

---

# Multi-Precision Policy Enforced Training (MuPPET): A precision-switching strategy for quantised fixed-point training of CNNs

---

Aditya Rajagopal<sup>\*1</sup> Diederik Adriaan Vink<sup>\*1</sup> Stylianos I. Venieris<sup>2</sup> Christos-Savvas Bouganis<sup>1</sup>

## Abstract

Large-scale convolutional neural networks (CNNs) suffer from very long training times, spanning from hours to weeks, limiting the productivity and experimentation of deep learning practitioners. As networks grow in size and complexity, training time can be reduced through low-precision data representations and computations, however, in doing so the final accuracy suffers due to the problem of vanishing gradients. Existing state-of-the-art methods combat this issue by means of a mixed-precision approach utilising two different precision levels, FP32 (32-bit floating-point) and FP16/FP8 (16-/8-bit floating-point), leveraging the hardware support of recent GPU architectures for FP16 operations to obtain performance gains. This work pushes the boundary of quantised training by employing a multilevel optimisation approach that utilises multiple precisions including low-precision fixed-point representations resulting in a novel training strategy MuPPET; it combines the use of multiple number representation regimes together with a precision-switching mechanism that decides at run time the transition point between precision regimes. Overall, the proposed strategy tailors the training process to the hardware-level capabilities of the target hardware architecture and yields improvements in training time and energy efficiency compared to state-of-the-art approaches. Applying MuPPET on the training of AlexNet, ResNet18 and GoogLeNet on ImageNet (ILSVRC12) and targeting an NVIDIA Turing GPU, MuPPET achieves the same accuracy as standard full-precision training with training-time speedup of up to  $1.84\times$  and an average speedup of  $1.58\times$  across the networks.

---

<sup>\*</sup>Equal contribution <sup>1</sup>Intelligent Digital Systems Lab, Department of Electrical Engineering, Imperial College London <sup>2</sup>Samsung AI Center, Cambridge. Correspondence to: Aditya Rajagopal <ar4414@ic.ac.uk>.

## 1. Introduction

Convolutional neural networks (CNNs) have demonstrated unprecedented accuracy in various machine learning tasks, from video understanding (Gan et al., 2015; He et al., 2018) to drone navigation (Loquercio et al., 2018; Kouris & Bouganis, 2018). To achieve such high levels of accuracy in inherently complex applications, current methodologies employ the design of large and complex CNN models (Szegedy et al., 2017; Huang et al., 2017) trained over large datasets (Deng et al., 2009; Lin et al., 2014). Despite the associated accuracy gains, the combination of large models with massive datasets also results in excessively long training times. This in turn leads to long turn-around times which set a limit to the productivity of deep learning practitioners and prohibit wider experimentation. For instance, automatic tuning and search of neural architectures (Cai et al., 2018; Zhong et al., 2018) is a rapidly advancing area where accelerated training enables the improvement of the produced networks.

To counteract these long turn-around times, substantial research effort has been invested in hyperparameter tuning for the acceleration of training, with a particular focus on batch size and content. One line of work aims to maximise memory and hardware utilisation by changing the batch size, in order to perform CNN training-specific prefetching, scheduling and dependency improvement (Chen et al., 2019; Rhu et al., 2016). Other works focus on altering batch size or reconstructing minibatches to improve the convergence rate while sustaining high hardware utilisation (Devarakonda et al., 2017; Johnson & Guestrin, 2018; Peng et al., 2019), demonstrating up to  $6.25\times$  training speedup.

In a different direction, a number of studies have focused on the use of reduced-precision training schemes. Reduced-precision arithmetic involves the utilisation of data formats that have smaller wordlengths than the conventional 32-bit floating-point (FP32) representation and is an approach for co-optimising processing speed, memory footprint and communication overhead. Existing literature can be categorised into i) works that use reduced precision to only accelerate the training stage while eventually yielding an FP32 model, and ii) those that produce networks with quantised weights. Regarding the former, Courbariaux et al. (2015) and Gupta et al. (2015) utilise dynamic quantisation and stochastic rounding respectively as a means to combat the accuracy

loss due to quantisation. Nevertheless, the effectiveness of the proposed schemes has only been demonstrated on small-scale datasets such as CIFAR-10 and MNIST, and on a limited set of networks. Furthermore, the range of quantisation levels that has been explored varies greatly, with a number of works focusing solely on mild quantisation levels such as half-precision floating-point (FP16) (Micikevicius et al., 2018), while others focus on lower precisions such as 8-bit floating-point (FP8) (Wang et al., 2018). Finally, quantisation has also been used as a means of reducing the memory and communication overhead in distributed training (De Sa et al., 2015; 2017; Alistarh et al., 2017).

At the same time, the characteristics of modern CNN workloads and the trend towards quantised models have led to an emergence of specialised hardware processors, with support for low-precision arithmetic at the hardware level. From custom designs such as Google’s TPUs (Jouppi et al., 2017) and Microsoft’s FPGA-based Brainwave system (Fowers et al., 2018) to commodity devices such as NVIDIA’s Turing GPUs (Burgess, 2019), existing platforms offer increased parallelism through native support for reduced-precision data types including 16-bit floating-point (FP16), as well as 8- (INT8) and 4-bit (INT4) fixed-point. Although these platforms have been mainly designed for the inference stage, the low-precision hardware offers significant opportunities for accelerating the time-consuming training stage. Consequently, there is an emerging need to provide training algorithms that can leverage these existing hardware optimisations and yield higher training speed.

This work tackles the field of reduced-precision training at an algorithmic level. Independently of the number of quantisation levels chosen, or how extreme the quantisation is, this work proposes a metric that estimates the amount of information each new training step obtains for a given quantisation level, by capturing the diversity of the computed gradients *across epochs*. This enables the design of a policy that, given a set of quantisation levels, decides *at run time* appropriate points to increase the precision of the training process at that current instant without impacting the achieved test accuracy compared to training in FP32. Due to its agnostic nature, it remains orthogonal and complementary to existing low-precision training schemes. Furthermore, by pushing the precision below the 16-bit bitwidth of existing state-of-the-art techniques, the proposed method is able to leverage the low-precision capabilities of modern processing systems to yield training speedups without penalising the resulting accuracy, significantly improving the time-to-accuracy trade-off.

