accesses, which cost hundreds of times more energy than computation: in a 45nm process, for example, accessing a 32-bit value from DRAM costs over 700× more energy than an 32-bit floating-point compute operation (640pJ vs. 0.9pJ,

Han et al., 2016), with an even larger gap at smaller process nodes. Existing pruning, quantization, and reduction techniques are only applied *after* training, and do not ameliorate this energy bottleneck.

Most of the off-chip memory references during training come from accessing activations and weights. The ratio depends on the amount of weight reuse available in the specific architecture being trained: for example, in an MLP weights are not reused within a single training sample and typically need more bandwidth than activations, while in a CNN each filter is reused for an entire layer and there can be an order of magnitude more activation accesses than weight accesses.

However, activations can be compressed 7× via techniques like gradient checkpointing (Chen et al., 2016) as well as accelerator architectures that elide zero-valued or small activations (Albericio et al., 2016; Judd et al., 2017). In addition, low-end devices with fewer compute resources may not even benefit from weight reuse across batches, as training on small batches and single images is often effective (Masters & Luschi, 2018). In this paper, therefore, we focus on reducing the number of *weights* that must be accessed during — and after — the training process.

We propose DropBack, a training-time pruning technique that relies on three observations:

- (a) that the parameters that have accumulated the highest total gradients account for most of the learning;
- (b) that the values of these parameters are predicated on the initialization values chosen for the remaining parameters; and

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

(c) that initialization values can be recomputed on-the-fly without needing to access memory.

Unlike prior pruning techniques, which train an unconstrained network, prune it, and retrain the pruned network, DropBack prunes the network from immediately at the start of training and does not require a retraining pass. Only the weights that remain require memory, and most of the weight set is never stored. Because DropBack recomputes initialization parameters, it can prune layers like batch normalization or parametric ReLU, which cannot be pruned using existing approaches.

DropBack outperforms best-in-class pruning-only methods (which require retraining) on network architectures that are already dense and have been found particularly challenging to compress (Liu et al., 2017; Louizos et al., 2017; Li et al., 2017). On Densenet, we achieve 5.86% validation error with 4.5 $\times$  weight reduction (vs. best prior 5.65% / 2.9 $\times$  reduction), and on WRN-28-10 we achieve 3.85%–4.20% error with 4.5 $\times$ –7.3 $\times$  weight reduction (vs. best prior 3.75% uncompressed and best pruned prior 16.6% / 4 $\times$  reduction). Using ResNet18 on ImageNet, we observe an 11.7 $\times$  weight reduction without accuracy loss, and up to 24.4 $\times$  with a small accuracy impact. Like other pruning techniques, it is orthogonal to, and can be combined with, quantization and value compression methods.

## 2 DROPBACK: TRAINING PRUNED NETWORKS

### 2.1 Approach and key insights

Existing deep neural network pruning techniques rely on the observation that many parameters do not contribute useful information to the final output, and can be removed. Because the exact values of the surviving weights still depend on the values of the zeroed parameters, the network must still be retrained, but can recover close to original accuracy with an order of magnitude fewer weights (Han et al., 2015a; Zhu et al., 2017; Ge et al., 2017; Masana et al., 2017; Luo et al., 2017; Ullrich et al., 2017; Yang et al., 2017, etc.). While they reduce the energy expended on memory accesses during inference, they require retraining and therefore *increase* the number of energy-consuming memory accesses at training time.

DropBack has a different purpose: we aim to reduce the memory footprints both *during* and *after* training, and so we must reduce the *number of parameters tracked* at training time. To choose which parameters to track, we are guided by three related but slightly different insights, described below.

**Track the highest accumulated gradients.** As post-training pruning techniques delete weights with the lowest *values*, it may seem natural to keep track of gradients for

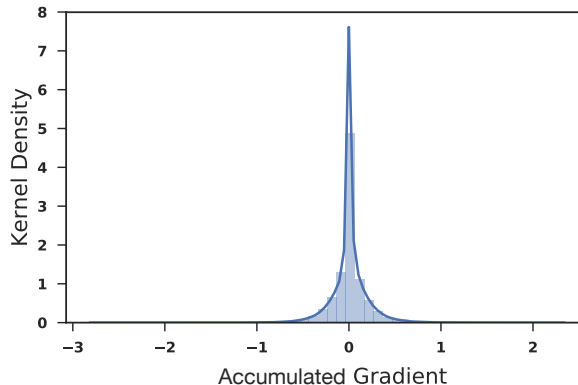


Figure 1. Distribution of accumulated gradients over 100 epochs of standard SGD training on MNIST using a 90,000-weight MLP.

the weights with the highest values. However, this naïve approach is not effective during the first few training iterations. Informally, this is because the initial value of each weight, typically drawn from a scaled normal distribution (LeCun et al., 1998), serves as scaffolding which gradient descent can amplify to train the network.

Instead, our approach is to track gradients for a fixed number of weights which have learned the most overall — that is, the weights with the *highest accumulated gradients*. Figure 1 shows that most weights move very little from their initial values and most accumulated gradients are near 0, suggesting that we only need to track gradients for a small fraction of the weights — *provided* we keep the remaining weights at their initial values.

When untracked weights are kept at their initialization values (see below), this approach results in a weight diffusion speed that is very close to that of the baseline SGD, a key factor in maintaining good generalization under different training regimes (Hoffer et al., 2017).

**Recompute initialization-time values for untracked weights.** While it may be intuitive to simply set untracked weights to 0, we observed that preserving the scaffolding provided by the initialization values is critical to the accuracy of the trained network. In our experiments on MNIST, we were able to reduce the tracked weights 60 $\times$  if initialization values were preserved, but only 2 $\times$  if untracked weights were zeroed. This is in line with the fact that prior pruning approaches require retraining after a subset of weights have been set to 0.

However, storing the initialization-time values would require accessing off-chip memory to retrieve them all during both the forward and backward passes of training — a costly proposition when a memory access consumes upwards of 700 $\times$  more energy than a floating-point operation.

To avoid storing these weights, we observe that in practice the initial values are generated using a pseudo-random number source that is initialized using a single seed value and postprocessed to fit a scaled normal distribution. Because each value only depends on the seed value and its index, it can be deterministically regenerated exactly when it is needed for computation, without ever being stored in memory.<sup>1</sup> Recomputing a normally distributed pseudo-random initialization value using the xorshift algorithm (Marsaglia, 2003) requires six 32-bit integer operations and one 32-bit floating point operation; this amounts to about 1.5pJ in a 45nm process, 427× less energy than a single off-chip memory access.

