
Progressive Stochastic Binarization of Deep Networks

David Hartmann
Institute of Computer Science
Johannes Gutenberg-University of Mainz
Staudingerweg 9, 55128 Mainz, Germany
dahartma@uni-mainz.de

Michael Wand
Institute of Computer Science
Johannes Gutenberg-University of Mainz
Staudingerweg 9, 55128 Mainz, Germany
mwand@uni-mainz.de

Abstract

A plethora of recent research has focused on improving the memory footprint and inference speed of deep networks by reducing the complexity of (i) numerical representations (for example, by deterministic or stochastic quantization) and (ii) arithmetic operations (for example, by binarization of weights).

We propose a stochastic binarization scheme for deep networks that allows for efficient inference on hardware by restricting itself to additions of small integers and fixed shifts. Unlike previous approaches, the underlying randomized approximation is progressive, thus permitting an adaptive control of the accuracy of each operation at run-time. In a low-precision setting, we match the accuracy of previous binarized approaches. Our representation is unbiased — it approaches continuous computation with increasing sample size. In a high-precision regime, the computational costs are competitive with previous quantization schemes. Progressive stochastic binarization also permits localized, dynamic accuracy control within a single network, thereby providing a new tool for adaptively focusing computational attention.

We evaluate our method on networks of various architectures, already pretrained on ImageNet. With representational costs comparable to previous schemes, we obtain accuracies close to the original floating point implementation. This includes pruned networks, except the known special case of certain types of separated convolutions. By focusing computational attention using progressive sampling, we reduce inference costs on ImageNet further by a factor of up to 33% (before network pruning).

1 Introduction

The ability to analyze large amounts of training data has been instrumental to progress in machine learning. This applies in particular to deep networks, whose renaissance has been triggered to some extent by the availability of computational resources [38; 56]. Hence, improving computational efficiency is an important endeavor, not only for using current methods on limited hardware such as embedded and mobile systems, but also for further progress in the field [61].

Current deep networks mostly compute scalar products (weighted sums), i.e., addition and multiplication of real numbers accounts for the largest proportion of computational costs. Of the two named, multiplication is more expensive at a gate-level, as hardware costs grow faster with precision. Nonetheless, floating-point additions also incur substantial costs due to the need for normalization of the mantissa. Thus, hardware vendors have already begun to support specialized low-precision integer and floating-point units. Further specialized hardware has been considered [5; 48]. The question of more elementary representations is similarly posed in “neuromorphic” computing, which targets non-CMOS hardware [57] or even machine learning through chemical and biological processes [19].

On the algorithmic side, recent research has explored different approaches for computational cost reduction. Broadly, these can be classified as quantization [21] (decrease bits per number), binarization [4] (single bit per number), stochastic computing [18] (random bit streams), and pruning [22] (dropping less important computations) approaches.

1.1 Our approach

Our approach combines ideas from stochastic computing, quantization and binarization: We represent activations as (small) integers. Weights are binarized stochastically, and multiplied by simple gating. Repeated addition (accumulation) increases the precision when needed. Such a simple stochastic approximation on its own, however, would lead to substantial variance and, thus, large sampling costs. We therefore introduce an importance sampling scheme that performs random choices between adjacent powers of two. While adding only moderate complexity (random 1-bit shifts), this reduces noise dramatically and makes the scheme competitive with continuous representations already at small sample sizes (16 to 64 gated additions already come close to 32-bit floating point accuracy).

Because this scheme permits trivial local and dynamic accuracy control, we subsequently employ a two-stage algorithm that first computes a rough estimate of accuracy demands (using 8 samples per number) and use the higher accuracy computation only sparsely. This step further reduces sampling costs by factors of 33% (on *ImageNet* [54] with a *ResNet50* (v2) [24]).

Our representation is particularly suited for the inference stage (forward pass), where the sampled factors (powers of two/shifts) are constants. We therefore evaluate our scheme for inference in common deep networks for image classification, including ResNet [24], DenseNet [29], Inception [62], Xception [10], MobileNet [27], and NasNetMobile [71]. Test results approach already at moderate levels of sampling (16 samples) higher precision, and converge to 32-bit floating point calculations for increasing numbers of samples. For the 50-layer ResNet (v2) architecture, for instance, we obtain 94.4% relative accuracy compared to full-precision calculations at 16 random samples and 98.6% at 64 samples. Less favorable convergence behavior occurs only in case of deep stacks of successive multiplications without accumulation. In our experiments, this occurred only for the MobileNet architecture (due the use of a specific type of separable convolution), which is a known issue for other quantization and binarization schemes and can easily be fixed by a slight architectural modification [60]. Simple pruning [22] of the network, which removes redundant computations, does not seem to affect the efficiency of our stochastic approximation scheme.

In summary, our number representation has the following advantages: For the first time, we are able to quantize pretrained networks without retraining and, at the same time, without significant loss of precision. Similarly, it is the first quantization scheme that allows for dynamic control of precision at run-time. This permits an adaptive sampling algorithm that further reduces costs in practice by a substantial factor. Generally, the method is unbiased and convergent to full precision in the sampling limit, therefore permitting high-precision inference on demand by (still moderately) increased computational costs.

2 Related Work

Binarization techniques, and our technique in particular, are at the center of the following topics, each of which has been studied recently.

Quantization By reducing the number of bits per operand, and avoiding expensive floating-point units by utilizing number formats with fixed scaling, computational costs get reduced considerably [14; 21; 68]. Direct quantization of full-precision pre-trained models, however, leads to a significant loss of accuracy in deep networks. Thus, these methods require a fine-tuning step. As a result, most techniques apply quantization during training or as fine-tuning of pre-trained full-precision models. One approach increases the significance of low-precision weights by adaptively rescaling intermediate results globally before and after the weight-multiplication [14]. Instead of optimising the scaling, it is possible to determine optimal global scaling factors for pretrained models by estimating the statistics of incoming data [42].

A second idea uses an iterative discretize-and-retrain technique. Combining this approach with threshold-based pruning produces networks with a much smaller memory footprint [23]. Another discretize-and-retrain algorithm translates the weights into representations consisting of only shift-

operations that are much easier to implement on hardware [67].

More recent works use the *local reparametrization trick* to include stochastic components into a network without losing the ability to estimate gradients [34]. One representative of this technique replaces weights by continuous relaxed random variables; stress on weights is increased not being on a specified discrete grid [45]. Closely related is the idea of sampling normal distributed preactivations during training instead of discrete random weights [49; 59].¹

Our method is an in-place quantization approach, i.e., additional fine-tuning of the model is not required as all weights transform bijectively into their new representation. Compared to other quantization methods, our method allows for choosing the quality of the representation at run-time.

Binarization The limit case of quantization is binarization, where either weights [51] or even both, weights and activations [30] are constrained to one of two possible values. This reduces costs dramatically, eliminating multiplications (binary weights) or even reducing all arithmetics to logical operations (full binarization, which has considerable impact on accuracy [4]).

The first variants used the *straight-through optimizer* [6] to enable gradient descent optimization on discrete variables. Fully binarized network architectures with binarized training processes exist [13; 30], however they tend to show a significant performance gap compared to full-precision networks [70].

Further work concentrates on different number representations. As in quantization, global real-valued scaling factors of binary weights represent numbers with greater variance [51]. Another approach learns number representations with binary coefficients by means of a sum of multiple globally scaled binary weights [43]. More recently, it has been shown that *ResNets* [24] are well suited for purely binary weights, if implemented with full-precision shortcuts [44].

