
Is Integer Arithmetic Enough for Deep Learning Training?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The ever-increasing computational complexity of deep learning models makes their
2 training and deployment difficult on various cloud and edge platforms. Replacing
3 floating-point arithmetic with low-bit integer arithmetic is a promising approach
4 to save energy, memory footprint, and latency of deep learning models. As such,
5 quantization has attracted the attention of researchers in recent years. However,
6 using integer numbers to form a fully functional *integer* training pipeline including
7 forward pass, back-propagation, and stochastic gradient descent is not studied
8 in detail. Our empirical and mathematical results reveal that *integer arithmetic*
9 is enough to train deep learning models. Unlike recent proposals, instead of
10 quantization, we directly switch the number representation of computations. Our
11 novel training method forms a fully integer training pipeline that does not change
12 the trajectory of the loss and accuracy compared to floating-point, nor does it
13 need any special hyper-parameter tuning, distribution adjustment, or gradient
14 clipping. Our experimental results show that our proposed method is effective in a
15 wide variety of tasks such as classification (including vision transformers), object
16 detection, and semantic segmentation.

17 1 Introduction

18 Recently, deep learning models have evolved to deliver acceptable performance in many areas such
19 as computer vision, natural language processing, and speech recognition. However, the number
20 of parameters and computational complexity of most deep learning applications have increased
21 significantly. This ever-increasing computational complexity makes the deployment of deep neural
22 networks harder on edge devices. Furthermore, the energy required to train a deep learning model
23 increases proportionately with computational complexity. Thus, the training cost of these large and
24 complex models also increases significantly on the high-end cloud servers. Although, quantization
25 techniques are commonly used to accelerate *inference* of deep learning models, here we target a
26 fully integer training pipeline (*i.e. forward propagation, back-propagation and SGD*). Unlike recent
27 proposals in quantized back-propagation, we directly change the number representation of floating-
28 point values. Furthermore, we show empirically and theoretically that our proposed integer training
29 methodology is effective in training deep learning models with *integer-only arithmetic*.

30 Despite the fact that accelerated training using integer back-propagation seems intriguing, there
31 are some challenges associated with integer gradient computation. For example, (i) gradients must
32 be scaled correctly in order to be adapted to the limited dynamic range of the integer number (e.g.
33 `int8`: $[-127, 126]$). (ii) The numerical error of the gradient must be *unbiased* in order to preserve
34 the convergence trajectory of the training algorithm. A small numerical error can be accumulated
35 through the course of training and change the convergence behaviour. (iii) The integer training must

36 be *distribution independent*, i.e. the training method should not depend on the distribution of the
37 gradients, weights, or training data.

38 We propose a novel integer training method designed to address aforementioned challenges (i), (ii)
39 and (iii) simultaneously. Having a linear dynamic fixed-point mapping coupled with a non-linear
40 inverse mapping is the key to full integer training.

41 To this end, we make the following contributions:

- 42 • A hardware-friendly integer training method is proposed based on extracting the maximum
43 floating-point exponent of the tensors as the *scale*. We propose linear fixed-point mapping
44 of tensors, while the corresponding inverse mapping for floating-point is non-linear. Our
45 proposed method addresses all the previously mentioned challenges: (i) it computes the
46 scales dynamically, and moreover the scale does not need to be adjusted in the course of
47 training; (ii) it provides an unbiased estimation of the gradients and consequently, its conver-
48 gence trajectory closely follows the floating-point version; (iii) our proposed representation
49 mapping does not depend on the distribution of the training data, weights, or gradients.
- 50 • We study the optimality gap of our proposed integer training algorithm and show it is
51 analogous to its floating-point counterpart in the course of training (**Theorem 1**). Our
52 analysis of stochastic gradient descent with our proposed fixed-point gradients shows the
53 original floating-point optimality gap is only shifted by a negligible amount (**Remark 3**).
- 54 • Our proposed method effectively performs all operations required in modern neural networks
55 using *integer-only arithmetic*. For instance, the computation of linear layer, convolutional
56 layer, batch-norm and layer-norm, residual connections and also stochastic gradient descent
57 (including gradients, momentum, weight decay, and weight update) are all performed in
58 integer arithmetic. To the best of our knowledge, this is the first time that *back-propagation*
59 of a batch-norm, and the computation of stochastic gradient descent (SGD) is performed in
60 integer arithmetic with negligible loss in the accuracy for large datasets such as ImageNet.

61 The rest of this paper is structured as follows. Section 2 reviews some previous works in the field of
62 quantized back-propagation and quantifies the similarities and differences with our *representation*
63 *mapping* method. Section 3 discusses our integer training methodology in detail. Theoretical aspects
64 of our integer training method on SGD are studied in Section 4. Experimental results supporting our
65 integer training methodology and convergence theory are presented in Section 5.