Regenerating untracked parameters to their initial values also works out-of-the-box for layers like Batch Normalization or Parametric ReLU, where the initialization strategy is typically a constant value (xorshift is not used for these). These layers are also pruned by DropBack, which to the best of our knowledge is unique.

### Freeze the set of tracked weights after a few epochs.

During training, gradients are still *computed* for the untracked weights, and those gradients can exceed the accumulated gradients of any weights that are being tracked; this is especially likely during the initial phases of training, when the optimization algorithm seeks the most productive direction. While the energy needed to compute the gradient is not significant, replacing the lowest tracked gradients with newly computed ones would require additional memory references, which in turn would expend additional energy.

Once the network has been trained for a few epochs, however, we would expect the accumulated gradients for the tracked weights to exceed any “new” gradients that could arise from the untracked weights. To verify this intuition, we trained a 90,000-parameter MLP on MNIST using standard SGD while keeping track of which parameters were in the top-2K gradients set. Figure 2 shows that the set of the highest-gradient weights stabilizes after the first few iterations. The “noise” of less than 0.04% of weights entering and leaving the highest-gradient set in the rest of the epochs remains throughout the training process, and has no effect on the final accuracy. This observation allows us to freeze the “tracked” parameter set after a small number of epochs, saving more energy-costly memory accesses.

## 2.2 The DropBack algorithm

Algorithm 1 shows the resulting DropBack training process. Weights are initialized from a scaled normal distribution (LeCun et al., 1998), generated by the xorshift pseudo-random

<sup>1</sup>In a processor or GPU chip, short-lived temporary values are stored in on-chip register files or scratchpads, reading which costs a fraction of the energy needed to access off-chip DRAM.

number generator. For clarity of exposition, the listing recomputes all gradients and sorts them to determine the highest-total-gradient set of  $k$  elements; in a practical implementation the tracked accumulated gradient set is stored a priority queue of size  $k$ , with incoming gradients higher than the stored minimum evicting the minimum elements. After a manually chosen iteration cutoff, the tracked set is “frozen” and gradients are only updated for the weights already tracked; this saves additional computation time and energy. Note that the tracked set  $T$  requires no storage: its elements are recomputed when needed from  $W^{t-1} - W^{(0)}$ . As in standard stochastic gradient descent, training ends once the network is considered to have converged.

---

**Algorithm 1:** DropBack training.  $N(0, \sigma)$  is generated from the xorshift pseudo-RNG.  $W_{trk}$  and  $W_{utr}$  = tracked and untracked weights;  $T$  and  $U$  = tracked and untracked accumulated gradients;  $S$  = sorted accumulated gradients;  $k$  = number of gradients to track and  $\lambda$  = lowest tracked cumulative gradient;  $\alpha$  = learning rate.  $mask$  indicates a boolean matrix with the same shape as the weights.  $\overline{mask}$  indicates the logical inverse of the mask.

---

**Initialization:**  $W^{(0)}$  with  $W^{(0)} \sim N(0, \sigma)$

**Output:**  $W^{(t)}$

**while not converged do**

$$T = \left\{ \left\| \sum_{i=0}^{t-1} \frac{\alpha \partial f(W^{(i-1)}; x^{(i-1)})}{\partial w} \right\| \text{ s.t. } w \in W_{trk} \right\}$$

if not frozen:

$$U = \left\{ \left\| \frac{\alpha \partial f(W^{(i-1)}; x^{(i-1)})}{\partial w} \right\| \text{ s.t. } w \in W_{utr} \right\}$$

else:

$$U = \{\}$$

$$S = \text{sort}(T \cup U)$$

$$\lambda = S_k$$

$$mask = \mathbb{1}(S > \lambda)$$

$$W^{(t)} = \overline{mask} \cdot (W^{(t-1)} - \alpha \nabla f(W^{(t-1)}; x^{(t-1)})) + \overline{mask} \cdot W^{(0)}$$

$$t = t + 1$$

**end**

---

**Differences from sparsity-based regularization.** Note that DropBack is very different from techniques like Dropout (Srivastava et al., 2014) or DSD (Han et al., 2017), which *temporarily* restrict the gradients that can be updated, essentially as regularization technique. Dropout restricts randomly selected gradient updates during each training iteration, while DSD repeatedly alternates sparse phases (where the lowest-absolute-value weights are deleted) and dense refinement phases (where all weights may be updated). In contrast, DropBack specifically restricts updates to gradients that have not substantially contributed to the overall optimization gradient, and does not involve retraining phases

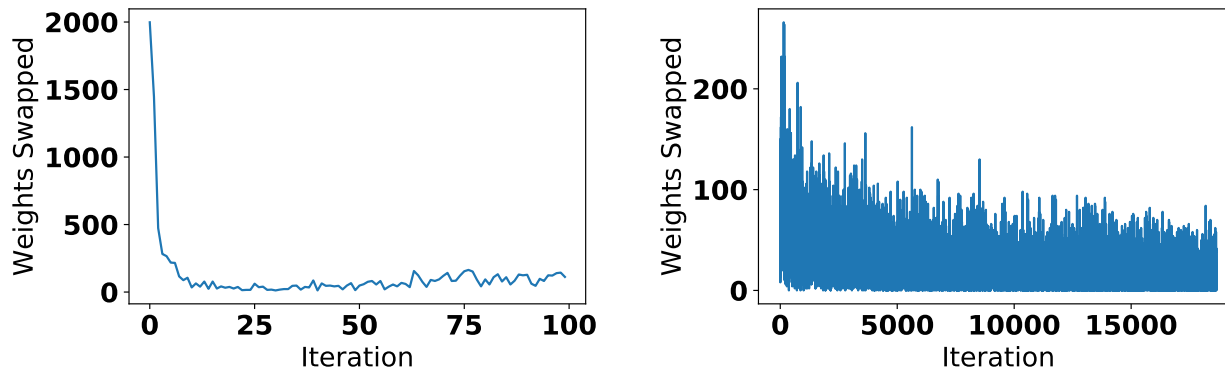


Figure 2. The number of weights added/removed to the top-2K gradient set in the first 10 mini-batches (left) and the remaining mini-batches (right) on MNIST. (Note different y-axis scales).

on either the pruned or dense network.

### 3 EXPERIMENTS

#### 3.1 Methods

We implemented DropBack using the Chainer deep neural network toolkit (Tokui et al., 2015); models were trained on an NVIDIA 1080Ti GPU. We compared DropBack to a baseline implementation without any pruning, as well as three representative pruning techniques:

- a straightforward magnitude-based pruning implementation where only the highest weights are kept after each iteration;
- variational dropout (Kingma et al., 2015), which can progressively weights during training; and
- network slimming (Liu et al., 2017), a modern train-prune-retrain pruning method that achieves state-of-the-art results on modern network architectures.