We propose a stochastic binarization that closes the performance gap of binary neural networks in most of our experiments. To the best of our knowledge, other binarization techniques require further hyperparameter tuning or add even more hyperparameters [7]. Our technique changes the number representations bijectively and thus the same hyperparameters apply. One has to mention that our method reevaluates stochastic weights to even out the randomness, while other methods use a single pass during inference. Our experiments show, however, that only a few accumulations are sufficient to achieve similar performance as previous methods. In practice, accumulations can be rolled out and computed in parallel.

Hardware Optimized Architectures The application to optimize networks for the use on embedded or mobile systems has led to architectures that allow for a trade off between latency and accuracy [27; 55]. Or, to go into low-level optimizations, discussions on new numerical representations exist [25]. Ternary quantization² [69], two-bit quantization [47] and integer only multiplications [31]³ are related interpretations of the restrictions of binarization that also allow for better implementation on hardware.

Concrete implementations prove faster inference of networks if implemented directly on hardware [37; 48]. Publications that discuss binary neural networks on FPGAs [41], improved inference on CPUs using SIMD-Instructions [64] and stochastic binary neural networks for near-sensor computing [39] show the interest in such techniques.

Our method uses only shift-operations on small integer numbers. For each such shift-instruction one random bit decides which of two shift-operations to use for accumulation. Although promising for hardware-implementations, it remains to show if hardware architectures have a substantial benefit from our proposed method compared to floating point multiplications.

Alternative Network Design A topic strongly related to our work is stochastic computation (SC) [3; 18]. Here, numbers are approximated by sequences of binary samples whose average corresponds to the intended number. This idea has recently been applied to hardware implementations of deep networks [5; 33; 53]. Similarly, stochastic quantization (see above) has also been proven to be a valuable tool to reduce the negative impact of reduced representational efforts.

Closely related to SC is the idea of Sum-Product-Networks (SPN), a family of probabilistic graphical models that encode joint distributions of their input random variables [50]. SPNs are equivalent to arithmetic circuits, and thus are also promising candidates for computationally efficient networks.

¹ The sum of i.i.d random variables tends towards a normal distribution due to the central limit theorem.

² Ternary weights are usually constrained to the values $-1, 0$ or 1 .

³ The authors use 8-bit weights, 32-bit biases and fixed-point numbers for intermediate results.

They have been used successfully for typical deep learning tasks such as natural language processing [9], image classification [58] and image segmentation [52].

The biggest difference of our work to SC and SPNs is the interpretation of data-streams and intermediate results. SC and SPNs interpret data as random data streams. Our approach, in contrast, uses fixed-point numbers for incoming data and intermediate results, while only weights are random variables. This reduces the variance of intermediate results significantly (see Section 4). In contrast to the compelling implementations of SC, we prove the feasibility of our approach beyond MNIST. To the best of our knowledge, SPNs have not yet been shown to work feasibly on large scale image data, like ImageNet, however, the joint-distribution nature of SPNs and their clear semantics allow for new unseen types of queries [8; 28; 63; 65].

3 Capacitor Units

Current deep networks typically consist of linear combinations of activations (or inputs) followed by ReLU non-linearities. For simplicity, we consider a single activation in one specific layer of the network.

$$y = \text{relu} \left(\sum_{i=1}^d x_i \cdot w_i \right) \quad (1)$$

Here, x denotes the d -dimensional activation of the previous layer and w_i denotes the corresponding i -th weight. We omit the bias term for clarity (w.l.o.g., we can set $x_0 \equiv 1$). Convolutional layers have the same basic computational structure except they reuse weights and limit their support.

In addition, architectures add typically batch normalization to each activation after each layer. W.l.o.g. batch normalization is a fixed affine map

$$\text{bn}(y) = a \cdot y + b, \quad (2)$$

where a, b are constants during inference, and dependent on x during training. This mapping can be “folded” into the preceding or following linear layer [32]. I.e, it is possible to determine a map $w \mapsto w'$ such that batch normalization can be omitted.

Our method replaces multiplication by random shifts. We will discuss in Section 4.3 that it is crucial for our method to fold successive multiplications, as multiplications of random variables induce much higher variance.

3.1 Stochastic Binarization with Importance Sampling

As discussed in Section 1, the most expensive operations in terms of chip area in standard implementations are multiplication and floating-point logic. In order to reduce this cost, we combine previous ideas from the quantization, binarization and stochastic computing literature.

Stochastic multiplication First, we observe that a multiplication $w \cdot x$ of a real number $w \in [0, 1]$ and an integer x can be estimated stochastically by $B_w \cdot x$ where $B_w \in \{0, 1\}$ is a Bernoulli random variable with probability $\Pr(B_w = 1) = w$. By construction the mean equals to the original multiplication,

$$E[B_w \cdot x] = w \cdot x.$$

Thus, we substitute all multiplications $w \cdot x$ after folding,

$$w \mapsto B_w \quad (3)$$

and obtain a statistical approximation of the preactivation Equation (1), the argument of the relu-function.

For the d -dimensional scalar product $\sum_{i=1}^d x_i \cdot w_i$, this simple estimator stochastically ignores components of the activation vector x , eliminating multiplication at the cost of introducing (substantial) variance.

Value-based importance Sampling A drawback of networks consisting of only binary weights is the limited range of possible values. A typical observation is that the probabilities w gather around the borders 0 and 1 [16]. This reduces the amount of obtainable information of each gradient descent step, as gradients that point outside the weight-domain $[0, 1]$ get discarded (various techniques exist

which try to reduce this problem [16; 43; 66].

We propose an unbiased encoding that replaces floating point multiplications completely. Floating point representations consist of three parts: $s \cdot 2^e \cdot m$. The sign $s \in \{-1, 1\}$ and the exponent $e \in \mathbb{N}$ adjust the coarse magnitude, the mantissa m (typically $m \in [0, 1]$) fixes the more accurate decimal places. Our stochastic number representation that is well-suited for typical deep-learning computation, replaces the mantissa by a stochastic version that takes one of the values: 1 or 2. We even out stochasticity of multiplication with these numbers using accumulators and call the new stochastic multiplication *capacitor-unit*. This design utilizes only low-precision integer addition and does not require multiplication.

In more detail, we therefore split each weight into an exponent $e \in \mathbb{Z}$, a sign $s \in -1, 1$ and a mantissa $p \in [0, 1]$ and reformulate any weight w as

$$w \mapsto \bar{w} := s \cdot 2^e \cdot (B_p + 1), \quad (4)$$

where

$$s := \text{sign}(w), \quad (5)$$

$$e := \lfloor \log_2 |w| \rfloor \text{ and} \quad (6)$$

$$p := \frac{|w|}{2^e} - 1. \quad (7)$$

Technically, one would choose stochastically the bitshift $\cdot 2^{e+1}$ with probability p and $\cdot 2^e$ with probability $1 - p$ (See Figure 1 (a) and (b)). By construction, the mean of the representation equals to w ,

$$E[\bar{w}] = s \cdot 2^e \cdot \left(\frac{|w|}{2^e} - 1 + 1 \right) = s \cdot |w| = w.$$

Capacitors The binarization scheme from Equation (4) has a high variance (See Section 3.2).