66 2 Related works

67 Banner et al. [1], proposed a bifurcated back-propagation method where the gradient of weights are
68 in floating-point format and input gradients are in 8-bit fixed-point format. Moreover, they proposed
69 a range-based batch-norm which is more tolerant to quantization noise. The main difference of
70 our proposed algorithm with Banner et al. [1] is that our integer training method does not demand
71 bifurcation. In other words, our gradients with respect to weights and inputs are represented and
72 computed using `int8` format. Furthermore, in our method, the back-propagation of batch-norm, as
73 well as stochastic gradient descent (SGD) algorithm, are also represented and computed using integer
74 values. Zhang et al. [2] have proposed a layer-wise “precision-adaptive” quantization method that
75 does not affect the distribution of the data. In their method, the quantization error is measured and
76 the quantization scale is adjusted over the course of training accordingly. Similarly, in Zhao et al. [3]
77 a distribution adaptive quantization method is proposed that takes into account the distribution of
78 gradients in the channel dimension. In addition, they introduced a gradient clipping strategy in order
79 to normalize the magnitude of the gradients in the back-propagation. The scale of the quantization
80 is adapted to the distribution of the gradients and adjusted iteratively over the course of training.
81 Unlike the distribution adaptive methods such as [2, 3], our proposed integer training method does
82 not depend on the distribution of training data, weights, and gradients. Thus, there is no need to
83 adjust the scale iteratively in the course of training. In Zhu et al. [4], a “direction sensitive gradient
84 clipping method” is proposed to avoid inappropriate updates in the back-propagation. Moreover, they
85 proposed a learning rate scaling that tackles the problem of unbiased gradient error accumulation. In
86 a similar work, Sakr and Shanbhag [5] proposed a precision assignment methodology to improve the
87 convergence of the quantized back-propagation. This precision assignment is based on some criterion
88 on internal accumulator noise, quantization noise of backward and forward propagation, and gradient

89 clipping. Unlike [4, 5], our proposed method enjoys having the advantage of *unbiased gradients*,
 90 hence the convergence trajectory closely follows the floating-point counterpart. Thus, our method
 91 does not need any specific gradient correction, gradient clipping strategy, and learning rate corrections
 92 as opposed to [4, 5]. Jin et al. [6] proposed a unified fixed-point and parameterized clipping activation
 93 to achieve high accuracy. Furthermore, they proposed a method that directly trains the fractional
 94 length (i.e. scale) of the fixed-point quantization. Additionally, they use a double forward batch-norm
 95 fusion to determine the scaling factor of the batch-norm layer in the forward propagation, while the
 96 back-propagation of batch-norm remains floating-point. In contrast, our proposed integer training
 97 method is capable of implementing integer variant of the batch-norm layer in forward propagation
 98 without double fusion as opposed to [6]. We also perform batch-norm’s back-propagation in integer
 99 arithmetic which was ignored in the previous works. Note that back-propagation of batch-norm
 100 and layer-norm is sensitive to quantization, as such, naive quantization leads to training divergence.
 101 Additionally, we perform stochastic gradient descent (SGD) computation including weight update,
 102 weight decay, and momentum in integer arithmetic with no significant loss of accuracy. In our
 103 proposed method, we used dynamic fixed-point number format [7] as the primary number format
 104 for integer training. This specific fixed-point number format (also known as block floating-point
 105 format) has been used previously to quantize *inference* of deep learning models in Courbariaux et al.
 106 [7], Drumond et al. [8], and Rouhani et al. [9].

107 3 Methodology

108 Common quantization methods, which use division and clipping techniques, are inefficient for back-
 109 propagation. Therefore, we decided to go one step deeper and change the number format directly.
 110 Here, we propose a hardware-friendly *number representation mapping* using the dynamic fixed-point
 111 number format where the scale is defined *per tensor*. In this approach, each tensor can be represented
 112 by its `int8` version while it is multiplied by a shared scale, as opposed to other methods that allow
 113 multiple shared scales for different partitions of the tensor (see Rouhani et al. [9, Figure 4]). One
 114 shared scale per tensor makes the computations easier in the computing hardware (e.g. CPU or GPU).
 115 However, the training algorithm diverges if the representation mapping is not executed properly. To
 116 tackle this problem, we suggest two subtle changes; (i) we propose to perform fixed-point mapping
 117 of a tensors in a *linear* fashion while the inverse mapping is *non-linear* as explained in Sections
 118 3.1 and 3.2. This is the key to success of our algorithm since a linear fixed-point mapping allows
 119 monotonic conversion of floating-point format to fixed-point, while a non-linear inverse mapping
 120 allows preserving information. (ii) We suggest to use stochastic rounding in the back-propagation
 121 in a way that preserves the expected value of the tensors as well as their vital statistics, such as the
 122 mean. Stochastic rounding in conjunction with representation mapping is crucial to keep the integer
 123 training loss trajectory close to its floating-point counterpart.