All networks were optimized using stochastic gradient descent with momentum, as all other optimization strategies cost significant extra memory.

We evaluated all techniques on the MNIST (LeCun, 1998), CIFAR-10 (Krizhevsky, 2009), and ImageNet (Russakovsky et al., 2015) datasets. For MNIST, we used both LeNet-300-100 (Lecun et al., 1998) and a simpler network with only 100 hidden units, which we refer to as MNIST-100-100. For CIFAR-10, we used three networks: Densenet (Huang et al., 2016), WRN-28-10 (Zagoruyko & Komodakis, 2016), and VGG-S, a reduced VGG-16-like model with dropout, batch normalization, and two FC layers of 512 neurons including the output layer (a total of 15M parameters vs. the 138M of VGG-16). We specifically chose Densenet and WRN because

they represent modern network architectures that are very challenging to prune with existing techniques (Liu et al., 2017; Louizos et al., 2017; Li et al., 2017). For ImageNet, we trained the ResNet18 (He et al., 2016) model; as we have been unable to access the test set evaluation server for ImageNet, we report results on the validation set.

In what follows, we report reductions in network size as the ratio of the unpruned weight count to the number of weights tracked during training. This does not include other reduction techniques such as quantization or compression, which are orthogonal to our approach.

#### 3.2 MNIST digit recognition

We first evaluated DropBack on the MNIST handwritten digits dataset using a small multi-layer perceptron (MLP) with approximately 90,000 weights, as well as the LeNet-300-100 MLP, which has approximately 266,600 weights. Training was allowed for up to 100 epochs, and the initial learning rate of 0.4 was exponentially reduced four times by a factor of 0.5. The best epoch was chosen by highest validation accuracy after 5 epochs of no improvement.

**Weight reduction and accuracy.** Table 1 shows the results for the baseline (unpruned) model and three configurations of DropBack, retaining respectively 50,000 weights ( $1.8\times$  reduction), 20,000 weights ( $4.5\times$  reduction), and 1,500 weights ( $60\times$  reduction). With MNIST-100-100 and a modest  $2\times$  reduction in weights, DropBack slightly exceeds the accuracy of the baseline model. This matches the trend reported in prior work such as DSD (Han et al., 2017), where a sparse model that omits 30%–50% weights outperforms the baseline dense (all-weights) model. However, DSD first trains the network to convergence on the *complete* parameter set, and only then prunes some weights and retrain the resulting sparse network. In contrast, DropBack achieves improvement *without* any training in the dense configuration,

**DropBack: full DNN training on a pruned weight budget**

MNIST-300-100	Validation Error	Weight reduction	Best Epoch	Freeze Epoch
Baseline 267K	1.41%	0×	65	N/A
DropBack 50k	1.51%	5.33×	24	100
DropBack 5k	2.58%	53.32×	32	20
DropBack 1.5k	3.84%	177.74×	97	40
MNIST-100-100	Validation Error	Weight reduction	Best Epoch	Freeze Epoch
Baseline 90K	1.70%	0×	47	N/A
DropBack 50k	1.58%	1.8×	24	5
DropBack 20k	1.70%	4.5×	32	5
DropBack 1.5k	3.78%	60×	26	30

Table 1. The MNIST digit dataset using LeNet-300-100 (top) and MNIST-100-100, a smaller MLP with 100 hidden neurons (bottom). DropBack 50k refers to a configuration where 50,000 gradients are retained during training, DropBack 5k to a configuration with 5,000 retained gradients, and so on.

and without repeatedly retraining the network to convergence. The larger MLP sees a slight drop in accuracy, but at 50,000 tracked weights is compressed many more times, and reaches maximum accuracy nearly 3× faster. Further reducing the model to 20,000 weights (4.5× reduction) results in nearly the same accuracy as the baseline, and convergence in a comparable number of epochs.

We also investigated an extreme reduction configuration with weight storage reduced drastically to 1,500 weights. Not surprisingly, the error rate increases (over 2×) and twice the number of epochs are needed to achieve convergence. The nearly 60× reduction in the weight storage requirement for the smaller model and 177× for the larger model potentially offers an attractive design point for low-power embedded accelerators in future mobile and edge devices.

**Tracked weight set freezing.** Figure 3 shows the effect of different freezing times for the tracked gradient set versus the tracked gradient counts. For both MNIST networks, freezing sooner to reduce the computational overhead results in lower achieved accuracy — especially for very high reduction ratios — but for smaller reduction ratios freezing early has little effect on the overall accuracy. Figure 4 shows the rate of convergence for DropBack and the baseline on the LeNet-300-100 network; despite different tracked parameter counts, both methods have similar convergence behaviour.

**Retained weight distribution.** Table 2 shows that the number of parameters retained per layer varies depending on the number of tracked weights. The smaller DropBack 1.5K network allocates a much higher amount of its weights to the later layers compared to the DropBack 10K network and the baseline — an indication that, in the smaller network, proportionally more neurons in the later layers are critical for decision-making.

### 3.3 CIFAR-10 image classification

We also studied the training performance of DropBack using VGG-S, Densenet, and WRN-28-10 on the CIFAR-10 dataset. This is a much more challenging task than MNIST, and the networks were specifically selected for being already relatively dense. Models were trained for 300 epochs on VGG-S and 500 epochs on Densenet and WRN; the best epoch was selected. No data augmentation of CIFAR-10 was performed. The learning rate was initialized to 0.4 and decayed 0.5× every 25 epochs.

**Weight reduction and accuracy.** Table 3 shows how DropBack compares to variational dropout, network slimming, and magnitude-based pruning on VGG-S, Densenet, and WRN-28-10. Overall, DropBack is able to achieve comparable (or even slightly improved) accuracy on VGG-S and Densenet with five-fold weight reduction, and up to 20×–30× if some accuracy is sacrificed. On WRN-28-10, DropBack achieves 5× and 7× reduction with less than 0.5% accuracy drop.

Note that these networks are challenging, as they are already quite dense for the accuracy level they achieve. Variational dropout works well only on VGG-S, and fails to converge on Densenet and WRN. Magnitude-based pruning does not achieve better accuracy than DropBack despite less weight reduction. Finally, network slimming achieves slightly better top accuracy on Densenet with 50% less reduction, but results in dramatic accuracy loss when applied to WRN; this corresponds to recent work which has shown that WRN is hard to compress more than about 2× without losing significant accuracy (Liu et al., 2017; Louizos et al., 2017; Li et al., 2017). DropBack, in contrast, is universally able to achieve weight reduction on the order of 5× with little to no accuracy loss.