## 2. Background and Related Work

The state-of-the-art method in training with reduced precision is mixed-precision training (Micikevicius et al., 2018).

The authors propose to employ low-precision FP16 computations in the training stage of high-precision CNNs that perform inference in FP32. Along the training phase, the algorithm maintains a high-precision FP32 copy of the network’s weights, known as a *master copy*. At each minibatch, the inputs and weights are quantised to FP16 with all computations of the forward and backward pass performed in FP16, yielding memory footprint and runtime savings. Under this scheme, each stochastic gradient descent (SGD) update step entails accumulating FP16 gradients into the FP32 master copy of the weights, with this process performed iteratively throughout the training of the network. Micikevicius et al. (2018) evaluate their scheme over a set of state-of-the-art models on ImageNet, and show that mixed-precision training with FP16 computations achieves comparable accuracy to standard FP32 training.

Wang et al. (2018) also presented a method to train an FP32 model using 8-bit floating-point (FP8). The authors propose a hand-crafted FP8 data type, together with a chunk-based computation technique, and employ strategies such as stochastic rounding to alleviate the accuracy loss due to training at reduced precision. For AlexNet, ResNet18 and ResNet50 on ImageNet, Wang et al. (2018) demonstrate comparable accuracy to FP32 training while performing computations in FP8.

Additionally, the works presented in (Zhou et al., 2016; Chen et al., 2017) approach the problem of reduced-precision training through fixed-point computations. FxpNet (Chen et al., 2017) was only evaluated on CIFAR-10, failing to demonstrate performance on more complex datasets such as ImageNet. DoReFa-Net (Zhou et al., 2016) was tested on ImageNet, but only ran on AlexNet missing out on state-of-the-art networks such as GoogLeNet and ResNet.

All related works focus on either accelerating the training of an FP32 model through reduced-precision floating-point computations or using reduced-precision fixed-point computations to train a fixed-point model. At the hardware level, 8-bit fixed-point multiplication uses  $18.5\times$  less energy and  $27.5\times$  less area with up to  $4\times$  lower runtimes than FP32 (Sze et al., 2017; Kouris et al., 2018). Consequently, this work attempts to push the boundaries of reduced-precision training by combining the processing performance gains of low-bitwidth fixed-point computations with the floating-point-level accuracy of training an FP32 model.

Preliminary tests (Sec. 4.4 for details) demonstrated that training solely in 8-bit fixed-point results in a significant degradation of validation accuracy compared to full FP32 training. This work aims to counteract this degradation by *progressively* increasing the precision of computations throughout training in an online manner determined by the proposed metric inspired by gradient diversity (Yin et al.,

2018). Additionally by operating in an online fashion, MuPPET tailors the training process to best suit the particular network-dataset pair at each stage of the training process.

Gradient diversity was introduced by Yin et al. (2018) as a metric of measuring the dissimilarity between sets of gradients that correspond to different minibatches. The gradient diversity of a set of gradients is defined as

$$\begin{aligned} \Delta_S(\mathbf{w}) &= \frac{\sum_{i=1}^n \|\nabla f_i(\mathbf{w})\|_2^2}{\|\sum_{i=1}^n \nabla f_i(\mathbf{w})\|_2^2} \\ &= \frac{\sum_{i=1}^n \|\nabla f_i(\mathbf{w})\|_2^2}{\sum_{i=1}^n \|\nabla f_i(\mathbf{w})\|_2^2 + \sum_{i \neq j} \langle \nabla f_i(\mathbf{w}), \nabla f_j(\mathbf{w}) \rangle} \end{aligned} \quad (1)$$

where  $\nabla f_i(\mathbf{w})$  is the gradient of weights  $\mathbf{w}$  for minibatch  $i$ .

The key point to note in Eq. (1) is that the denominator contains the inner product between two gradients from *different minibatches*. Thus, orthogonal gradients would result in high gradient diversity, while similar gradients would yield low gradient diversity. The proposed framework, MuPPET, enhances this concept by considering gradients between minibatches *across epochs* and proposes the developed metric as a proxy for the amount of new information gained in each training step. Section 3 further expands on how gradient diversity is incorporated into the MuPPET algorithm.

### 3. Methodology - MuPPET Algorithm

#### 3.1. Multilevel Optimisation for Training CNNs

Conventionally, the training process of a CNN can be expressed as in Eq. (2). Given a CNN model  $f$  parametrised by a set of weights  $\mathbf{w} \in \mathbb{R}^D$ , where  $D$  is the number of weights of  $f$ , training involves a search for weight values that minimise the task-specific empirical loss,  $Loss$ , on the target dataset. Typically, a fixed arithmetic precision is employed across the training algorithm with FP32 currently being the *de facto* representation used by the deep learning community.

$$\min_{\mathbf{w}^{(FP32)} \in \mathbb{R}^D} Loss(f(\mathbf{w}^{(FP32)})) \quad (2)$$

The proposed method follows a different approach by introducing a multilevel optimisation scheme (Migdalas et al., 2013) that leverages the performance gains of reduced-precision arithmetic. The single optimisation problem of Eq. (2) is transformed into a series of optimisation problems with each one employing different precision for computations, but maintaining weights storage at FP32 precision. Under this scheme, an  $N$ -level formulation comprises  $N$  sequential optimisation problems to be solved, with each level corresponding to a “finer” model.