In order to reduce the variance, we could draw multiple network samples and average the outcome of every sample. Deep networks are, however, non-linear;⁴ therefore, the statistical estimates obtained from plugging Bernoulli samples into Eq. 1 are not consistent for multi-layer networks that include non-linearities. Especially gradients w.r.t. the parameters of a probabilistic variable are hard to estimate [45]. The training process can be augmented to account for this [17], or to reduce the variance of the gradients [66], but even then, inference can suffer from high variance of intermediate results.⁵

Our solution is as follows: We run statistical averaging layer-wise, averaging over multiple samples *before* applying the non-linearity. We call this step a *capacitor*, in allusion to the analog component. Implementation costs are not impacted, assuming that computations are ordered accordingly and computed in parallel. This is the case for wide networks operating on, for example, lots of pixels, independent samples or activations.

Thus, we adapt Equation (4) to

$$w \mapsto \bar{w}_n := s \cdot 2^e \cdot \left(\frac{B_{n,p}}{n} + 1 \right), \quad (8)$$

where n denotes the sample size, and $B_{n,p}$ the n -fold binomial distribution with probability p .

In practice one would continue to sample from bernoulli distributions, for example by rewriting the multiplication to hardware-optimized operations,

$$x \cdot \bar{w}_{2^n} = s \cdot \left(\sum_{i=1}^{2^n} x \ll (e + B_p^{(i)}) \right) \gg n, \quad B_p^{(i)} \text{ i.i.d} \quad (9)$$

and where $\ll k$ denotes the k -bit left (multiplicational) shift and $\gg k$ denotes the k -bit right (divisional) shift.

We provide further notes on our implementation in the supplementary material.

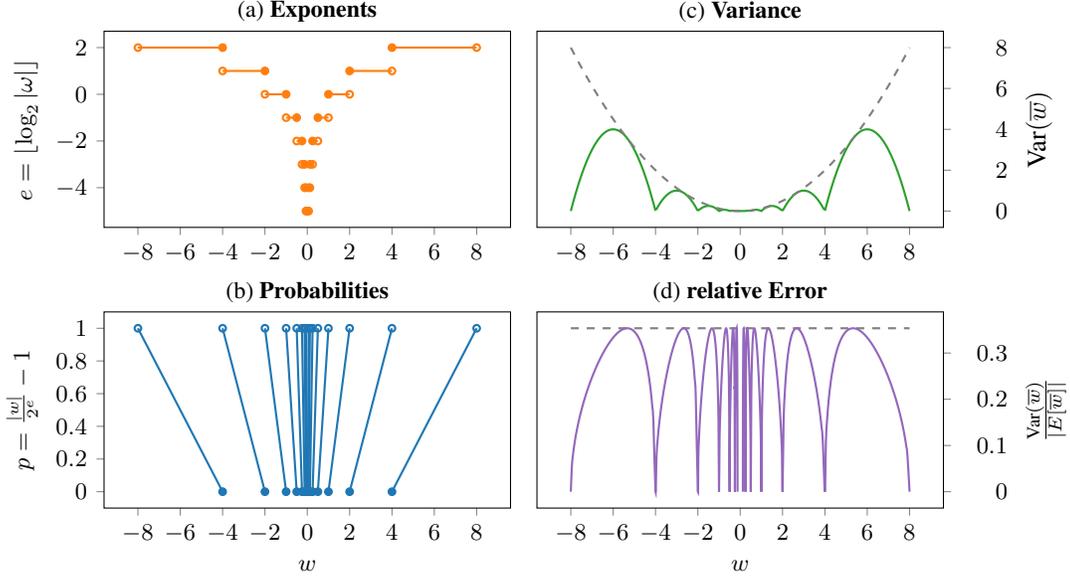


Figure 1: Figure (a) and (b) show of the values of the exponent of the components of the number system. Figure (b) and (c) show the variance and the relative error. In practice, the values near 0 are only hypothetical; too many shifts of integers always result in the number 0.

3.2 Properties

First, we consider the single-sampled case \bar{w} . Even though the mean of \bar{w} equals to w , the value varies significantly for large values of w (Figure 1 (c)),

$$\text{Var}(\bar{w}) \leq \frac{w^2}{8}$$

With multiple samples n , the bound decreases antiproportionally,

$$\text{Var}(\bar{w}_n) \leq \frac{w^2}{n \cdot 8}. \quad (10)$$

Sampling from this number system leads to high variances between two shifts, i.e. it is locally maximal whenever $p = 0.5$. For instance, the representation for $w = 3$ is $(e = 1, p = 0.5)$. Weights in these regions have high variances (see upper bound above). As our experiments show, however, this is not a problem. We will discuss this in more detail in Section 4.

It follows that the *relative variance* of the number system is constant (Figure 1 (d)):

$$\frac{\sigma_{\bar{w}_s}}{|E[\bar{w}_s]|} \leq \frac{1}{\sqrt{s \cdot 8}}. \quad (11)$$

4 Experimental Results

4.1 Setup

We implemented the method outlined above in TensorFlow. In all tests, we simulate our method using Equation (8), thus we sample the corresponding filter directly. We quantize all intermediate

⁴ This is an important difference to traditional statistical well-known integration approaches such as Monte-Carlo light-transport in computer graphics [11]. In this setting linearity allows for averaging of end-results.

⁵ This is intuitively clear: binary samples at the input of a highly nonlinear computational network lead to strongly varying outputs.

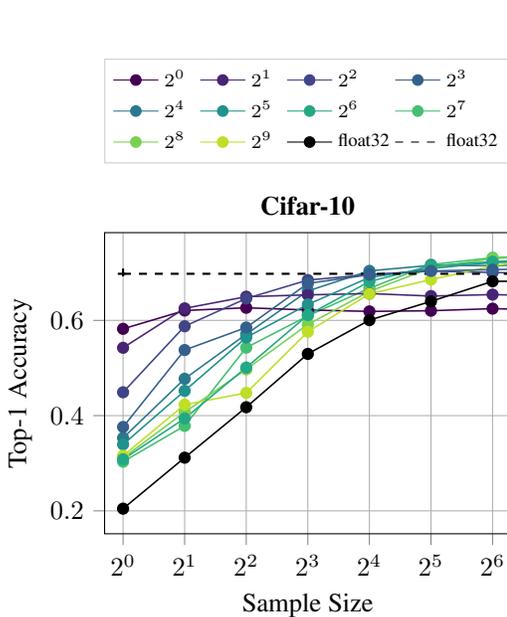


Figure 2: Training Cifar-10 using stochastic binarization with different sample sizes. After training, we use the network adaptively with other sample sizes. For reference, we trained the same network using *float32* (the black dashed line) and evaluated it afterwards again using stochastic binarization (the black solid line)

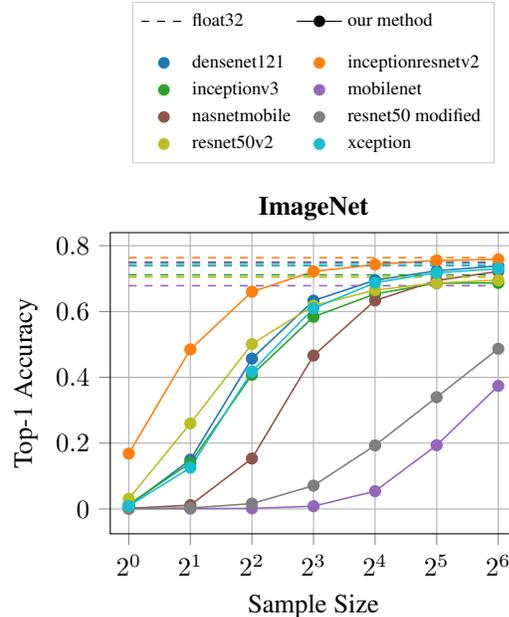


Figure 3: We evaluate pretrained ImageNet models after binarizing using stochastic binarization with different sample sizes. The dashed lines show the results for their unbinarized floating-point representation.

results to 16-bit integers ranging from -32 to 32. Absolute run-times are therefore not informative. Instead, we analyse sampling-accuracy trade-offs of our method.