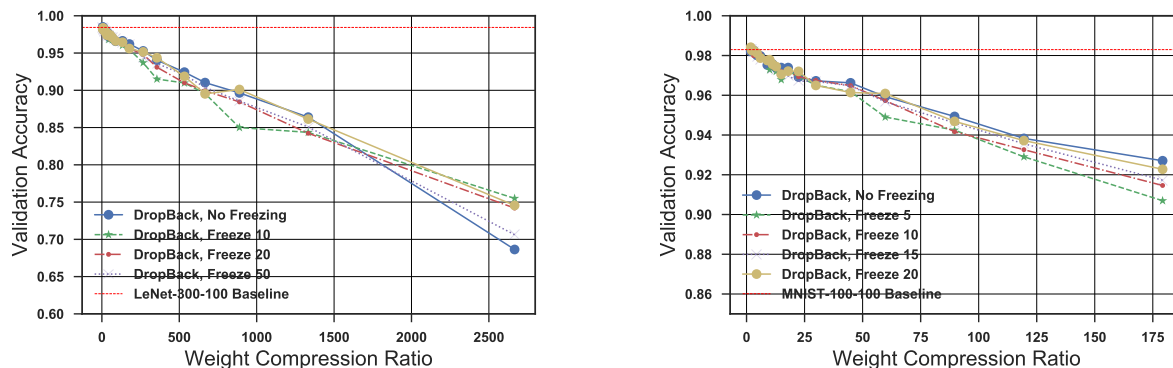


Figure 3. Tracked parameter tradeoff vs. network accuracy, with several parameter selection freezing points for LeNet-300-100 (left) and MLP-100-100 (right). Freeze 0 represents a network trained without freezing.

layer	DropBack 1500	DropBack 10000	Baseline
fc1 (100×784)	734 (52.3%)	7223 (72.2%)	78500 (87.6%)
fc2 (100×100)	512 (34.1%)	2128 (21.3%)	10100 (11.3%)
fc3 (100×10)	254 (16.9%)	549 (5.5%)	1010 (1.1%)
Total	1500	10000	89610

Table 2. Number of gradients for each layer retained in the final trained MNIST network

CIFAR-10	Validation error	Weight reduction	Best epoch	Freeze epoch
Baseline 15M	10.08%	0×	214	NA
VGG-S DropBack 5M	9.75%	3×	127	5
VGG-S DropBack 3M	9.90%	5×	128	20
VGG-S DropBack 0.75M	13.49%	20×	269	35
VGG-S DropBack 0.5M	20.85%	30×	201	15
VGG-S Var. Dropout	13.50%	3.4×	200	NA
VGG-S Mag. Pruning .80	9.42%	5.0×	182	NA
VGG-S Slimming	11.08%	3.8×	196	NA
Densenet Baseline 2.7M	6.48%	0×	382	NA
Densenet DropBack 600k	5.86%	4.5×	409	NA
Densenet DropBack 100k	9.42%	27×	307	NA
Densenet Var. Dropout	9.0%	NA×	NA	NA
Densenet Mag. Pruning .75	6.41%	4.0×	480	NA
Densenet Slimming	5.65%	2.9×	NA	NA
WRN-28-10 Baseline 36M	3.75%	0×	326	NA
WRN-28-10 DropBack 8M	3.85%	4.5×	384	NA
WRN-28-10 DropBack 7M	4.02%	5.2×	417	NA
WRN-28-10 DropBack 5M	4.20%	7.3×	304	NA
WRN-28-10 Var. Dropout	9.0%	NA×	NA	NA
WRN-28-10 Mag. Pruning .75	26.52%	4×	109	NA
WRN-28-10 Slimming .75	16.640%	4×	173	NA

Table 3. Validation accuracy and reduction ratio of several pruned networks on CIFAR-10

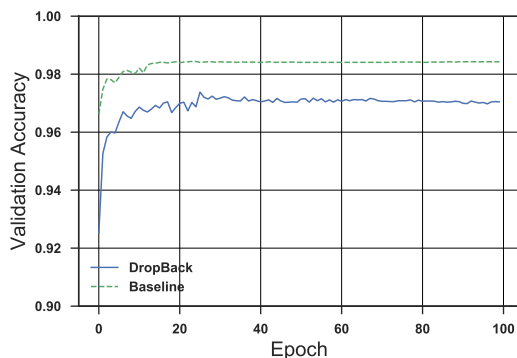


Figure 4. Rate of convergence for LeNet-300-100 for our technique and the baseline model. Note that the y-axis starts as 0.90, and are within 1% of each other.

**Effects of freezing.** Figure 5(right) shows how freezing the tracked parameter set after 10, 20, and 50 epochs affects the reduction and accuracy on VGG-S. At under 10 $\times$  weight pruning, accuracy does not suffer against the uncompressed baseline, and freezing does not affect either reduction or accuracy. With more reduction, however, the tracked gradient set struggles to maintain accuracy, and freezing the gradients early results in substantial accuracy drops at 30 $\times$  reduction.

**Convergence.** Figure 5(left) shows that DropBack initially learns slightly more slowly than the uncompressed baseline, but exhibits the same convergence behaviour after about 20 epochs (VGG-S). Variational dropout, in contrast, learns more quickly initially but converges on a substantially lower accuracy.

### 3.4 ImageNet image classification

Finally, we studied whether DropBack can be applied in the more realistic setting of the ImageNet dataset. This dataset is substantially more challenging than CIFAR-10 — there are 1,000 instead of 10 categories and the images are 64 $\times$  larger — and represents the high end of tasks for which inference is possible in a mobile-device setting.

ResNet18 was trained for 100 epochs using stochastic gradient descent with momentum, and the best epoch was selected. The initial learning rate of 0.2 was decayed by 0.97 $\times$  in each epoch. Freezing was not applied during DropBack training, and no data augmentation was performed.

**Weight reduction and accuracy.** Table 4 and Figure 6(left) show top-1 results on the ImageNet dataset using the ResNet18 network (11 million weights uncompressed) for different weight reduction ratios. DropBack is able to match or improve the accuracy of the baseline uncompressed

network with weight reduction ratios of up to 11.7 $\times$ , and can reduce the weight count over 23 $\times$  with a modest drop in accuracy (39.5% vs 31.6% baseline). In comparison, magnitude-based pruning stops improving after epoch 6 and never matches baseline accuracy: even a 4 $\times$  reduction in weight count incurs the same cost in accuracy as DropBack pays for a 23.4 $\times$  reduction.