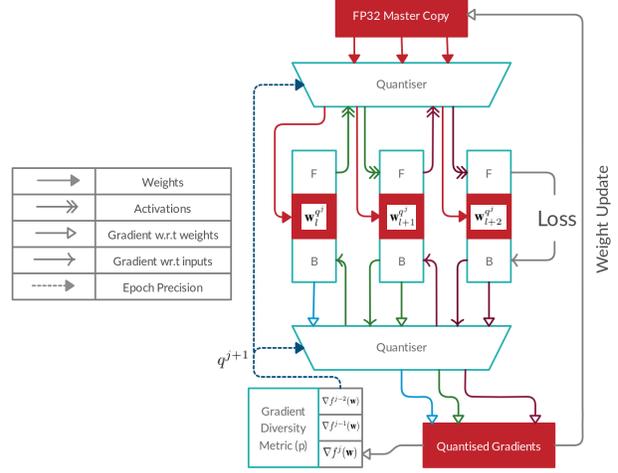


Figure 1. Quantised training scheme of MuPPET.

Overall, this formulation adds a hierarchical structure to the training stage, with increasing arithmetic precision across the hierarchy of optimisation problems. Starting from the  $N$ -th problem, the inputs, weights and activations of the CNN model  $f$  are quantised with precision  $q^N$ , which is the lowest precision in the system and represents the coarsest version of the model. Each of the  $N$  levels progressively employs higher precision until the first level is reached, which corresponds to the original problem of Eq. (2). Formally, at the  $i$ -th level, the optimisation problem is formulated as

$$\min_{\mathbf{w}^{(q^i)} \in \mathcal{V}} Loss(f(\mathbf{w}^{(q^i)})) \quad \text{s.t. } \mathcal{V} = \left\{ \mathbf{w}^{(q^i)} \in [\text{LB}, \text{UB}]^D \right\} \quad (3)$$

where LB and UB are the lower and upper bound in the representational range of precision  $q^i$ . The target CNN model  $f$  uses a set of weights quantised with precision  $q^i$  and hence the solution of this optimisation problem can be interpreted as an approximation to the original problem of Eq. (2). To transition from one level to the next, the result of each level of optimisation is employed as a starting point for the next level, up to the final outermost optimisation that reduces to Eq. (2).

#### 3.2. Quantised Training

Fig. 1 presents the process of training a CNN using the proposed algorithm. At epoch  $j$ , MuPPET performs mixed-precision training where the weights are stored in an FP32 master copy and are quantised to the desired fixed-point precision ( $q^j$ ) on-the-fly. At epoch  $j$ , the computations for the forward and backward passes ( $F$  and  $B$  blocks respectively) are performed at the current quantised precision ( $q^j$ ) and the activations as well as the gradients obtained from each layer are quantised by the quantiser module before being passed on to the next layer, or stored. After each minibatch,

the full-precision master copy of the weights is updated using a quantised gradient matrix. As discussed in Section 3.1, the quantisation level is gradually increased over the period of the training. To not compromise the final validation accuracy, switching between these optimisation levels at the correct times is crucial and controlled by a proposed precision-switching policy described in Section 3.3.

### 3.2.1. QUANTISATION STRATEGY

In order to implement quantised training, a quantisation strategy needs to be defined. The proposed quantisation strategy utilises block floating-point arithmetic (also known as dynamic fixed-point), where each fixed-point number is represented as a pair of a  $WL^{\text{net}}$ -bit signed integer  $x$  and a scale factor  $s$ , such that the value is represented as  $x \times 2^{-s}$ .

During the forward and backward passes of the training process, the weights and feature maps are both quantised, and the multiplication operations are performed at the same low precision. The quantisation method employs a stochastic rounding methodology (Gupta et al., 2015). The accumulation stage of the matrix-multiply operation is accumulated into a 32-bit fixed-point value to prevent overflow on the targeted networks.<sup>1</sup> The result of this matrix multiplication is subsequently quantised to the target wordlength before being passed as input to the next layer. Following the block floating-point scheme, quantisation is performed such that each weight and feature map matrix in the network has a single scale factor shared by all values within the matrix. The quantisation configuration for the  $i$ -th level of optimisation and the  $l$ -th layer,  $q_l^i$ , and the full set of configurations,  $q^i$ , are given in Eq. (4) respectively.

$$q_l^i = \left\langle WL^{\text{net}}, s_l^{\text{weights}}, s_l^{\text{act}} \right\rangle^i, \quad (4)$$

$$\forall l \in [1, |\mathcal{L}|] \quad \text{and} \quad q^i = \langle q_l^i \mid \forall l \in [1, |\mathcal{L}|] \rangle$$

where  $|\mathcal{L}|$  is the number of layers of the target network,  $WL^{\text{net}}$  is the fixed wordlength across the network,  $s_l^{\text{weights}}$  and  $s_l^{\text{act}}$  are the scaling factors for the weights and activations respectively, of the  $l$ -th layer for the  $i$ -th level of optimisation. As a result, for  $N$  levels, there are  $N$  distinct quantisation schemes;  $N - 1$  of these schemes are with varying fixed-point precisions and the finest level of quantisation,  $q^1$ , is single-precision floating-point (FP32). The scaling factor for a matrix  $\mathbf{X}$  is first calculated as shown in Eq. (5) and individual elements are quantised as in Eq. (6).

$$s^{\{\text{weights,act}\}} = \left\lceil \log_2 \left( \min \left( \frac{UB + 0.5}{\mathbf{X}_{\max}^{\{\text{weights,act}\}}}, \frac{LB - 0.5}{\mathbf{X}_{\min}^{\{\text{weights,act}\}}} \right) \right) \right\rceil \quad (5)$$

<sup>1</sup>The accumulator wordlength is large enough to accommodate the current CNN models without overflow.

$$x_{\text{quant}}^{\{\text{weights,act}\}} = \left\lfloor x^{\{\text{weights,act}\}} \cdot 2^{s^{\{\text{weights,act}\}}} + \text{Unif}(-0.5, 0.5) \right\rfloor \quad (6)$$

where  $\mathbf{X}_{\{\text{max,min}\}}^{\{\text{weights,act}\}}$  is either the maximum or minimum value in the weights or feature maps matrix of the current layer, LB and UB are the lower and upper bound of the current wordlength  $WL^{\text{net}}$ , and  $\text{Unif}(a,b)$  represents sampling from the uniform distribution in the range  $[a,b]$ . Eq. (5) adds 0.5 and  $-0.5$  to UB and LB respectively to ensure maximum utilisation of  $WL^{\text{net}}$ .