124 3.1 Linear fixed-point mapping

125 Our proposed integer training method is based on manipulating the floating-point number format.
 126 This is a very simple and effective way of converting floating-point values to fixed-point values as
 127 opposed to other commonly used methods that involve division operation. Our method essentially
 128 uses *shift* and *round* operations to convert floating-point to fixed-point format, see Figure 1(a).

129 In this method, a tensor comprising n floating-point numbers (f_1, f_2, \dots, f_n) is converted to a fixed-
 130 point tensor. First, as shown in Figure 1(a), sign, exponent, and mantissa of each floating-point
 131 number are extracted using the *unpack to integer* function. For instance, f_1 is unpacked to s_1, e_1, m_1
 132 where (s, e, m) is used to denote (*sign, exponent, mantissa*). We simply find the *maximum* element of
 133 e_1, e_2, \dots, e_n i.e. $e_{\max} = \max(e_1, e_2, \dots, e_n)$ to extract the shared scale of the tensor. Subsequently,
 134 the original mantissas are shifted to the right according to the difference of their exponent and e_{\max} ,
 135 to get the individual mantissas. For instance, m_1 is shifted to right by $e_{\max} - e_1$ which is denoted by
 136 $m_1 \gg (e_{\max} - e_1)$. By shifting mantissas to right, we intentionally push the small values to the
 137 *sub-normal* region, where the most significant bit of mantissa is not 1 in binary format i.e. $(1)_2$, see
 138 Zuras et al. [10]. Pushing the mantissas to the sub-normal region is performed to align their exponents
 139 with e_{\max} and this unifies the scale for fixed-point values in a tensor. In the next step, to create
 140 `int8` mantissas, the 24-bit single floating-point mantissas are further rounded to 7-bit mantissas to
 141 construct signed `int8` values i.e 7-bit unsigned integer and 1 sign bit.

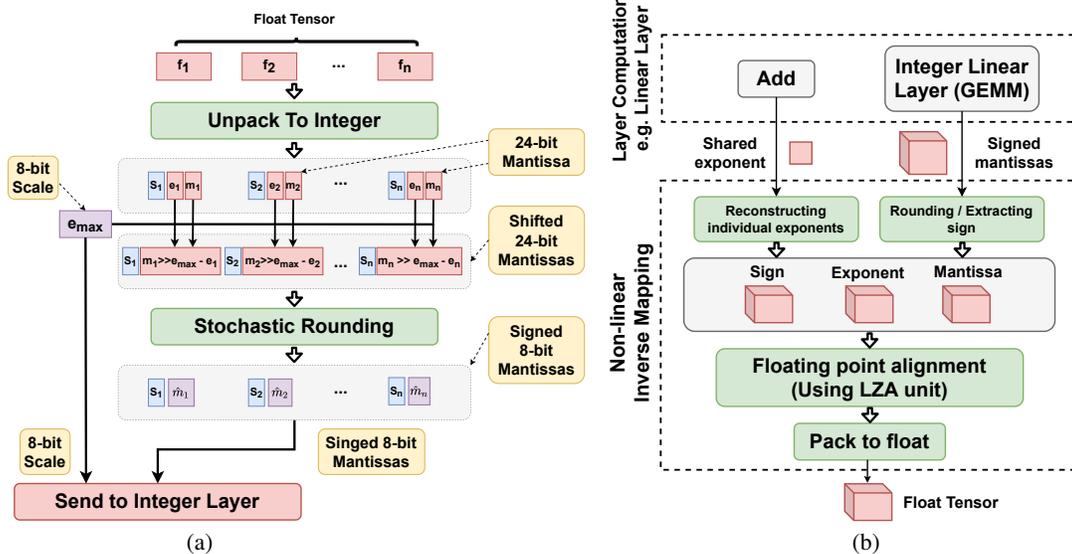


Figure 1: (a) Linear fixed-point mapping, (b) Non-linear inverse mapping.

142 3.2 Non-linear inverse mapping

143 Our proposed fixed-point mapping method is performed by pushing mantissa values to the sub-normal
 144 region and rounding them stochastically. Hence, the inverse mapping must be performed using a
 145 floating-point alignment module that normalizes the mantissas. A *normalized* floating-point mantissa
 146 is required to start with the most significant bit $(1)_2$. For instance, an alignment module converts
 147 $2^{127} \times (0.0101)_2$ to $2^{127-2} \times (1.0100)_2$. The alignment module is a very well-known logic circuit
 148 that is available in the commodity hardware using Leading Zero Anticipator (LZA) block [11].
 149 Coupling a linear fixed-point mapping with a non-linear inverse mapping in this way keeps the
 150 number format close to floating-point. Moreover, combining them with stochastic rounding, results
 151 in having unbiased fixed-point variant that forces the trajectory of the integer training closely follow
 152 the floating-point counterpart.