**Convergence.** Figure 6(right) shows that DropBack learns at the same rate as the uncompressed baseline up to 5.85 $\times$  weight reduction. At 11.7 $\times$  weight reduction, the model initially learns more slowly and takes substantially longer to converge, but eventually matches the baseline accuracy.

## 4 DISCUSSION

To investigate why DropBack consistently achieves good reduction/accuracy tradeoffs across a wide range of networks, we considered the training process analysis due to Hoffer et al. (2017). Briefly, they observe that when DNNs are trained using SGD, the  $\ell^2$  distance of weights  $\mathbf{w}_t$  at time  $t$  from the initial weights  $\mathbf{w}_0$  grows logarithmically, i.e.,  $\|\mathbf{w}_t - \mathbf{w}_0\| \sim \log t$  (this is the same as the gradient accumulated until time  $t$ ). They therefore model SGD as a random walk on a random potential surface, which exhibits the same logarithmic distance effect (known as ultra-slow diffusion). Hoffer et al. (2017) argue that SGD configurations (for them, batch sizes) that preserve the ultra-slow diffusion effect result in models that generalize well.

We reasoned that DropBack maximally preserves the  $\ell^2$  diffusion distance of the baseline training scheme because (a) DropBack tracks the highest gradients, and (b) most of the remaining gradients are very close to zero (cf. Fig. 1). To verify this intuition, we measured the diffusion distance for the baseline uncompressed network, DropBack, variational dropout, and magnitude-based pruning, all on the MNIST-100-100 network.<sup>2</sup> Figure 7(left) shows that under DropBack weights diffuse very similarly to the baseline training scheme, with the overall  $\ell^2$  distance negligibly lower because the untracked weights remain at their initialization values. In contrast, magnitude-based pruning begins with a large  $\ell^2$  distance (because many initialization weights are zeroed), and does not provide enough scaffolding structure for SGD to train well. Finally, variational dropout favours dropping individual weights, and so diffuses much faster than baseline and DropBack, which corresponds to its failure to generalize on the denser networks and to ImageNet (see Section 3).

To visualize how the weight values evolve under DropBack compared to the baseline and the two pruning techniques, we projected the parameter space down to 3D using principal

<sup>2</sup>Network slimming, being a train-prune-retrain technique, is not amenable to this type of analysis.

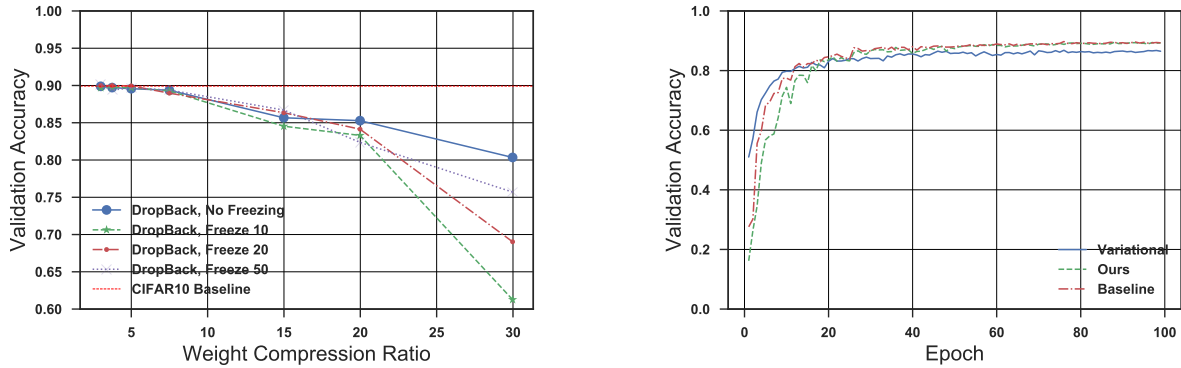


Figure 5. VGG-S CIFAR-10. Left: tracked parameter tradeoff vs. network accuracy, with several parameter selection freezing points (note the y-axis range). Right: epoch vs. validation accuracy for our method at 5M tracked parameters, variational dropout, and the baseline model.

ResNet18	Validation error	Weight reduction	Best epoch
Baseline 11M	31.59%	1.00×	26
DropBack 4M	30.31%	2.92×	41
DropBack 2M	29.95%	5.85×	44
DropBack 1M	32.01%	11.69×	49
DropBack 0.5M	39.53%	23.39×	47
Var. Dropout	FAIL	FAIL	FAIL
Mag. Pruning 2.75M	38.43%	4×	6

Table 4. Validation accuracy and reduction ratio of DropBack applied to ResNet18 on the ImageNet dataset. DropBack is able to reduce the weight count by 11.7× without accuracy loss, and nearly 24× with some degradation. In comparison, magnitude-based pruning at 4× reduction suffers the same loss as DropBack at 23.4×. All of our Variational Dropout experiments on ResNet18 failed due to numerical instability.

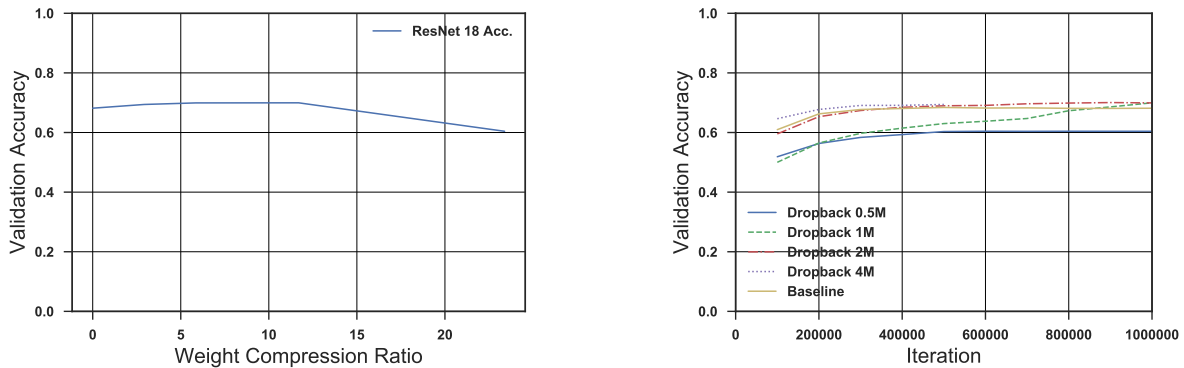


Figure 6. ResNet18 on ImageNet. Left: tracked parameter tradeoff vs. network accuracy, with several parameter selection freezing points (note the y-axis range). Right: epoch vs. validation accuracy at different weight reduction ratios, as well as the uncompressed baseline (which has 11M weights).