### 3.2.2. INFORMATION TRANSFER BETWEEN LEVELS

Employing multilevel training for CNNs requires an appropriate mechanism for transferring information between levels. To achieve this, the proposed optimiser maintains a master copy of the weights in full precision (FP32) throughout the optimisation levels. Similar to mixed-precision training (Micikevicius et al., 2018), at each level the SGD update step is performed by accumulating a fixed-point gradient value into the FP32 master copy of the weights. Starting from the coarsest quantisation level  $i = N$ , to transfer the solution from level  $i$  to level  $i - 1$ , the master copy is quantised using the quantisation scheme  $q^{i-1}$ . With this approach, the weights are maintained in FP32 and are quantised on-the-fly during run time in order to be utilised in each training step.

### 3.3. Precision-Switching Policy

MuPPET computes  $\Delta_S(\mathbf{w})$  (Eq. (1)) between gradients obtained across epochs as a proxy to measure the information that is obtained during the training process; the lower the diversity between the gradients, the less information this level of quantisation provides towards the training of the model. Therefore, the proposed method comprises a novel normalised inter-epoch version of the gradient diversity along with a run-time policy to determine the epochs in which to switch precision.

The following policy is employed to determine when a precision switch is to be performed. For a network with layers  $\mathcal{L}$  and a quantisation scheme  $q^i$  that was switched into at epoch  $e$ :

1. For each epoch  $j$  and each layer  $l \in \mathcal{L}$ , the last mini-batch's gradient,  $\nabla f_l^j(\mathbf{w})$ , is stored.
2. After  $r$  (resolution) number of epochs, the inter-epoch gradient diversity at epoch  $j$  is

$$\Delta_S(\mathbf{w})^j = \frac{\sum_{\forall l \in \mathcal{L}} \sum_{k=j-r}^j \|\nabla f_l^k(\mathbf{w})\|_2^2}{\left| \sum_{k=j-r}^j \|\nabla f_l^k(\mathbf{w})\|_2^2 \right|} \quad (7)$$

3. At an epoch  $j$ , given a set of gradient diversities

$$\mathcal{S}(j) = \left\{ \Delta_{\mathcal{S}}(\mathbf{w})^i \mid \forall e \leq i < j \right\},$$

the ratio  $p = \frac{\max \mathcal{S}(j)}{\Delta_{\mathcal{S}}(\mathbf{w})^j}$  is calculated.

4. An empirically determined decaying threshold  $T = \alpha + \beta e^{-\lambda j}$  (8) is placed on the ratio  $p$ .
5. If  $p$  violates  $T$  more than  $\gamma$  times, a precision switch is triggered and  $\mathcal{S}(j) = \emptyset$ .

As long as the gradients across epochs remain diverse,  $\Delta_{\mathcal{S}}(\mathbf{w})^j$  (Eq.(7)) at the denominator of  $p$  sustains a high value and the value of  $p$  remains low. However, when the gradients across epochs become similar,  $\Delta_{\mathcal{S}}(\mathbf{w})^j$  decreases and the value of  $p$  becomes larger. *Generalisability across epochs* is obtained as  $p$  accounts for the change in information relative to the maximum information available since the last precision change. Hence, the metric acknowledges the presence of temporal variations in information provided by the gradients. *Generalisability across networks and datasets* is maintained as  $p$  measures a ratio. Consequently, the absolute values of gradients which could vary between networks and datasets, matter less. Overall, MuPPET employs the metric  $p$  as a mechanism to trigger a precision switch whenever  $p$  violates threshold  $T$  more than  $\gamma$  times.

The likelihood of observing  $r$  gradients across  $r$  epochs that have low gradient diversity, especially at early stages of training is low. The intuition applied here is that when this does happen at a given precision, it may be an indication that information is being lost due to quantisation and thus corresponds to a high  $p$  value, which argues to move to a higher bitwidth.

### 3.3.1. HYPERPARAMETERS

The hyperparameters for the proposed MuPPET algorithm are the following: 1) values of  $\alpha$ ,  $\beta$ , and  $\lambda$  that define the decaying threshold from Eq. (8), 2) the number of threshold violations allowed before the precision change is triggered ( $\gamma$ ), 3) the resolution  $r$ , 4) the set of precisions at which training is performed, and 5) the epochs at which the learning rate is changed. The values of  $\alpha$ ,  $\beta$ ,  $\lambda$ ,  $r$ , and  $\gamma$  were set at 1, 1.5, 0.1, 3, and 2 respectively after empirical cross-validation. These were tuned by running training on AlexNet and ResNet20 on the CIFAR-10 dataset. All MuPPET hyperparameters remain the same regardless of network or dataset. Regarding training hyperparameters, batch size was increased from 128 to 256 going from CIFAR-10 to ImageNet. All other training hyperparameters, including learning rate, remained constant. Analysis of generalisability and the training hyperparameters used are presented in Section 4.1. The empirically-chosen quantised precisions at which training was performed were 8-, 12-, 14- and 16-

bit fixed-point. Precisions below this did not result in any progress towards convergence for any network.

Overall, MuPPET introduces a policy that decides at run time an appropriate point to switch between quantisation levels. After training at 16-bit fixed-point, the rest of the training is performed at FP32 until the desired validation accuracy is reached. Decaying the learning rate causes a finer exploration of the optimisation space as does increasing the quantisation level. Therefore, the learning rate was kept constant during quantised training and was decayed only after switching to FP32.