153 Figure 1(b) shows the inverse mapping unit. This unit receives a tensor of mantissas and a single
 154 value of shared exponent computed by an integer layer. Then a tensor of exponents is reconstructed
 155 by repeating the value of the shared exponent. The size of this tensor is equal to the mantissa tensor's
 156 size. Moreover, the mantissa tensor is rounded and its sign is extracted. Note that the rounding is
 157 used here because the output of the previous layer might have some excessive mantissa bits; this
 158 normally happens in matrix multiplication and convolution operations. Then the sign, exponent, and
 159 mantissa tensors are sent to an alignment unit that shifts the mantissa and adjust the exponent in order
 160 to normalize the floating-point number format [10]. Finally, the result is packed and carried out to the
 161 next layer.

162 3.3 Integer layer computations

163 When the fixed-point mapping of floating-point values is completed, the fixed-point values are used to
 164 perform the *integer-only* layer computations. For instance, let us consider a linear layer which consists
 165 of a General Matrix Multiplication (GEMM) operation, but, the idea can be generalized to other
 166 types of layers. Figure 2 demonstrates the layer-wise integer computations of a linear layer where the
 167 shared exponents and integer mantissas are treated separately. As shown in the Figure 2, in order
 168 to perform an integer-only GEMM operation, we need to multiply scales $(2^{e_{\max 1}} \times 2^{e_{\max 2}})$. This
 169 multiplication performs an integer addition operation on the exponents $(e_{\max 1} + e_{\max 2})$. Furthermore,
 170 the integer mantissas are sent to an integer GEMM module to compute the output mantissa tensor.

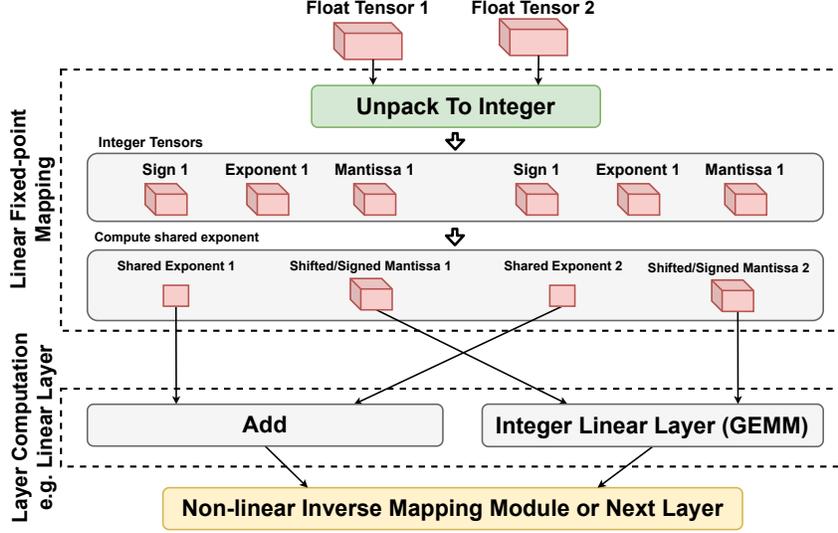


Figure 2: A fully integer linear layer.

171 3.4 Understanding of representation mapping

172 Here we provide an intuition of how our proposed representation mapping approach works. Let
 173 us denote A as a tensor and \hat{A} as its fixed-point version in the way that is introduced earlier. In
 174 addition, random variables A_i and \hat{A}_i are i^{th} element of those tensors. Also note that A_i and \hat{A}_i
 175 are different *representations* of the same real number, but one is in floating-point format and the other
 176 is in dynamic fixed-point format. Thus, one can relate A_i and \hat{A}_i with a random error term δ_i as
 177 $\hat{A}_i = A_i + \delta_i$. Since in our training method we used stochastic rounding [12], $\mathbb{E}\{\hat{A}_i\} = A_i$, or
 178 equivalently $\mathbb{E}\{\delta_i\} = 0$ (see Appendix A.1). In other words, the fixed-point value is on the average
 179 equal to the floating-point value. Note that here we consider single precision floating-point values as
 180 a surrogate of real values.