Our experiments are as follows: First, we study the effect of different globally chosen sampling sizes during training on Cifar-10 (See Section 4.2), then we study our method as a binarization scheme on pretrained ImageNet (ILSVRC) models (Section 4.3). For the latter test case, we additionally examine the effect of probability discretization or naive pruning (Section 4.4). Lastly, we implement a naive attention-mechanisms that optimize the overall computational attention, i.e. choosing the sample size layerwise or spatially local (Section 4.5).

4.2 Cifar-10

The Cifar-10 dataset consists of 60,000 color images with 32×32 pixels and ten classes of natural objects. Each class consists of 5,000 training and 1,000 test examples [35].

We train a simple eight-layer convolutional network (a stack of eight 3×3 convolutions followed by a batch-normalization and a ReLU-activation) on the Cifar-10 dataset. This low-complexity test permits the study of significant changes of accuracy for this simple task. For the Cifar-10 experiments, we use a normal-distribution based initialization (LeCun initialization). We perform training using Adam optimizer with learning-rate $5e^{-3}$ with exponential decay after each 10 epochs by a factor of 0.1, weight decay of $5e^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1.0$ for 35 epochs.

In Figure 2, we compare the small convolutional model described above trained with floating point arithmetics to the same model incorporating capacitor units instead. First, we train a full precision network as a baseline (black dashed line in Figure 2). As our number representation applies to any weight, we use the pretrained full precision model and evaluate it using our number representation (black solid line) with varying number of samples. In contrast to that, we evaluate the effect of training directly on our progressive stochastic binarization, by training with sample-sizes ranging from 2^0

to 2^6 and then evaluating the pretrained weights with other sample sizes of stochastic binarization. Results show that when training using progressive stochastic binarization, the accuracy increased further in contrast to using pretrained models on full precision weights for stochastic quantization in test mode. This is in accordance with other quantization schemes.

4.3 ImageNet

Next we evaluate the binarization on several on ImageNet pretrained architectures in Figure 3. ImageNet is a large natural color image dataset, containing more than 14 million images and over 20,000 image categories. [54] We choose the commonly used ILSVRC-Subset. It consists of 1000 object classes and contains approximately 1.2 million training images, 50 thousand validation images and 100 thousand test images.

For evaluation, we use pretrained state-of-the-art architectures,⁶ mostly *Inception*- or *ResNet*-Types [24; 62]. This constraint is in favour of our implementation, as these architectures provide simple inner structures, namely easily foldable batch-normalizations after convolutional layers. For every model, we include the full-precision performance as a dashed line (see Figure 3).

In most cases, the binarization scheme already yields half of the accuracy of an unbinarized network with only four samples. With increasing sample size, the results approach floating point calculations. One of two exceptions is *MobileNet* [27]. *MobileNet* uses separable convolutions to reduce computational costs. But, in contrast to the other mobile-optimized networks *NasNetMobile* [71] and *Xception* [10] that incorporate separable convolutions as well, it separates depthwise from pixelwise convolutions by relu-nonlinearities, inducing a much higher error due to clipping stochastic results. This is in accordance with previous work that combines *MobileNet* and quantization [60]. Also, *NasNetMobile* and *Xception* accumulates different intermediate layers, reducing the stochastic errors further.

To emphasize the requirement of reducing direct multiplication of stochastic numbers (without convolutions of sufficient kernel size) in long operation chains, we added a modified *ResNet* model, *Resnet50 modified* with Batch-Normalizations after each shortcut.⁷ In this modification, data has to pass multiple Batch-Normalizations, thus, the data-stream that follows the shortcut undergoes multiple batch normalizations that we did not fold into the preceding layers. Instead, every batch normalizations acts as a separate multiplication of a stochastic number, additionally to the preceding stochastic convolution layers. In these shortcut-paths of the graph, the product of successive stochastic approximations increases the variance of the result directly.

4.4 Weight Pruning and Discretization of Probabilities

Next, we evaluate typical graph modifications to allow for even more computational or memory-wise efficiency on hardware. For these tests, we focus on a single representative network pretrained on ImageNet. We chose the *ResNet50* (v2) model as a commonly used representative for image-recognition tasks. For simplicity, we evaluate only float32, and, our discretization with 16 samples. Table 1

First, we apply a straight-forward magnitude-based threshold pruning method [22] to reduce 90% and 99% of all weights close to zero. In Table 1, we observe that pruning of 90% of weights does affect stochastic computation reasonably, however, for over-pruning (without retraining) of 99% of weights, disadvantages are more noticeable for stochastic computation.

Next, we reduce the memory footprint of the network. As stated in Section 3, using a high-precision probability-value of the number system, p , does not impact computational performance under the assumption that random numbers are generated efficiently. The reason for that is, that our method uses each probability value to generate one single bit to choose between one of two shifts. A drawback is, however, that high-precision probability-values do impact the memory-footprint of the network. Thus, we evaluate a strict quantization of probabilities in Table 1. We reduce the number of bits for each probability to 6, 4, 3, 2 bits and discrete quantization (i.e. 1 bits), where the quantization of probabilities is regular, including the boundary $p = 0$ and excluding the boundary $p = 1$ (the right boundary would result in a higher exponent). The results indicate that for the weight discretization, the accuracy does not change much in a stochastic setting, but drops significantly for the discrete case (1-bit). As we use 16-bits fixed point numbers for all intermediate results, we conclude that 4-bit

⁶The pretrained models are provided by https://github.com/qubvel/classification_models

⁷ The original paper calls this modification “BN after addition”.

exponents and 4-bit probabilities are sufficient for the use of our progressive stochastic binarization scheme on typical image recognition tasks.

It follows that our method can be reformulated to a deterministic version for tasks that require only a limited upper precision. For instance, for a 4-bit quantization of probabilities, we cannot gain any accuracy of the number system for more than 16 samples. Thus, instead of sampling the probability $p = 3/16$, we could use deterministically the smaller shift in 3 of 16 cases. This deterministic version, however, does not allow for a dynamic control of higher precision calculations.

Experiment	Number System	Accuracy Top-1 [%]
no modification	float32	70.43
	psb64	69.50
	psb32	68.56
	psb16	66.76
	psb8	61.86
pruning 90%	float32	70.34
	psb16	66.50
pruning 99%	float32	45.88
	psb16	39.10
6-bit probs	psb16	66.64
4-bit probs	psb16	66.80
3-bit probs	psb16	66.23
2-bit probs	psb16	50.29
1-bit probs	psb16	19.09
attention	psb8/16	65.74
	psb16/32	68.44
combined	psb8/16	66.15
	psb16/32	68.29

Table 1: Classification error (%) for the same (float32-)pretrained ResNet50 (v2) using the following modifications; plain inference, inference with pruning, inference with probability-discretization, inference using our attention mechanism and inference of a combination of the named techniques. We used two number systems: float32 and ours (psb16, 16-fold sampling) for inference.