## 4. Evaluation of MuPPET

### 4.1. Generalisability

The MuPPET framework was evaluated on its applicability across epochs, networks and datasets. Fig. 2 shows the value of the metric  $p$  over the epochs in blue, and the decaying threshold described in Eq. (8) in orange. The number of epochs for which training in each precision was performed is shown by the various overlay colours. The first violation is denoted by a red dot and the second violation is not seen as it occurs exactly at the point of switching. The graphs show that across various networks and datasets, the values of  $p$  stay relatively similar, backing the choice of a universal decaying factor. Furthermore, empirical results for CIFAR-10 indicated that changing from one fixed-point precision to another too early in the training process had a negative impact on the final validation accuracy. Using a decaying threshold ensures that the value of  $p$  needs to be much higher in the initial epochs to trigger a precision change due to the volatility of  $p$  in early epochs of training.

### 4.2. Performance Evaluation

The accuracy results presented in this section utilised the proposed stochastic quantisation strategy. The methodology was developed using PyTorch. As the framework does not natively support low-precision implementations, all quantisation and computations corresponding to 8-, 12-, 14-, and 16-bit precisions were performed through emulation on floating-point hardware<sup>2</sup>. All hyperparameters not specified below were left as PyTorch defaults. For all networks, an SGD optimiser was used with batch sizes 128 on CIFAR-10/100 or 256 on ImageNet, momentum of 0.9 and weight decay of  $1e^{-4}$ .

As a baseline, an FP32 model with identical hyperparameters (except for batch size) was trained. The baseline FP32 training was performed by training for 150 epochs and reducing the learning rate by a factor of 10 at epochs 50 and 100. In order to achieve comparable final validation accuracy

<sup>2</sup>[https://github.com/ICIdsl/pytorch\\_training.git](https://github.com/ICIdsl/pytorch_training.git)

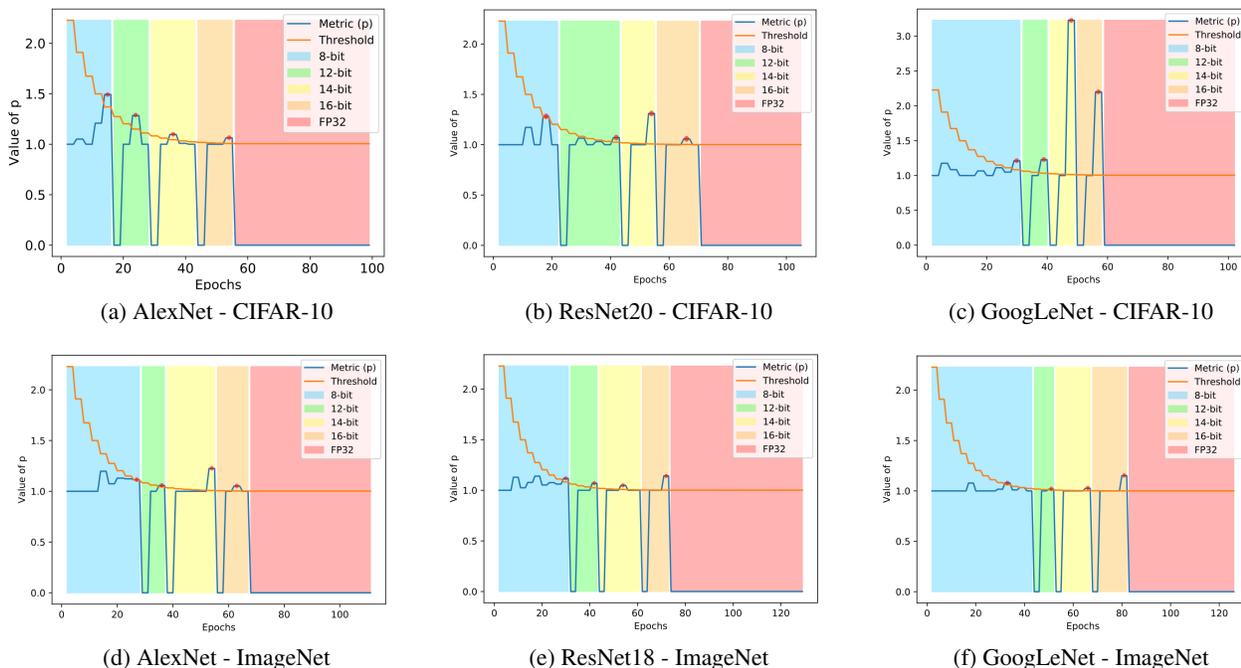


Figure 2. Demonstration of the generalisability of  $p$  over networks, datasets and epochs. Figures show the metric  $p$  (blue) and the decaying threshold  $T$  (orange) over epochs of training.

to the FP32 baseline, once MuPPET triggered a precision change out of 16-bit fixed-point, 45 training epochs at FP32 precision were performed. The learning rate was reduced by a factor of 10 every 15 FP32 training epochs. For AlexNet, ResNet18, ResNet20, and GoogLeNet, the initial learning rate was set to 0.01, 0.1, 0.1, and 0.001 respectively. The ImageNet training and validation loss curves can be seen in the Fig.3. For each of the graphs, the light grey lines indicate the epoch at which precision was switched in the MuPPET run. The green lines follow MuPPET runs and the blue FP32 training. Solid lines show validation loss and dashed training loss.

Table 1. Top-1 test accuracy (%) on CIFAR-10/100 and ImageNet (ILSVRC12 Validation Set) for FP32 baseline and MuPPET.