181 **Linear and convolutional layers:** For these two types of layers, both forward and backward
 182 propagation computations are based on inner products. Thus, it is easy to see that our proposed
 183 fixed-point inner product $\hat{C}_{ij} = \sum_k \hat{A}_{ik} \hat{B}_{kj}$ is an unbiased estimator of the floating-point inner
 184 product

$$\mathbb{E}\{\hat{C}_{ij}\} = \mathbb{E}\left\{\sum_k \hat{A}_{ik} \hat{B}_{kj}\right\} = \mathbb{E}\left\{\sum_k (A_{ik} + \delta_{ik}^A)(B_{kj} + \delta_{kj}^B)\right\} = \sum_k A_{ik} B_{kj} = C_{ij}. \quad (1)$$

185 **Residual connections:** A residual connection involves element-wise addition of two tensors. Each
 186 fixed-point element $\hat{C}_{ij} = \hat{A}_{ij} + \hat{B}_{ij}$ is an unbiased estimator of its floating-point version

$$\mathbb{E}\{\hat{C}_{ij}\} = \mathbb{E}\{\hat{A}_{ij} + \hat{B}_{ij}\} = \mathbb{E}\{(A_{ij} + \delta_{ij}^A) + (B_{ij} + \delta_{ij}^B)\} = A_{ij} + B_{ij} = C_{ij}. \quad (2)$$

187 **Batch-norm:** A batch-norm layer is defined as

$$\hat{\omega} = \frac{\hat{A} - \hat{\mu}}{\sqrt{\hat{\sigma}^2 + \epsilon}} + \hat{\beta}. \quad (3)$$

188 For this type of layer, it is important to compute signal statistics such as mean $\hat{\mu}$ and variance $\hat{\sigma}^2$
 189 correctly in integer arithmetic. Thus, the fixed-point mean $\hat{\mu}$ is also an unbiased estimator of the true
 190 mean μ

$$\mathbb{E}\{\hat{\mu}\} = \mathbb{E}\left\{\frac{\sum_{i=1}^N \hat{A}_i}{N}\right\} = \frac{1}{N} \mathbb{E}\left\{\sum_{i=1}^N (A_i + \delta_i)\right\} = \mu. \quad (4)$$

191 Then the following derivation holds for fixed-point variance $\hat{\sigma}^2$

$$\mathbb{E}\{\hat{\sigma}^2\} = \mathbb{E}\left\{\frac{\sum_{i=1}^N (\hat{A}_i - \hat{\mu})^2}{N}\right\} = \frac{1}{N}\mathbb{E}\left\{\sum_{i=1}^N [(A_i - \mu)^2 + \delta_i^2]\right\} = \sigma^2 + \sigma_\delta^2, \quad (5)$$

192 where σ^2 is the true floating-point variance of the batch and σ_δ^2 is the variance of the noise introduced
193 by the linear mapping to fixed-point. Note that the error variance σ_δ^2 is rather small and can be
194 integrated to ϵ in the denominator of equation (3).

195 4 Theoretical analysis of SGD using the proposed method

196 The generic equation of weight update in the k^{th} iteration of SGD is

$$w_{k+1} = w_k + \alpha_k g(w_k, \xi_k), \quad (6)$$

197 where $g(w_k, \xi_k)$ is the estimated gradients of random samples of the batch generated by the seed ξ_k ,
198 and α_k is the learning rate. We make the following common assumption in the sequel.

199 **Assumption 1 (Lipschitz-continuity).** The loss function $\mathcal{L}(w)$ is continuously differentiable and its
200 gradients $\nabla\mathcal{L}(w)$ satisfies the following inequality where $L > 0$ is the Lipschitz constant

$$\mathcal{L}(w) \leq \mathcal{L}(\bar{w}) + \nabla\mathcal{L}(\bar{w})^\top (w - \bar{w}) + \frac{1}{2}L\|w - \bar{w}\|_2^2; \quad \forall w, \bar{w} \in \mathbb{R}^d. \quad (7)$$

Assumption 2. (i) $\mathcal{L}(w_k)$ is bounded. (ii) Estimated gradients $g(w_k, \xi_k)$ is an unbiased estimator of
the true gradients of the loss function

$$\nabla\mathcal{L}(w_k)^\top \mathbb{E}_{\xi_k}\{g(w_k, \xi_k)\} = \|\nabla\mathcal{L}(w_k)\|_2^2 = \|\mathbb{E}_{\xi_k}\{g(w_k, \xi_k)\}\|_2^2,$$

201 and (iii,a) there exist scalars $M \geq 0$ and $M_V \geq 0$ such that for all iterations of SGD
202 $\mathbb{E}_{\xi_k}\{g(w_k, \xi_k)\} \leq M + M_V\|\nabla\mathcal{L}(w_k)\|_2^2$.

Note that here we define

$$\mathbb{V}_{\xi_k}\{g(w_k, \xi_k)\} := \mathbb{E}_{\xi_k}\{\|g(w_k, \xi_k)\|_2^2\} - \|\mathbb{E}_{\xi_k}\{g(w_k, \xi_k)\}\|_2^2.$$

203 Also from *Assumption 2. (ii) and (iii,a)*, the second moment bound can be derived

$$\mathbb{E}_{\xi_k}\{\|g(w_k, \xi_k)\|_2^2\} \leq M + M_G\|\nabla\mathcal{L}(w_k)\|_2^2 \quad \text{with } M_G := 1 + M_V. \quad (8)$$

204 **Effect of gradient variance on convergence:** The quality of the estimated gradients $g(w_k, \xi_k)$
205 directly affects the convergence of the SGD. The effect of first and second moments of gradient are
206 already studied on *real numbers* in the literature.