4.5 Computational Attention

In this section, we take advantage of the adaptive sampling property of our method. As above, we use a ResNet50 (v2) model, pretrained on ImageNet. First, we discuss if computational effort benefits from a layer-wise adaption of sample accuracy. Then, we evaluate a simple spatial sampling adaption method.

One way of doing so is to adapt the sample size of each layer separately. In Figure 4, we analyse the approximation error, $|\frac{x_{psb2} - x_{float32}}{x_{float32}}|$, of the network using our method (psb2) compared to floating point calculations. We use our number system in a low-precision regime of only 2 samples per multiplication. For a given image (Figure 4a), we estimate the mean pixelwise approximation error using 100 samples of inference up to the first convolutional layer (Figure 4b) and the last convolutional layer (Figure 4c).

We evaluated if the distribution of errors changes with the layers (i.e. larger errors in lower layers or vice versa), however, we did not notice any such characteristic.

The approximation error seems to follow local features, at least in the last layer. In the following, we inspect the feasibility of spatial adaption of computation. Thus, we optimize the overall computational

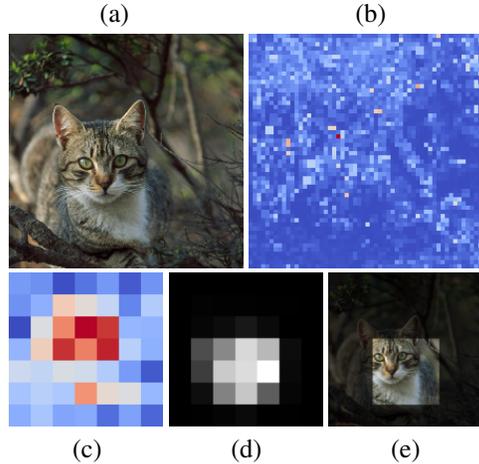


Figure 4: For a given image (a), we show the mean pixelwise approximation errors of 2 samples of our number system compared to float32 calculations in the first (b) and the last convolutional layer, L , (c), in the ResNet50 (v2) network.

We estimate the importance of the region by computing the pixelwise entropy in the layer L , (d), and its resulting mask, (e), for 8 samples. Blue or dark areas represent smallest values in the layer, red or light areas represent the biggest values in the layer.

attention by evaluating the network in a high-precision mode only on the regions of high activation in the image. We obtain those regions by using the same network in a low-precision mode on the full image (8 samples per computation). In more detail, we evaluate the pixelwise entropy of the image in the last convolutional layer (Figure 4d): We estimate the entropy h_{xy} by

$$h_{xy} := \sum_c -\text{softmax}(a_{xyc}) \cdot \log(\text{softmax}(a_{xyc})),$$

where a_{xyc} denotes the activation of the last layer in the pixel (x, y) and the channel c . In conclusion, regions of lowest pixelwise entropy have high magnitudes in only a few channels, regions of highest pixelwise entropy have equal activations in all channels. Thus, we have to refine high-entropy regions as those have higher probability of returning the wrong classification. For our attention mechanism, we use a hard threshold at the mean entropy in the image to estimate the interesting regions (fig. 4e). For the ImageNet testset this results in a ratio of about 35% of interesting regions of higher entropy and 65% of regions of lower entropy.

We use these two types of regions to separate the activations spatially (on all layers), to adapt computational attention of each filter individually (Table 1, attention modification). In the first experiment, we use 8 samples for regions outside computational interest and 16 samples on these regions (psb8/16). In total, we reduce computational effort by a factor of 33% compared to psb16. We expect this factor to increase significantly for datasets with higher image dimensions and thus, smaller regions of interest, like the MSCOCO-Dataset. The second experiment (psb16/32) increases computational effort by a factor of 33% compared to psb16, but results in performance close to psb32, a sampling mode of twice as many samples as psb16.

We reduce computational costs further by applying the previous techniques, probability discretization (to 4-bit probabilities) and pruning (90% threshold based pruning of weights) without losing significant amounts of precision (Table 1), resulting in a total of 68.29 % top-1 accuracy on the ImageNet test dataset.

5 Conclusion

We have introduced progressive stochastic binarization (psb), an unbiased stochastic binarization scheme for deep networks that replaces all floating-point multiplications. Our method converts network weights into stochastic shift operations. We even out stochasticity using layer-wise sampling. Progressive stochastic binarization therefore allows for efficient inference on hardware by restricting itself to additions of small integers and fixed shifts. The individual operations employed can be implemented on CMOS hardware at reduced costs (a full evaluation on actual FPGAs or ASICs, however, is still subject to future work), and the key operation of gating of continuous information might also be a useful primitive for unconventional non-CMOS computing processes such as neuro-morphic circuits. We applied the method for inference of pre-trained image classification networks, and also during training of newly initialized models. Performance wise we match the accuracy of previous binarized approaches in a low-precision setting, and in a high-precision regime, our method is accuracy-wise competitive with previous quantization schemes. The major limitation is that although our scheme applies successfully to a wide range of ResNet architectures, we found that it cannot be applied to specific types of separable convolutions. We examined the best practices for architectures by modifying a previously well-working ResNet architecture. Nonetheless, the method also permits localized, dynamic accuracy control within a single network, providing a new tool for adaptively focusing computational attention; only few quantization schemes provide this feature. We use that feature to direct computational costs adaptively using the network itself as an attention proposal mechanism for better classification results.

References

- [1] *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017. IEEE Computer Society. ISBN 978-1-5386-0457-1.
- [2] *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, 2018. IEEE Computer Society.
- [3] Alaghi, A. and Hayes, J. P. Fast and accurate computation using stochastic circuits. In Fettweis, G. P. and Nebel, W. (eds.), *Design, Automation & Test in Europe Conference & Exhibition*,