	CIFAR-10			CIFAR-100			ImageNet		
	FP32	MuPPET	Diff (pp)	FP32	MuPPET	Diff (pp)	FP32	MuPPET	Diff (pp)
AlexNet	75.45	74.49	-0.96	39.20	38.19	-0.99	56.21	55.33	-0.88
ResNet	90.08	90.86	0.78	64.60	65.80	1.20	69.48	69.09	-0.39
GoogLeNet	89.23	89.47	0.24	62.90	65.70	2.80	59.15	63.70	4.55

Table 1 presents the achieved Top-1 validation accuracy of MuPPET and the FP32 baseline, together with the accuracy difference in percentage points (pp). As shown on the table, MuPPET is able to provide comparable Top-1 validation accuracy to standard FP32 training across both networks and datasets. Due to a sub-optimal training setup of GoogLeNet on ImageNet, the baseline and MuPPET training severely

underperformed compared to the reported state-of-the-art works. Nevertheless, the results demonstrate the quality of training with MuPPET using identical hyperparameters. As a result, MuPPET’s performance demonstrates the effectiveness of the precision-switching strategy in achieving significant acceleration of training time (Section 4.3) at negligible cost in accuracy by running many epochs at lower precision, particularly on very large datasets.

### 4.3. Wall-clock Time Improvements

This section explores the gains in estimated wall-clock time of the current implementation of MuPPET (Current Impl.) with respect to baseline FP32 training, Mixed Precision by Micikevicius et al. (2018) and MuPPET’s ideal implementation (Table 2). For all performance results, the target platform was an NVIDIA RTX 2080 Ti GPU.

At the moment, deep learning frameworks, such as PyTorch, do not provide native support for reduced-precision hardware except for FP16 used by Micikevicius et al. (2018). Consequently, the wall-clock times on Table 2 were estimated using a performance model developed with NVIDIA’s CUTLASS library (Kerr et al., 2018) for reduced-precision general matrix-multiplication (GEMM) employing the latest Turing architecture GPUs. The entire training process in PyTorch was profiled including the following MuPPET specific operations: quantisation of weights, activations and gradients as well as all calculations pertaining to  $\Delta_S(\mathbf{w})^j$ .

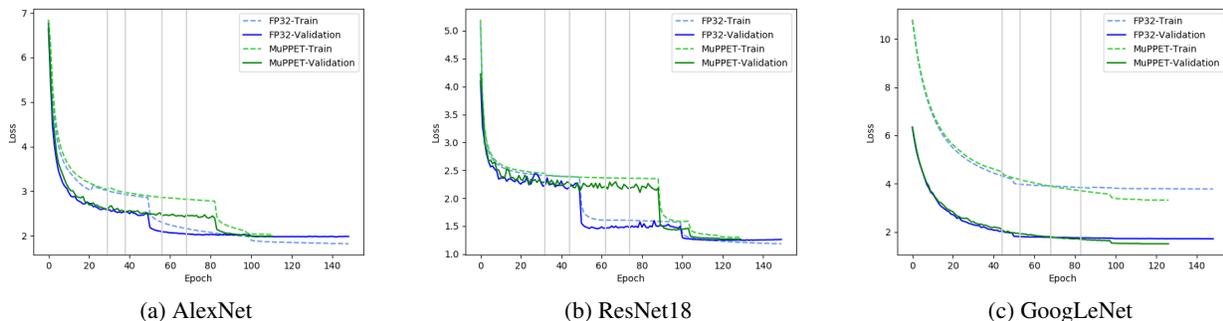


Figure 3. FP32 vs MuPPET training and validation loss graphs for different networks on the ILSVRC12 dataset.

The GEMMs that were accelerated were in the convolutional and fully-connected layers of each network. INT8 hardware was used to profile the 8-bit fixed-point computations, while FP16 hardware was used to profile 12-, 14-, and 16-bit fixed-point computations as well as Mixed Precision (Micikevicius et al., 2018) wall-clock time. CUTLASS (Kerr et al., 2018) natively implements bit-packing to capitalise on improved memory-bandwidth utilisation. To exploit further performance benefits by utilising 12- and 14-bit fixed-point computations, custom hardware such as FPGAs or ASICs would be required. The model for the current implementation is limited by the fact that current frameworks force quantisation of activations inbetween layers to happen to and from FP32. For the MuPPET (Ideal) scenario and Mixed Precision (Micikevicius et al., 2018), the model assumes native hardware utilisation for all deployed data representations which would eliminate the current overheads.

As shown on Table 2, MuPPET consistently achieves 1.25-1.32x speedup over the FP32 baseline across the networks when targeting ImageNet on the given GPU. With respect to Mixed Precision, the proposed method outperforms it on AlexNet by 1.23x and delivers comparable performance for ResNet18 and GoogLeNet. Currently, the absence of native quantisation support, and hence the necessity to emulate quantisation and the associated overheads, is the limiting factor for MuPPET to achieve higher processing speed. In this respect, MuPPET run on native hardware would yield 1.05x and 1.48x speedup for ResNet18 and GoogLeNet respectively compared to Mixed Precision. As a result, MuPPET demonstrates consistently faster time-to-accuracy (Coleman et al., 2019) compared to Mixed Precision across the benchmarks. Additionally, while Mixed Precision has already reached its limit by using FP16 on FP16-native GPUs, the 8-, 12-, 14- and 16-bit fixed-point computations enabled by MuPPET leave space for further potential speedup when targeting next- and current-generation (Fowers et al., 2018) precision-optimised fixed-point platforms. Similar to the analysis in Section 4.2, Micikevicius et al. (2018) and Wang et al. (2018) compare their schemes

Table 2. Wall-clock time (GPU hours:mins) & relative acceleration for networks targeting ImageNet.