207 **Lemma 2.** Suppose *Assumption 2* is true, then we have

$$\mathbb{E}_{\xi_k}\{\mathcal{L}(w_{k+1})\} - \mathcal{L}(w_k) \leq -(1 - \frac{1}{2}\alpha_k LM_G)\alpha_k\|\nabla\mathcal{L}(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 LM. \quad (9)$$

208 **Proof:** See Bottou et al. [13, Lemma 4.4].

209 Inequality (9) shows the effect of gradient variance bounds, M and M_G , on each iterate of SGD, and
210 shows the greater the variance, the more deterioration in the quality of SGD steps. The first term,
211 $-(1 - \frac{1}{2}\alpha_k LM_G)\alpha_k\|\nabla\mathcal{L}(w_k)\|_2^2$ contributes to the decrease of the loss function while the second
212 term, $\frac{1}{2}\alpha_k^2 LM$, prevents it. Upon choosing the correct gradient estimates, the right hand side of
213 inequality (9) is bounded by a deterministic quantity, and asymptotically ensures sufficient descent of
214 the loss $\mathcal{L}(w)$. Note that the expectation \mathbb{E}_{ξ_k} is taken over random samples with seed ξ_k .

215 4.1 Fixed-point gradient noise

216 When performing representation mapping in the back-propagation, the quality of the gradients
217 deteriorate. Thus, there is a need to consider the effect of fixed-point mapping variance. Then,
218 *Assumption 2 (iii,a)* should be modified accordingly.

219 **Remark 1.** Note that *Assumption 2 (i), (ii)* still hold after fixed-point mapping because of stochastic
 220 rounding, i.e. the fixed-point gradient remains an unbiased estimator of gradient (refer to Appendix
 221 A.1).

Assumption 2 (iii,b). When the gradients are in fixed-point format i.e. $\hat{g}(w_k, \xi_k)$, there exist scalars
 $M \geq 0$, $M_V \geq 0$, $M^q \geq 0$ and $M_V^q \geq 0$ such that for all iterations of SGD

$$\mathbb{V}_{\xi_k} \{\hat{g}(w_k, \xi_k)\} \leq M + M^q + (M_V + M_V^q) \|\nabla \mathcal{L}(w_k)\|_2^2.$$

222 If $\tilde{M} := M + M^q$ and $\tilde{M}_V := M_V + M_V^q$, then *Assumption 2 (iii,b)* takes the exact form of
 223 *Assumption 2 (iii,a)* i.e. $\mathbb{V}_{\xi_k} \{g(w_k, \xi_k)\} \leq \tilde{M} + \tilde{M}_V \|\nabla \mathcal{L}(w_k)\|_2^2$. However, here we separated M^q
 224 and M_V^q to emphasize the effect of fixed-point mapping on the true gradients.

225 **Remark 2.** If *Assumption 2 (iii,b)* holds true, inequality (9) can be transformed to its fixed-point
 226 version

$$\begin{aligned} \mathbb{E}_{\xi_k} \{\mathcal{L}(w_{k+1})\} - \mathcal{L}(w_k) &\leq -\left(1 - \frac{1}{2} \alpha_k L (M_G + M_G^q)\right) \alpha_k \|\nabla \mathcal{L}(w_k)\|_2^2 + \frac{1}{2} \alpha_k^2 L (M + M^q) \\ &\quad \text{with } M_G^q := 1 + M_V^q. \end{aligned} \tag{10}$$

227 Inequality (10) shows the effect of **added** representation mapping variance with bounds, M^q and M_G^q ,
 228 on each iterate of SGD. This observation shows that fixed-point mapping degrades the convergence
 229 of SGD unless its variance bounds are relatively small, or controlled by the learning rate. As an
 230 example, refer to Appendix A.2 for analytical derivations of M^q and M_G^q for the back-propagation of
 231 a linear layer involving a fixed-point inner product.

232 4.2 Strongly convex and locally convex loss

233 **Assumption 3 (Strong convexity).** The loss function $\mathcal{L}(w)$ is differentiable and strongly convex. We
 234 recall that strong convexity for differentiable functions is equivalent to the following inequality with
 235 constant $c > 0$

$$\mathcal{L}(w) \geq \mathcal{L}(\bar{w}) + \nabla \mathcal{L}(\bar{w})^\top (w - \bar{w}) + \frac{1}{2} c \|w - \bar{w}\|_2^2; \quad \forall w, \bar{w} \in \mathbb{R}^d. \tag{11}$$