- DATE 2014, Dresden, Germany, March 24-28, 2014*, pp. 1–4. European Design and Automation Association, 2014. ISBN 978-3-9815370-2-4. doi: 10.7873/DATE.2014.089.
- [4] Alizadeh, M., Fernández-Marqués, J., Lane, N. D., and Gal, Y. An empirical study of binary neural networks’ optimisation. In *International Conference on Learning Representations*, 2019.
- [5] Ardakani, A., Leduc-Primeau, F., Onizawa, N., Hanyu, T., and Gross, W. J. VLSI implementation of deep neural network using integral stochastic computing. *IEEE Trans. VLSI Syst.*, 25 (10):2688–2699, 2017. doi: 10.1109/TVLSI.2017.2654298.
- [6] Bengio, Y., Léonard, N., and Courville, A. C. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013.
- [7] Bethge, J., Yang, H., Bartz, C., and Meinel, C. Learning to train a binary neural network. *CoRR*, abs/1809.10463, 2018.
- [8] Bueff, A., Speichert, S., and Belle, V. Tractable querying and learning in hybrid domains via sum-product networks. *CoRR*, abs/1807.05464, 2018.
- [9] Cheng, W., Kok, S., Pham, H. V., Chieu, H. L., and Chai, K. M. A. Language modeling with sum-product networks. In Li, H., Meng, H. M., Ma, B., Chng, E., and Xie, L. (eds.), *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014*, pp. 2098–2102. ISCA, 2014.
- [10] Chollet, F. Xception: Deep learning with depthwise separable convolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017 DBL [1]*, pp. 1800–1807. ISBN 978-1-5386-0457-1. doi: 10.1109/CVPR.2017.195.
- [11] Cook, R. L., Porter, T. K., and Carpenter, L. C. Distributed ray tracing. In Christiansen, H. (ed.), *Proceedings of the 11th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH 1984, Minneapolis, Minnesota, USA, July 23-27, 1984*, pp. 137–145. ACM, 1984. ISBN 0-89791-138-5. doi: 10.1145/800031.808590.
- [12] Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M., and Garnett, R. (eds.). *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, 2015.
- [13] Courbariaux, M., Bengio, Y., and David, J. Binaryconnect: Training deep neural networks with binary weights during propagations. In Cortes et al. [12], pp. 3123–3131.
- [14] Courbariaux, M., Bengio, Y., and David, J.-P. Training deep neural networks with low precision multiplications. *International Conference on Learning Representations Workshop*, 2015.
- [15] Dally, B. Efficient methods and hardware for deep learning. In *Proceedings of the Workshop on Trends in Machine-Learning (and impact on computer architecture), TIML@ISCA 2017, Toronto, ON, Canada, June 25, 2017*, pp. 2. ACM, 2017. ISBN 978-1-4503-5564-3. doi: 10.1145/3149166.3149168.
- [16] Darabi, S., Belbahri, M., Courbariaux, M., and Nia, V. P. BNN+: improved binary network training. *CoRR*, abs/1812.11800, 2018.
- [17] Fu, M. C. Gradient estimation. *Handbooks in operations research and management science*, 13: 575–616, 2006.
- [18] Gaines, B. R. *Stochastic computing systems*, pp. 37–172. Springer, 1969.
- [19] Gkoupidenis, P., Koutsouras, D. A., and Malliaras, G. G. Neuromorphic device architectures with global connectivity through electrolyte gating. *Nature communications*, 8:15448, 2017.
- [20] Gumbel, E. J. *Statistical theory of extreme values and some practical applications: a series of lectures*, volume 33. US Government Printing Office, 1954.

- [21] Gupta, S., Agrawal, A., Gopalakrishnan, K., and Narayanan, P. Deep learning with limited numerical precision. In Bach, F. R. and Blei, D. M. (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1737–1746. JMLR.org, 2015.
- [22] Han, S., Pool, J., Tran, J., and Dally, W. J. Learning both weights and connections for efficient neural network. In Cortes et al. [12], pp. 1135–1143.
- [23] Han, S., Mao, H., and Dally, W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *International Conference on Learning Representations (ICLR)*, 2016.
- [24] He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In Leibe et al. [40], pp. 630–645. ISBN 978-3-319-46492-3. doi: 10.1007/978-3-319-46493-0_38.
- [25] Hill, P., Zamirai, B., Lu, S., Chao, Y., Laurenzano, M., Samadi, M., Papaefthymiou, M. C., Mahlke, S. A., Wenisch, T. F., Deng, J., Tang, L., and Mars, J. Rethinking numerical representations for deep neural networks. *CoRR*, abs/1808.02513, 2018.
- [26] Horowitz, M. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, pp. 10–14, Feb 2014. doi: 10.1109/ISSCC.2014.6757323.
- [27] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017.
- [28] Hsu, W., Kalra, A., and Poupart, P. Online structure learning for sum-product networks with gaussian leaves. *CoRR*, abs/1701.05265, 2017.
- [29] 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 2017, Honolulu, HI, USA, July 21-26, 2017* DBL [1], pp. 2261–2269. ISBN 978-1-5386-0457-1. doi: 10.1109/CVPR.2017.243.
- [30] Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., and Bengio, Y. Binarized neural networks. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 4107–4115, 2016.
- [31] Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A. G., Adam, H., and Kalenichenko, D. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018* DBL [2], pp. 2704–2713. doi: 10.1109/CVPR.2018.00286.
- [32] Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A. G., Adam, H., and Kalenichenko, D. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018* DBL [2], pp. 2704–2713. doi: 10.1109/CVPR.2018.00286.
- [33] Kim, K., Kim, J., Yu, J., Seo, J., Lee, J., and Choi, K. Dynamic energy-accuracy trade-off using stochastic computing in deep neural networks. In *Proceedings of the 53rd Annual Design Automation Conference, DAC 2016, Austin, TX, USA, June 5-9, 2016*, pp. 124:1–124:6. ACM, 2016. ISBN 978-1-4503-4236-0. doi: 10.1145/2897937.2898011.
- [34] Kingma, D. P., Salimans, T., and Welling, M. Variational dropout and the local reparameterization trick. In Cortes et al. [12], pp. 2575–2583.
- [35] Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [36] Kung, H. T. Why systolic architectures? *IEEE Computer*, 15(1):37–46, 1982. doi: 10.1109/MC.1982.1653825.

- [37] Lacey, G., Taylor, G. W., and Areibi, S. Deep learning on fpgas: Past, present, and future. *CoRR*, abs/1602.04283, 2016.
- [38] LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *Nature*, pp. 436–444.
- [39] Lee, V. T., Alaghi, A., Hayes, J. P., Sathe, V., and Ceze, L. Energy-efficient hybrid stochastic-binary neural networks for near-sensor computing. In Atienza, D. and Natale, G. D. (eds.), *Design, Automation & Test in Europe Conference & Exhibition, DATE 2017, Lausanne, Switzerland, March 27-31, 2017*, pp. 13–18. IEEE, 2017. ISBN 978-3-9815370-8-6. doi: 10.23919/DATE.2017.7926951.
- [40] Leibe, B., Matas, J., Sebe, N., and Welling, M. (eds.). *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, volume 9908 of *Lecture Notes in Computer Science*, 2016. Springer. ISBN 978-3-319-46492-3. doi: 10.1007/978-3-319-46493-0.
- [41] Li, Y., Liu, Z., Xu, K., Yu, H., and Ren, F. A gpu-outperforming FPGA accelerator architecture for binary convolutional neural networks. *JETC*, 14(2):18:1–18:16, 2018. doi: 10.1145/3154839.
- [42] Lin, D. D., Talathi, S. S., and Annapureddy, V. S. Fixed point quantization of deep convolutional networks. In Balcan, M. and Weinberger, K. Q. (eds.), *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 2849–2858. JMLR.org, 2016.
- [43] Lin, X., Zhao, C., and Pan, W. Towards accurate binary convolutional neural network. In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pp. 344–352, 2017.
- [44] Liu, Z., Wu, B., Luo, W., Yang, X., Liu, W., and Cheng, K. Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm. In Ferrari, V., Hebert, M., Sminchisescu, C., and Weiss, Y. (eds.), *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XV*, volume 11219 of *Lecture Notes in Computer Science*, pp. 747–763. Springer, 2018. ISBN 978-3-030-01266-3. doi: 10.1007/978-3-030-01267-0_44.
- [45] Maddison, C. J., Mnih, A., and Teh, Y. W. The concrete distribution: A continuous relaxation of discrete random variables. *CoRR*, abs/1611.00712, 2016.
- [46] Marsaglia, G. Random number generators. *Journal of Modern Applied Statistical Methods*, 2(1):2, 2003.
- [47] Meng, W., Gu, Z., Zhang, M., and Wu, Z. Two-bit networks for deep learning on resource-constrained embedded devices. *CoRR*, abs/1701.00485, 2017.
- [48] Mittal, S. A survey of fpga-based accelerators for convolutional neural networks. *Neural computing and applications*, pp. 1–31, 2018.
- [49] Peters, J. W. T. and Welling, M. Probabilistic binary neural networks. *CoRR*, abs/1809.03368, 2018.
- [50] Poon, H. and Domingos, P. M. Sum-product networks: A new deep architecture. In *IEEE International Conference on Computer Vision Workshops, ICCV 2011 Workshops, Barcelona, Spain, November 6-13, 2011*, pp. 689–690. IEEE Computer Society, 2011. ISBN 978-1-4673-0062-9. doi: 10.1109/ICCVW.2011.6130310.
- [51] Rastegari, M., Ordonez, V., Redmon, J., and Farhadi, A. Xnor-net: Imagenet classification using binary convolutional neural networks. In Leibe et al. [40], pp. 525–542. ISBN 978-3-319-46492-3. doi: 10.1007/978-3-319-46493-0_32.