	FP32 (Baseline)	Mixed Prec (Micikevicius et al., 2018)	MuPPET (Current Impl.)	MuPPET (Ideal)
AlexNet	30:13 (1x)	29:20 (1.03x)	23:52 (1.27x)	20:25 (1.48x)
ResNet18	132:46 (1x)	97:25 (1.36x)	100:19 (1.32x)	92:43 (1.43x)
GoogLeNet	152:28 (1x)	122:51 (1.24x)	122:13 (1.25x)	82:38 (1.84x)

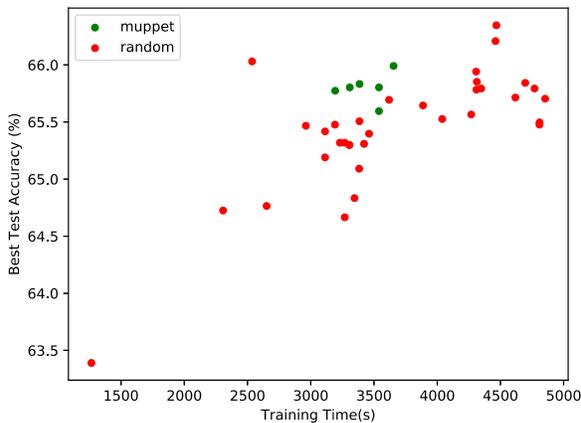


Figure 4. Test accuracy vs time trade-off for ResNet20 MuPPET runs on CIFAR-100.

to baseline FP32 training performed by them. The reported results demonstrate that their methods achieve similar accuracy results to our method by lying close to the respective FP32 training accuracy. As Wang et al. (2018) do not provide any results for gains in wall-clock times, propose a custom FP8 data type and computation, provide emulated results from an in-house (not open source) machine learning framework, and could build only a subset of their proposed ideas on an ASIC, their work could not be directly compared to our method.

#### 4.4. Precision Switching

To evaluate the ability of MuPPET to effectively choose an epoch to switch precision at, AlexNet and ResNet20 were

Table 3. Evaluation of MuPPET’s ability to tailor to dataset and network combination using Top-1 test accuracy on CIFAR-100.

	ResNet20		GoogLeNet	
	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-100
<b>ResNet20</b>	65.01	<b>65.80</b>	-	65.0
<b>GoogLeNet</b>	-	64.00	64.70	<b>65.70</b>

first trained using MuPPET on the CIFAR-100 dataset. The hyperparameters for MuPPET were kept the same across all runs. From the results it was noted that training at reduced precision and not switching at all causes a drop in validation accuracy of 1.4% and 1.3% for AlexNet and ResNet20 respectively, hence demonstrating the need to switch precisions when training at bitwidths as low as 8-bit fixed-point.

To demonstrate the benefits of a precision-switching methodology, two further sets of experiments were conducted on ResNet20 using CIFAR-100. First, 34 training runs were performed (34 red dots in Fig. 4), where for each run four epochs along the standard training duration were randomly selected and used as the switching points. Fig. 4 shows the best test accuracy achieved by each of the runs and the training time as estimated by our performance model described in Section 4.3. It shows that for a given time-budget, MuPPET runs (6 green dots) outperform on average all other experiment sets, demonstrating the need for a precision-switching policy that is real-time in order to achieve a good accuracy-to-training-time trade-off. The accuracy and time have mean  $\pm$  standard deviation of  $(65.54\% \pm 0.06, 3429.19s \pm 46.1)$  and  $(65.44\% \pm 0.09, 3679.15s \pm 141.1)$  for the green and red points respectively.

To demonstrate that MuPPET tailors the training process to the network-dataset pair while using a policy that is agnostic of the network-dataset combination, the experiments in Table 3 were performed. The switching strategy generated by MuPPET for both ResNet20 and GoogLeNet on CIFAR-10 and CIFAR-100 were recorded. The experiments shown in Table 3 used the CIFAR-100 dataset. The bold values show the accuracy achieved by the schedule proposed by MuPPET for the corresponding network-dataset pair. For the ResNet20 row, column 1 shows the accuracy achieved if the precision schedule proposed by MuPPET on the CIFAR-10 dataset for ResNet20 was applied to CIFAR-100 training (dataset tailoring). In the same row, column 4 shows the accuracy obtained when CIFAR-100 training was performed on ResNet20 using GoogLeNet’s CIFAR-100 MuPPET schedule (network tailoring). Similar results are shown for GoogLeNet. As all experiments underperform compared to the MuPPET training process, Table 3 demonstrates the ability of MuPPET to tailor the process to a network-dataset combination.

## 5. Conclusion

This paper proposes MuPPET, a novel low-precision CNN training scheme that combines the use of fixed-point and floating-point representations to produce a network trained for FP32 inference. By introducing a precision-switching mechanism that decides at run time an appropriate transition point between different precision regimes, the proposed framework achieves Top-1 validation accuracies comparable to that achieved by state-of-the-art FP32 training regimes while delivering significant speedup in terms of training time. Quantitative evaluation demonstrates that MuPPET’s training strategy generalises across CNN architectures and datasets by adapting the training process to the target CNN-dataset pair during run time. Overall, MuPPET enables the utilisation of the low-precision hardware units available on modern specialised processors, such as next-generation GPUs, FPGAs and TPUs, to yield improvements in training time and energy efficiency without impacting the resulting accuracy. Future work will focus on applying the proposed framework to the training of LSTMs, where the training process is more sensitive to gradient quantisation (Rizakis et al., 2018), as well as on the extension of MuPPET to include batch size and learning rate as part of its hyperparameters. Furthermore, we will explore improved quantisation techniques that could enable training convergence for bitwidths even lower than 8-bit fixed-point.

## 6. Acknowledgements

The support of the EPSRC Centre for Doctoral Training in High Performance Embedded and Distributed Systems (HiPEDS, Grant Reference EP/L016796/1) is gratefully acknowledged.