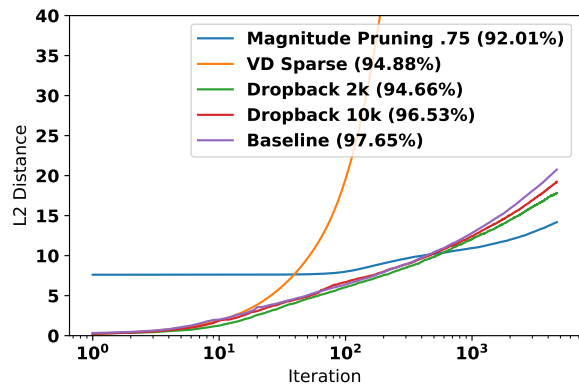


Figure 7. Diffusion ( $\ell^2$ ) distance vs. training time on MNIST-100 (note log time scale). Unlike magnitude-based pruning and variational dropout, DropBack weights remain at nearly the same  $\ell^2$  distance from their initial point as the uncompressed baseline throughout the training process.

component analysis. Figure 8 shows that under DropBack, the principal components of the trained weight vector stay very close to those of the baseline-trained weight vector, whereas those of magnitude-based pruning and variational dropout diverge significantly. If we imagine the training path of the baseline uncompressed configuration to be optimal, DropBack results in a near-optimal evolution.

## 5 RELATED WORK

### 5.1 Pruning

Pruning networks to enforce sparsity after training has been effective in Han et al. (2015a); Zhu et al. (2017); Ge et al. (2017); Masana et al. (2017); Luo et al. (2017); Ullrich et al. (2017); LeCun et al. (1990); Srivastava et al. (2014); Wan et al.; Yang et al. (2017), while other work has focused on reducing the rank of the parameter matrices during training (Alvarez & Salzmann, 2017) for later reduction. In contrast to DropBack, pruning post-training requires a retraining step to regain lost accuracy, and both pruning and low rank constraints still require full dense backpropagation during the training phase; the memory and energy cost of the extra backpropagation step are undesirable in embedded systems. DropBack requires no retraining steps, and during both training and inference only ever stores the small subset of weights that are tracked. As a result, it never incurs the memory access costs of backpropagation during the training process.

Other work has focused on pruning while training to either improve accuracy or to increase the eventual sparsity of the trained network. Zhu & Gupta (2017) gradually increase the number of weights masked from contributing to the network, while Molchanov et al. (2017) extend variational

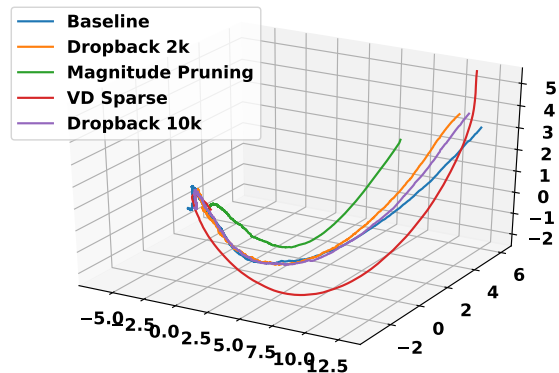


Figure 8. Evolution of weights under SGD projected into 3D space using principal component analysis for DropBack, baseline, magnitude based pruning, and variational dropout. DropBack remains closest to the uncompressed baseline.

dropout (Kingma et al., 2015) with per-parameter dropout rates to increase sparsity. (Babaeizadeh et al., 2016) inject random noise into a network to find and merge the most correlated neurons. Finally, (Langford et al., 2008) decay weights every  $k$  steps (for a somewhat large value of  $k$ ), inducing sparsity gradually. Unlike DropBack, all of these techniques require at least as much memory to train as the unpruned network initially, and all of them take longer to converge.

### 5.2 Quantization

Reducing storage costs is often approached through quantizing weights to smaller bit widths. This can be performed either after training, as in Ge et al. (2017); Wu et al. (2015); Choi et al. (2016); Han et al. (2015a); Ullrich et al. (2017); Gysel (2016), during an iterative retraining process (Zhou et al., 2017), or during training, as in Cai et al. (2017); Zhou et al. (2016); Courbariaux et al. (2016); Rastegari et al. (2016); Gupta et al. (2015); Mishra et al. (2017); Hubara et al. (2016); Courbariaux et al. (2014); Simard & Graf (1994); Holt & Baker. Out of all of these methods, only Gupta et al. (2015); Mishra et al. (2017); Hubara et al. (2016); Courbariaux et al. (2014) use reduced precision *while training* to lower the train time storage costs; other methods other methods store the full, non-quantized weights during backpropagation just like the post-training methods. Quantization is orthogonal to, and has been combined with, pruning techniques (Han et al., 2015a); as such, it is also orthogonal to DropBack.

### 5.3 Other compression methods

A few compression methods do not naturally fall into the quantization categories. HashedNets (Chen et al., 2015) use a

hash function to group neuron connections into buckets with the same weight value. In effect, this is a kind of quantization, where the quantized values are limited in number but each is a full 32-bit floating-point number. Huffman coding (Huffman) has also been used for compression: for example, Deep Compression (Han et al., 2015a) applies Huffman coding to a pruned and quantized network. Both techniques require the full network to be trained first, and, like quantization, are orthogonal to DropBack.

## 6 CONCLUSION

In this paper, we introduce DropBack, a pruning technique that reduces weight storage both *during* and *after* training by (a) tracking only the weights with the highest accumulated gradients, and (b) recomputing the remaining weights on the fly.

Because DropBack prunes exactly those weights that have learned the least, its weight diffusion profile during training is very close to that of standard (unconstrained) SGD, in contrast to other techniques. This allows DropBack to achieve better accuracy and weight reduction than prior methods on dense modern networks like Densenet, WRN, and ResNet, which have proven challenging to prune using existing techniques: a  $5\times$ – $11\times$  weight count reduction with no accuracy loss.

Perhaps most attractively, DropBack is the first technique to dramatically reduce the memory footprint and memory bandwidth needed to store and access weights *during* training. Because of this, DropBack can potentially be used to train networks  $5\times$ – $10\times$  larger than currently possible with typical hardware, or to train/retrain standard-size networks on small mobile and embedded devices, something not possible with current training techniques.