- [52] Rathke, F., Desana, M., and Schnörr, C. Locally adaptive probabilistic models for global segmentation of pathological OCT scans. In Descoteaux, M., Maier-Hein, L., Franz, A. M., Jannin, P., Collins, D. L., and Duchesne, S. (eds.), *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017 - 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I*, volume 10433 of *Lecture Notes in Computer Science*, pp. 177–184. Springer, 2017. ISBN 978-3-319-66181-0. doi: 10.1007/978-3-319-66182-7_21.
- [53] Ren, A., Li, Z., Ding, C., Qiu, Q., Wang, Y., Li, J., Qian, X., and Yuan, B. SC-DCNN: highly-scalable deep convolutional neural network using stochastic computing. In Chen, Y., Temam, O., and Carter, J. (eds.), *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2017, Xi'an, China, April 8-12, 2017*, pp. 405–418. ACM, 2017. ISBN 978-1-4503-4465-4. doi: 10.1145/3037697.3037746.
- [54] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M. S., Berg, A. C., and Li, F. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [55] Sandler, M., Howard, A. G., Zhu, M., Zhmoginov, A., and Chen, L. Mobilenetv2: Inverted residuals and linear bottlenecks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018 DBL [2]*, pp. 4510–4520. doi: 10.1109/CVPR.2018.00474.
- [56] Schmidhuber, J. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015. doi: 10.1016/j.neunet.2014.09.003.
- [57] Schuman, C. D., Potok, T. E., Patton, R. M., Birdwell, J. D., Dean, M. E., Rose, G. S., and Plank, J. S. A survey of neuromorphic computing and neural networks in hardware. *CoRR*, abs/1705.06963, 2017.
- [58] Sguerra, B. M. and Cozman, F. G. Image classification using sum-product networks for autonomous flight of micro aerial vehicles. In *5th Brazilian Conference on Intelligent Systems, BRACIS 2016, Recife, Brazil, October 9-12, 2016*, pp. 139–144. IEEE Computer Society, 2016. ISBN 978-1-5090-3566-3. doi: 10.1109/BRACIS.2016.035.
- [59] Shayer, O., Levi, D., and Fetaya, E. Learning discrete weights using the local reparameterization trick. *CoRR*, abs/1710.07739, 2017.
- [60] Sheng, T., Feng, C., Zhuo, S., Zhang, X., Shen, L., and Aleksic, M. A quantization-friendly separable convolution for mobilenets. *CoRR*, abs/1803.08607, 2018.
- [61] Sze, V., Chen, Y., Yang, T., and Emer, J. S. Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12):2295–2329, 2017. doi: 10.1109/JPROC.2017.2761740.
- [62] Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. A. Inception-v4, inception-resnet and the impact of residual connections on learning. In Singh, S. P. and Markovitch, S. (eds.), *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pp. 4278–4284. AAAI Press, 2017.
- [63] van de Wolfshaar, J. and Pronobis, A. Deep convolutional sum-product networks for probabilistic image representations. *CoRR*, abs/1902.06155, 2019.
- [64] Vanhoucke, V., Senior, A., and Mao, M. Z. Improving the speed of neural networks on cpus. In *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop*, volume 1, pp. 4. Citeseer, 2011.
- [65] Vergari, A., Peharz, R., Mauro, N. D., Molina, A., Kersting, K., and Esposito, F. Sum-product autoencoding: Encoding and decoding representations using sum-product networks. In McIlraith, S. A. and Weinberger, K. Q. (eds.), *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial*

Intelligence (IAAI-18), and the 8th AAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, pp. 4163–4170. AAAI Press, 2018.

- [66] Yin, M. and Zhou, M. ARM: augment-reinforce-merge gradient for discrete latent variable models. *CoRR*, abs/1807.11143, 2018.
- [67] Zhou, A., Yao, A., Guo, Y., Xu, L., and Chen, Y. Incremental network quantization: Towards lossless cnns with low-precision weights. *CoRR*, abs/1702.03044, 2017.
- [68] Zhou, S., Ni, Z., Zhou, X., Wen, H., Wu, Y., and Zou, Y. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *CoRR*, abs/1606.06160, 2016.
- [69] Zhu, C., Han, S., Mao, H., and Dally, W. J. Trained ternary quantization. *CoRR*, abs/1612.01064, 2016.
- [70] Zhuang, B., Shen, C., Tan, M., Liu, L., and Reid, I. Rethinking binary neural network for accurate image classification and semantic segmentation. *CoRR*, abs/1811.10413, 2018.
- [71] Zoph, B., Vasudevan, V., Shlens, J., and Le, Q. V. Learning transferable architectures for scalable image recognition. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018* DBL [2], pp. 8697–8710. doi: 10.1109/CVPR.2018.00907.

Supplementary Material for: Progressive Stochastic Binarization of Deep Networks

David Hartmann
Institute of Computer Science
Johannes Gutenberg-University of Mainz
Staudingerweg 9, 55128 Mainz, Germany
dahartma@uni-mainz.de

Michael Wand
Institute of Computer Science
Johannes Gutenberg-University of Mainz
Staudingerweg 9, 55128 Mainz, Germany
mwand@uni-mainz.de

1 Notes on our Implementation

Forward pass: We obtain the forward pass by substituting the matrix products that linearly combine activations using the weights by capacitor units, as described in Section 3.1 of the paper.

Our implementation simulates the method, thus it performs all computations using `float32` precision. In order to account for quantization, we quantize to 16-bit fixed-point numbers, ranging from -32 to 32 .⁸ On all of our experiments, 16-bits were enough, however, using less bits was sufficient for some architectures as the values of the activations did not reach the boundaries ± 32 .

Backward pass: Our method supports an optional training step to improve performance for networks that have not been adapted to stochastic quantization noise. We currently perform training on standard hardware (GPUs or CPU). In the backward pass we just optimize continuous weights and convert them into the exponent-probability representation after each step of gradient descent (regardless of the specific method such as Adam, Momentum or SGD). We compute gradients without any modification; as we introduce a new number representation we calculate the gradients as if no modification was made to the weights.