## References

- Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1709–1720. 2017.
- Burgess, J. RTX ON – The NVIDIA TURING GPU. In *IEEE Hot Chips 31 Symposium (HCS)*, pp. 1–27, 2019.
- Cai, H., Chen, T., Zhang, W., Yu, Y., and Wang, J. Efficient Architecture Search by Network Transformation. In *AAAI Conference on Artificial Intelligence*, 2018.
- Chen, X., Hu, X., Xu, N., Zhou, H., Zhou, H., and Xu, N. FxpNet: Training deep convolutional neural network in fixed-point representation. In *International Joint Conference on Neural Networks (IJCNN 2017)*, May 2017.
- Chen, X., Chen, D. Z., Han, Y., and Hu, X. S. moDNN:

- Memory Optimal Deep Neural Network Training on Graphics Processing Units. *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, 30(3):646–661, 2019.
- Coleman, C., Kang, D., Narayanan, D., Nardi, L., Zhao, T., Zhang, J., Bailis, P., Olukotun, K., Ré, C., and Zaharia, M. Analysis of DAWN-Bench, a Time-to-Accuracy Machine Learning Performance Benchmark. *SIGOPS Oper. Syst. Rev.*, 53(1):14–25, 2019.
- Courbariaux, M., Bengio, Y., and David, J.-P. Training deep neural networks with low precision multiplications. In *Workshop in International Conference in Learning Representations*, 2015.
- De Sa, C., Feldman, M., Ré, C., and Olukotun, K. Understanding and Optimizing Asynchronous Low-Precision Stochastic Gradient Descent. In *International Symposium on Computer Architecture (ISCA)*, pp. 561–574, 2017.
- De Sa, C. M., Zhang, C., Olukotun, K., Ré, C., and Ré, C. Taming the Wild: A Unified Analysis of Hogwild-Style Algorithms. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 2674–2682, 2015.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. ImageNet: A Large-Scale Hierarchical Image Database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255, 2009.
- Devarakonda, A., Naumov, M., and Garland, M. AdaBatch: Adaptive Batch Sizes for Training Deep Neural Networks. *arXiv*, 2017.
- Fowers, J., Ovtcharov, K., Papamichael, M., Massengill, T., Liu, M., Lo, D., Alkalay, S., Haselman, M., Adams, L., Ghandi, M., Heil, S., Patel, P., Sapek, A., Weisz, G., Woods, L., Lanka, S., Reinhardt, S. K., Caulfield, A. M., Chung, E. S., and Burger, D. A Configurable Cloud-Scale DNN Processor for Real-Time AI. In *International Symposium on Computer Architecture (ISCA)*, pp. 1–14, 2018.
- Gan, C., Wang, N., Yang, Y., Yeung, D.-Y., and Hauptmann, A. G. DevNet: A Deep Event Network for Multimedia Event Detection and Evidence Recounting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2568–2577, 2015.
- Gupta, S., Agrawal, A., Gopalakrishnan, K., and Narayanan, P. Deep Learning with Limited Numerical Precision. In *32nd International Conference on Machine Learning (ICML)*, pp. 1737–1746, 2015.
- He, K., Gkioxari, G., Dollar, P., and Girshick, R. Mask R-CNN. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q. Densely Connected Convolutional Networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269, 2017.
- Johnson, T. B. and Guestrin, C. Training Deep Models Faster with Robust, Approximate Importance Sampling. In *Advances in Neural Information Processing Systems 31*, pp. 7265–7275, 2018.
- Jouppi, N. P. et al. In-Datcenter Performance Analysis of a Tensor Processing Unit. In *44th Annual International Symposium on Computer Architecture (ISCA)*, pp. 1–12, 2017.
- Kerr, A., Merrill, D., Demouth, J., and Tran, J. CUTLASS: Fast Linear Algebra in CUDA C, Sep 2018. URL <https://devblogs.nvidia.com/cutlass-linear-algebra-cuda/>.
- Kouris, A. and Bouganis, C. Learning to Fly by MySelf: A Self-Supervised CNN-Based Approach for Autonomous Navigation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1–9, 2018.
- Kouris, A., Venieris, S. I., and Bouganis, C. CascadeCNN: Pushing the Performance Limits of Quantisation in Convolutional Neural Networks. In *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*, pp. 155–1557, 2018.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision (ECCV)*, pp. 740–755, 2014.
- Loquercio, A., Maqueda, A. I., del Blanco, C. R., and Scaramuzza, D. DroNet: Learning to Fly by Driving. *IEEE Robotics and Automation Letters*, 3(2):1088–1095, 2018.
- Micikevicius, P., Narang, S., Alben, J., Diamos, G., Elsen, E., Garcia, D., Ginsburg, B., Houston, M., Kuchaiev, O., Venkatesh, G., and Wu, H. Mixed Precision Training. In *International Conference on Learning Representations (ICLR)*, 2018.
- Migdalas, A., Pardalos, P. M., and Värbrand, P. *Multilevel Optimization: Algorithms and Applications*, volume 20. Springer Science & Business Media, 2013.
- Peng, X., Li, L., and Wang, F.-Y. Accelerating Minibatch Stochastic Gradient Descent using Typicality Sampling. *arXiv:1903.04192 [cs, math, stat]*, 2019.
- Rhu, M., Gimelshein, N., Clemons, J., Zulfiqar, A., and Keckler, S. W. vDNN: Virtualized Deep Neural Networks for Scalable, Memory-efficient Neural Network Design. In *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 18:1–18:13, 2016.

- Rizakis, M., Venieris, S. I., Kouris, A., and Bouganis, C. Approximate FPGA-based LSTMs under Computation Time Constraints. In *14 International Symposium on Applied Reconfigurable Computing (ARC)*, 2018.
- 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.
- Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. In *AAAI Conference on Artificial Intelligence*, 2017.
- Wang, N., Choi, J., Brand, D., Chen, C.-Y., and Gopalakrishnan, K. Training Deep Neural Networks with 8-bit Floating Point Numbers. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 7675–7684. 2018.
- Yin, D., Pananjady, A., Lam, M., Papailiopoulos, D., Ramchandran, K., and Bartlett, P. Gradient Diversity: a Key Ingredient for Scalable Distributed Learning. In *21st International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1998–2007, 2018.
- Zhong, Z., Yan, J., Wu, W., Shao, J., and Liu, C. Practical Block-Wise Neural Network Architecture Generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2423–2432, 2018.
- 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*, 2016.