## REFERENCES

- Albericio, J., Judd, P., Hetherington, T., Aamodt, T., Jerger, N. E., and Moshovos, A. Cnvlutin: ineffectual-neuron-free deep neural network computing. In *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*, pp. 1–13. IEEE, 2016.
- Alvarez, J. M. and Salzmann, M. Compression-aware Training of Deep Networks. *arXiv:1711.02638 [cs]*, November 2017. arXiv: 1711.02638.
- Apple. An On-device Deep Neural Network for Face Detection. *Apple Machine Learning Journal*, November 2017.
- Babaeizadeh, M., Smaragdis, P., and Campbell, R. H. Noise-Out: A simple way to prune neural networks. 2016.
- Cai, Z., He, X., Sun, J., and Vasconcelos, N. Deep Learning with Low Precision by Half-wave Gaussian Quantization. *arXiv:1702.00953 [cs]*, February 2017. arXiv: 1702.00953.
- Chen, T., Xu, B., Zhang, C., and Guestrin, C. Training Deep Nets with Sublinear Memory Cost. *arXiv:1604.06174 [cs]*, April 2016. arXiv: 1604.06174.
- Chen, W., Wilson, J. T., Tyree, S., Weinberger, K. Q., and Chen, Y. Compressing neural networks with the hashing trick. *CoRR*, abs/1504.04788, 2015. URL <http://arxiv.org/abs/1504.04788>.
- Choi, Y., El-Khamy, M., and Lee, J. Towards the Limit of Network Quantization. *arXiv:1612.01543 [cs]*, December 2016. arXiv: 1612.01543.
- Courbariaux, M., Bengio, Y., and David, J.-P. Training deep neural networks with low precision multiplications. *arXiv:1412.7024 [cs]*, December 2014. arXiv: 1412.7024.
- Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., and Bengio, Y. Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1. *arXiv:1602.02830 [cs]*, February 2016. arXiv: 1602.02830.
- Ge, S., Luo, Z., Zhao, S., Jin, X., and Zhang, X. Y. Compressing deep neural networks for efficient visual inference. In *2017 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 667–672, July 2017. doi: 10.1109/ICME.2017.8019465.
- Gupta, S., Agrawal, A., Gopalakrishnan, K., and Narayanan, P. Deep learning with limited numerical precision. 2015.
- Gysel, P. Ristretto: Hardware-Oriented Approximation of Convolutional Neural Networks. *arXiv:1605.06402 [cs]*, May 2016. arXiv: 1605.06402.
- Han, S., Mao, H., and Dally, W. J. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. *arXiv:1510.00149 [cs]*, October 2015a. arXiv: 1510.00149.
- 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, 2015b.
- 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 *ACM/IEEE International Symposium on Computer Architecture (ISCA)*, 2016.
- Han, S., Pool, J., Narang, S., Mao, H., Gong, E., Tang, S., Elsen, E., Vajda, P., Paluri, M., Tran, J., Catanzaro, B.,

- 550 and Dally, W. J. DSD: Dense-Sparse-Dense Training for  
551 Deep Neural Networks. In *International Conference on*  
552 *Learning Representations (ICLR)*, 2017.
- 553
- 554 Hassibi, B., Stork, D. G., and Wolff, G. J. Optimal brain sur-  
555 geon and general network pruning. In *IEEE International*  
556 *Conference on Neural Networks, 1993.*, pp. 293–299.  
557 IEEE, 1993.
- 558
- 559 He, K., Zhang, X., Ren, S., and Sun, J. Deep Residual  
560 Learning for Image Recognition. In *The IEEE Conference*  
561 *on Computer Vision and Pattern Recognition (CVPR)*,  
562 June 2016.
- 563
- 564 Hoffer, E., Hubara, I., and Soudry, D. Train longer, gener-  
565 alize better: closing the generalization gap in large batch  
566 training of neural networks. In *NIPS*, 2017.
- 567
- 568 Holt, J. L. and Baker, T. E. Back propagation simulations  
569 using limited precision calculations. In *IJCNN-91-Seattle*  
570 *International Joint Conference on Neural Networks*, vol-  
571 ume ii, pp. 121–126 vol.2.
- 572
- 573 Huang, G., Liu, Z., van der Maaten, L., and Weinberger,  
574 K. Q. Densely Connected Convolutional Networks.  
575 *arXiv:1608.06993 [cs]*, August 2016. URL [http://](http://arxiv.org/abs/1608.06993)  
576 [arxiv.org/abs/1608.06993](http://arxiv.org/abs/1608.06993). arXiv: 1608.06993.
- 577
- 578 Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., and  
579 Bengio, Y. Quantized neural networks: Training neural  
580 networks with low precision weights and activations.  
581 2016.
- 582
- 583 Huffman, D. A. A method for the construction of minimum-  
584 redundancy codes. 40(9):1098–1101. ISSN 0096-8390.  
585 doi: 10.1109/JRPROC.1952.273898.
- 586
- 587 Intel. Intel® Movidius™ Myriad™ X VPU, 2017.
- 588
- 589 Judd, P., Delmas, A., Sharify, S., and Moshovos, A. Cn-  
590 vlutin?: Ineffectual-activation-and-weight-free deep neural  
591 network computing. *arXiv preprint arXiv:1705.00125*,  
592 2017.
- 593
- 594 Kingma, D. P., Salimans, T., and Welling, M. Variational  
595 dropout and the local reparameterization trick. 2015.
- 596
- 597 Kingsley-Hughes, A. Inside Apple’s new A11 Bionic pro-  
598 cessor. ZDNet, September 2017.
- 599
- 600 Krizhevsky, A. Learning multiple layers of features from  
601 tiny images. 2009.
- 602
- 603 Langford, J., Li, L., and Zhang, T. Sparse online learning  
604 via truncated gradient. 2008.
- 605
- 606 LeCun, Y., Denker, J. S., and Solla, S. A. Optimal brain  
607 damage. In Touretzky, D. S. (ed.), *Advances in Neural*  
608 *Information Processing Systems 2*, pp. 598–605. Morgan-  
609 Kaufmann, 1990.
- 610
- 611 Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-  
612 based learning applied to document recognition. 86(11):  
613 2278–2324, 1998. ISSN 0018-9219.
- 614
- 615 LeCun, Y., Bottou, L., Orr, G. B., and Müller, K.-R. Efficient  
616 BackProp. In *Neural Networks: Tricks of the Trade, This*  
617 *Book is an Outgrowth of a 1996 NIPS Workshop*, pp.  
618 9–50, London, UK, UK, 1998. Springer-Verlag. ISBN  
619 3-540-65311-2.
- 620
- 621 Li, H., De, S., Xu, Z., Studer, C., Samet, H., and Goldstein,  
622 T. Training Quantized Nets: A Deeper Understanding.  
623 *arXiv:1706.02379 [cs, stat]*, June 2017. URL [http://](http://arxiv.org/abs/1706.02379)  
624 [arxiv.org/abs/1706.02379](http://arxiv.org/abs/1706.02379). arXiv: 1706.02379.
- 625
- 626 Liu, Z., Li, J., Shen, Z., Huang, G., Yan, S., and Zhang,  
627 C. Learning Efficient Convolutional Networks through  
628 Network Slimming. *arXiv:1708.06519 [cs]*, August  
629 2017. URL [http://](http://arxiv.org/abs/1708.06519)  
630 [arxiv.org/abs/1708.06519](http://arxiv.org/abs/1708.06519). arXiv:  
631 1708.06519.
- 632
- 633 Louizos, C., Welling, M., and Kingma, D. P. Learning  
634 Sparse Neural Networks through  $\$L_0\$$  Regularization.  
635 *arXiv:1712.01312 [cs, stat]*, December 2017. URL [http://](http://arxiv.org/abs/1712.01312)  
636 [arxiv.org/abs/1712.01312](http://arxiv.org/abs/1712.01312). arXiv: 1712.01312.
- 637
- 638 Luo, J.-H., Wu, J., and Lin, W. ThiNet: A Filter Level  
639 Pruning Method for Deep Neural Network Compression.  
640 *arXiv:1707.06342 [cs]*, July 2017. arXiv: 1707.06342.
- 641
- 642 Marsaglia, G. Xorshift rngs. *Journal of Statistical Software*,  
643 8(14):1–6, 2003.
- 644
- 645 Masana, M., van de Weijer, J., Herranz, L., Bagdanov,  
646 A. D., and Alvarez, J. M. Domain-adaptive deep network  
647 compression. *arXiv:1709.01041 [cs]*, September 2017.  
648 arXiv: 1709.01041.
- 649
- 650 Masters, D. and Luschi, C. Revisiting Small Batch Training  
651 for Deep Neural Networks. 2018. URL [http://](http://arxiv.org/abs/1804.07612)  
652 [arxiv.org/abs/1804.07612](http://arxiv.org/abs/1804.07612).
- 653
- 654 Mishra, A., Nurvitadhi, E., Cook, J. J., and Marr, D. WRPN:  
655 Wide reduced-precision networks. 2017.
- 656
- 657 Molchanov, D., Ashukha, A., and Vetrov, D. Variational  
658 dropout sparsifies deep neural networks. 2017.
- 659
- 660 Qualcomm. Snapdragon Neural Processing Engine for AI,  
661 December 2017.
- 662
- 663 Rastegari, M., Ordonez, V., Redmon, J., and Farhadi, A.  
664 XNOR-Net: ImageNet Classification Using Binary Convo-  
665 lutional Neural Networks. *arXiv:1603.05279 [cs]*, March  
666 2016. arXiv: 1603.05279.