236 A strongly convex function has a unique minimum point at w_* with the loss value $\mathcal{L}_* = \mathcal{L}(w_*)$.

237 **Theorem 1.** Suppose *Assumptions 1, 2(i), 2(ii), 2 (iii,b), 3* are all true, then a SGD method running
 238 with fixed-point gradients i.e. $\hat{g}(w_k, \xi_k)$ and a fixed learning rate $0 < \bar{\alpha} \leq \frac{1}{L(M+M^q)}$ satisfies the
 239 following bound for its optimality gap with the minimum loss \mathcal{L}_* at the k^{th} iteration

$$\begin{aligned} \mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\} &\leq \frac{\bar{\alpha} L (M + M^q)}{2c} + (1 - \bar{\alpha} c)^{(k-1)} \left(\mathcal{L}(w_1) - \mathcal{L}_* - \frac{\bar{\alpha} L (M + M^q)}{2c} \right) \\ &\xrightarrow{k \rightarrow \infty} \frac{\bar{\alpha} L (M + M^q)}{2c}. \end{aligned} \tag{12}$$

240 **Proof.** See Appendix A.3.

241 **Remark 3.** Note that when $k \rightarrow \infty$, $\frac{\bar{\alpha} L (M)}{2c}$ is the original optimality gap (i.e. $\mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\}$) with
 242 *real* gradients, see Bottou et al. [13, Theorem 4.6], and then, the optimality gap is also increased by
 243 $\frac{\bar{\alpha} L (M^q)}{2c}$ due to fixed-point representation. Here we argue that by keeping the variance bound M^q
 244 relatively small, we can theoretically achieve the original performance. On the other hand, optimality
 245 gap is related to $\bar{\alpha}$, which means smaller learning rates leads to smaller optimality gap.

246 **Remark 4 (Local convexity).** Strongly convex loss is not a realistic assumption in large deep learning
 247 models. However, local convexity is a more realistic assumption i.e. \mathcal{L} is convex around the minimum
 248 point w_* . Hence, inequality (12) still holds around the minimum point w_* of a locally convex loss.
 249 We essentially train our ImageNet ResNet18 classification in the same way: first we train the network
 250 with a fixed learning rate until reaching a certain optimality gap; then the learning rate is reduced to a
 251 smaller fixed value in order to shrink the optimality gap.

252 **Empirical evidence:** Figure 3(a) demonstrates the locally convex loss landscape of ResNet18 training
 253 on CIFAR10 dataset using floating-point computation. Figure 3(b) shows the same loss landscape in

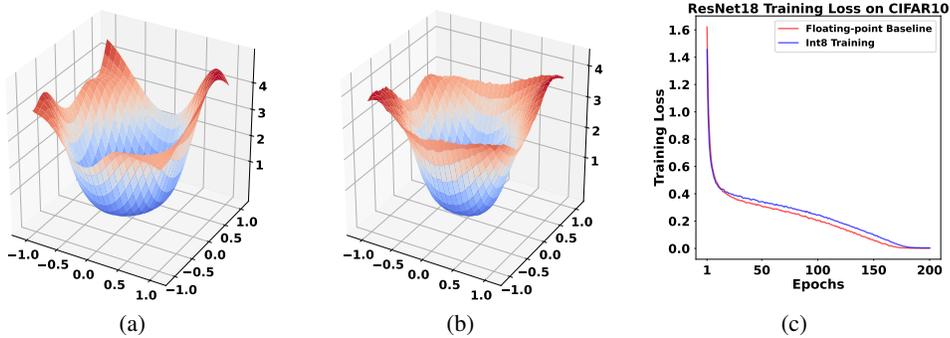


Figure 3: (a) Floating-point loss landscape, (b) Fixed-point int8 loss landscape (c) Training loss trajectory comparison. **Integer training setup:** int8 linear layer, int8 convolutional layer and int8 batch-norm layer.

254 the fixed-point int8 format. We perturbed the weights around the w_* in x and y axes using Gaussian
 255 noise and evaluated the loss \mathcal{L} in the z axis to acquire Figures 3(a) and 3(b). Comparing these two
 256 figures pronounces our assumption of local convexity in both integer and floating-point tests. Figure
 257 3(c) shows a comparison of the loss trajectory for floating-point and integer training. The fixed-point
 258 gradients are unbiased estimators of the true gradients, so the trajectory of the integer training closely
 259 follows the trajectory of its floating-point counterpart.

260 **Remark 5.** We implemented an integer weight update, hence, the computations of equation (6) is
 261 also performed in integer arithmetic. It is shown in Appendix A.4 that integer weight update with
 262 stochastic rounding is an unbiased estimator of the true weight update.