Efficient simulation of the forward pass: A technical problem is that GPU hardware does not benefit from our optimizations; to the contrary, despite already running `float32`-computations, additional costs incur for quantization and format conversion. In addition, if not parallelized, repeated sampling slows down the process linearly, without achieving any speed benefits for each individual pass. We therefore augment the forward pass to directly sample from Binomial distributions. Rather than accumulating n Bernoulli experiments with probability p , we sample from

$$B_p(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad (12)$$

using the Gumbel-Max-Trick [20];

$$B_p(i) = \arg \max_k \left[\log \left(\binom{n}{k} p^k (1-p)^{n-k} \right) - \log(-\log U_k), \right] \quad (13)$$

$$(14)$$

where U_k denotes i.i.d. uniform variables $U_k \sim U(0, 1)$. We reformulate this to be numerically stable by using log-rules,

$$B_p(i) = \arg \max_k \left[\log \binom{n}{k} + k \cdot \log(p) + (n-k) \cdot \log(1-p) - \log(-\log U_k), \right]. \quad (15)$$

operation	chip area [μm^2]	chip area, relative to fp32 mul	energy [pJ]
int8 add	36	0.005	0.03
int16 add	67	0.01	0.06
int32 add	137	0.02	0.10
int8 mul	282	0.04	0.20
int32 mul	3,495	0.45	1.10
fp16 add	1,360	0.18	0.40
fp16 mul	1,640	0.21	1.10
fp32 add	4,184	0.54	0.90
fp32 mul	7,700	1	3.70

Table 2: Hardware costs for common arithmetic units (45nm-process, [15; 26]).

1.1 Hardware Implementation Perspectives

As the computational structure of deep network training and inference is highly data-parallel, a reduction in chip area and power consumption per operation directly translates to substantial performance advantages (see Table 2).

The representation outlined in the main part of the paper is motivated by gaining potential for a more efficient implementation on custom CMOS hardware.

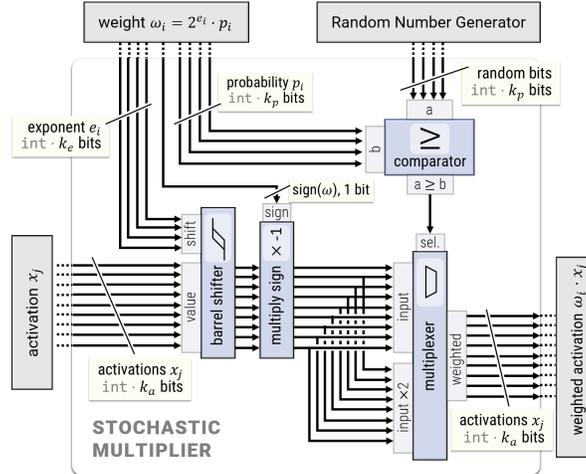


Figure 5: Schematic circuit diagram of a stochastic multiplier.

It performs multiplication by randomly selecting between two shifted versions of the activation; all shifts are by constants. One single random bit with probability p chooses which of those two shifts (e or $e + 1$) is used. The accuracy of intermediate integer representation restricts the magnitude of the exponents, the precision of the probability values is restricted by the memory requirements and accuracy requirements of the network.

Stochastic multiplier: In Section 4, we quantized probabilities using k_p integer numbers. Sampling of random bits with (quantized) probability $p \in [0, 1]$ requires a k_p -bit comparator, which corresponds to an accordingly sized integer subtraction unit. This subunit returns one single random bit. As no other operation depends on the concrete value of p , the number of bits k_p is only bound by memory restrictions. In Section 4 of the main paper, we have discussed the implications on performance if reducing this number dramatically.

Handling exponents (stored as k_e -bit integers) for every weight requires a barrel-shifter (i.e., the shift is selected at run-time by turning on or off shifts by powers of two by control-bits); the additional shift by two is obtained by simple wiring. Finally, a multiplexing gate, controlled by the randomly sampled bit, determines the outcome. In Section 4, we observe that 4-bit exponents and 4-bit probabilities are sufficient in all our experiments.

⁸ We provide a tensorflow-tool that first folds all batchnorm-convolution pairs, then converts all remaining convolutions and batch-normalization layers to our stochastic number representation, and quantizes all bias terms and intermediate results to fixed-point numbers.

Activations and Accumulation: We store all activations and perform all accumulations in integer format (fixed point). As addition units are cheap compared to multiplications, this appears to be a better trade-off than reintroducing adaptive shifts of a full floating-point scheme. Our Experiments show that a 16-bit integer representation is sufficient to avoid significant accuracy penalties. Figure 5 shows a rough data-flow diagram of a full capacitor unit.

ReLU, Batch-Norm, and Pooling: ReLU units are easy to realize in hardware (using a gate that depends on the sign bit). Batch-norm requires a more sophisticated approach; during training, one could quantize factors to powers of two. When using pretrained networks, this is not possible (in particular when short-cut connections void scaling invariance). Thus, the better alternative is to “fold” the mapping into the linear weights as stated in Section 3. Note that we can fold the batch-norm into the weights before the calculation of the encoding, i.e. we transform e and p accordingly. Similarly, non-power-of-two scaling factors in average-pooling layers can be folded for pretrained networks as well, if needed, or implemented as multiplication with a global stochastic number (both variations did work in our experiments). Max-pooling needs additional comparator units and gating.

Random number generation: Stochastic computation requires random numbers that are sufficiently independent and evenly distributed. Simple linear feedback shift registers are sufficient for this purpose; however, they add some overhead (numbers are generated by recurrent fixed shifts, additions, and XOR with constants) [46]. The current implementation of tensorflow uses *XORWOW* as the standard random number generator on graphics cards and *MT19937_64* for CPU-operations. We tested both and did not recognize any differences in the resulting performance.

Classification Layer: For inference, we can ignore the softmax-layer, as we are only interested in the maximum value of the classification. A possible average pooling in the last layer can be ignored w.l.o.g., as the normalization can be trivially folded to the preceding layer. Training, however, may need a softmax-layer depending on the loss function. We did not cover this in the scope of this work; for simplicity, we did not replace these layers in our implementation.

Systolic Arrays: Modern Computational Units like GPUs and TPUs include systolic arrays for faster matrix-matrix multiplication [36]. It appears plausible that our method can be included in various ways directly into systolic arrays, as the accumulation step of independent samples and matrix-multiplication are interchangeable due to linearity.

Thus, one could directly replace the innermost floating point multiplication of systolic arrays with a capacitor unit. For more flexibility sampling-wise, however, it is beneficial to replace the multiplication with a single stochastic multiplier without accumulation (Figure 5). For multiple samples, one could either sum any number of results of the single-sampled systolic array, or concatenate any number of same weight tensors as the input.

Discussion: The proposal leaves out many details required for an actual hardware design, such as allocating the right ratio of accumulation and multiplication units, organizing the accumulation loop, pipelining, and memory management (storage for intermediate results, avoiding energy costs for non-local memory, and many more). Our paper focuses on the main representation and its accuracy properties, in particular in practical experiments. Weighting the costs against alternative quantization schemes might lead to trade-offs that can only be assessed in a full hardware implementation. The idea of stochastic accumulation is, however, of conceptual interest beyond digital CMOS circuits. Neuromorphic computing with (partially) analog elements could similarly benefit from it, and accumulation could be performed by capacitors (electrical or otherwise), and sources of random noise might be available without digital computation.