- 605 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S.,  
606 Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein,  
607 M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale  
608 Visual Recognition Challenge. *International Journal of*  
609 *Computer Vision (IJCV)*, 115(3):211–252, 2015. doi:  
610 10.1007/s11263-015-0816-y.
- 611 Samsung. Samsung Optimizes Premium Exynos 9 Series  
612 9810 for AI Applications and Richer Multimedia Content,  
613 January 2018.
- 614 Simard, P. Y. and Graf, H. P. Backpropagation without multi-  
615 plication. In *Advances in Neural Information Processing*  
616 *Systems*, pp. 232–239, 1994.
- 617  
618 Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I.,  
619 and Salakhutdinov, R. Dropout: A Simple Way to Prevent  
620 Neural Networks from Overfitting. *Journal of Machine*  
621 *Learning Research*, 15(1):1929–1958, 2014.
- 622  
623 Tokui, S., Oono, K., Hido, S., and Clayton, J. Chainer: a  
624 next-generation open source framework for deep learning.  
625 In *Proceedings of workshop on machine learning systems*  
626 *(LearningSys) in (NIPS)*, volume 5, 2015.
- 627  
628 Ullrich, K., Meeds, E., and Welling, M. Soft Weight-Sharing  
629 for Neural Network Compression. *arXiv:1702.04008 [cs,*  
630 *stat]*, February 2017. arXiv: 1702.04008.
- 631  
632 Wan, L., Zeiler, M., Zhang, S., Cun, Y. L., and Fergus, R.  
633 Regularization of neural networks using DropConnect. In  
634 *PMLR*, pp. 1058–1066.
- 635  
636 Wu, J., Leng, C., Wang, Y., Hu, Q., and Cheng, J. Quant-  
637 ized Convolutional Neural Networks for Mobile De-  
638 vices. *arXiv:1512.06473 [cs]*, December 2015. arXiv:  
1512.06473.
- 639  
640 Yang, T.-J., Chen, Y.-H., and Sze, V. Designing Energy-  
641 Efficient Convolutional Neural Networks Using Energy-  
642 Aware Pruning. In *CVPR*, 2017.
- 643  
644 Zagoruyko, S. and Komodakis, N. Wide Residual Networks.  
645 *arXiv:1605.07146 [cs]*, May 2016. URL <http://arxiv.org/abs/1605.07146>. arXiv: 1605.07146.
- 646  
647 Zhou, A., Yao, A., Guo, Y., Xu, L., and Chen, Y. Incremental  
648 Network Quantization: Towards Lossless CNNs with Low-  
649 Precision Weights. *arXiv:1702.03044 [cs]*, February 2017.  
650 arXiv: 1702.03044.
- 651  
652 Zhou, S., Wu, Y., Ni, Z., Zhou, X., Wen, H., and Zou,  
653 Y. DoReFa-Net: Training Low Bitwidth Convolutional  
654 Neural Networks with Low Bitwidth Gradients.  
*arXiv:1606.06160 [cs]*, June 2016. arXiv: 1606.06160.
- 655  
656 Zhu, J., Jiang, J., Chen, X., and Tsui, C.-Y. SparseNN:  
657 An Energy-Efficient Neural Network Accelerator Exploit-  
658 ing Input and Output Sparsity. *arXiv:1711.01263 [cs]*,  
659 November 2017. arXiv: 1711.01263.
- Zhu, M. and Gupta, S. To prune, or not to prune: exploring  
the efficacy of pruning for model compression. 2017.