263 5 Experimental results

264 **Image classification:** Table 1 reports the experimental results of our proposed integer training method
 265 on the conventional vision classification models. In this set of models, we used int8 linear layer,
 266 int8 convolutional layer, int8 batch-norm layer, and int16 SGD to form a *fully integer training*
 267 pipeline. Note that all forward and back-propagation computations of layers are performed in *integer*
 268 format. Moreover, we have chosen smaller vision models such as ResNet18 and MobileNetV2
 269 knowing that smaller models are harder to train. Our experimental results show that our proposed
 270 integer training method is as good as floating-point with negligible loss of accuracy ($\leq 0.5\%$).

271 **Remark 6.** We compare our integer training results with official Pytorch floating-point training. For
 272 instance, in conventional vision models (Table 1), we compare the results directly with the accuracy of
 273 *torchvision* models that are reported on the official Pytorch website¹. We further emphasize that there
 274 is no change in hyper-parameters of the training. The full set of hyper-parameters of our experiments
 275 are reported in Appendix A.5. We did not compare our results with quantized back-propagation
 276 literature such as [1–4] for two reasons: (i) they often chose to use *mixed-precision* design where
 277 some parts of the model such as batch-norm and SGD are in floating-point (ii) they have changed
 278 hyper-parameters or used gradient clipping which is not required for our proposed integer training
 279 method.

280 **Vision transformer:** We validated the applicability of our proposed integer training on the original
 281 vision transformer model, notably ViT-B-16-224 [14]. We took the floating point checkpoint pre-
 282 trained on ImageNet21K². We used Huggingface [15] to fine-tune the model on CIFAR10. In this
 283 experiment, we used int8 linear layer, int8 matrix multiplication, int8 convolutional layer, and
 284 int8 layer-norm for our integer training pipeline. The result is reported in Table 1 demonstrating
 285 negligible loss of accuracy ($\leq 0.5\%$).

¹Refer to Pytorch Models: <https://pytorch.org/vision/stable/models.html>

²Refer to: <https://huggingface.co/google/vit-base-patch16-224>

Table 1: Classification

Model	Dataset	Accuracy	
		int8 (top1 - top5)	Pytorch baseline float (top1 - top5)
Conventional Vision Models			
ResNet18	CIFAR10	94.84 - N/A	95.34 - N/A
ResNet18	ImageNet	69.25 - 88.79	69.75 - 89.07
MobileNetV2	ImageNet	72.80 - 90.83	71.87 - 90.28
Vision Transformer			
ViT-B (fine-tuning)	CIFAR10	98.8 - N/A	99.1 - N/A

Table 2: Semantic Segmentation

Method	Dataset	mIOU	
		int8	baseline float
DeepLab-V1	VOC	74.73	75.0
DeepLab-V1	COCO	34.8	34.7
DeepLab-V2	VOC	77.65	77.71
DeepLab-V2	COCO	35.90	36.0

Table 3: Object Detection

Method	Dataset	mAP	
		int8	baseline float
Faster R-CNN	COCO	37.4	37.8
Faster R-CNN	VOC07+12	80.14	80.31
Faster R-CNN	Cityscapes	39.50	40.00
SSD	COCO	42.5	43.6

286 **Semantic segmentation:** To validate our proposed integer training on semantic segmentation task,
 287 we used DeepLabV1/V2³ with ResNet-101 as the backbone. We trained the models on PASCAL
 288 VOC-2012 [16] and MS COCO 10K [17]. For the PASCAL dataset, as suggested by Chen et al.
 289 [18], the models were initialized with MS COCO checkpoint and data augmentation was used during
 290 training. Likewise, for COCO dataset we used PASCAL checkpoint for initialization. The batch-norm
 291 layers are frozen in our experiments as suggested in [18]. We used linear and convolutional layers in
 292 int8 format. Table 2 shows that the mean intersect of union (mIOU) of our method closely matches
 293 the floating-point baseline.

294 **Object detection:** We used Faster R-CNN[19] with ResNet50 as its backbone and Single Shot
 295 Detector (SSD)[20] which relies on VGG-16 for its backbone. In these experiments, we used int8
 296 linear and convolutional layers, while batch-norm layers are frozen. We integrated our integer training
 297 method with MMDetection[21] toolbox and performed our experiments on MS COCO 10K[17],
 298 PASCAL VOC-2007[22], VOC-2012[16], and Cityscapes[23] datasets. As shown in Table 3, the
 299 mean average precision (mAP) of our proposed integer training method is very close to floating-point
 300 baseline.

301 6 Conclusion

302 We proposed a hardware-friendly integer training method based on coupling a linear fixed-point
 303 mapping with a non-linear inverse mapping. This method performs the representation mapping
 304 directly on the floating-point number format. Furthermore, our method uses stochastic rounding and
 305 produces unbiased estimators of gradients in the back-propagation. Thus, benefiting from unbiased
 306 gradients and effective representation mapping, the training loss trajectory closely follows its floating-
 307 point version. Moreover, there is no need to tune hyper-parameters or perform gradient clipping to
 308 correct the inappropriate gradient and weight updates. Using our proposed technique, we designed
 309 the integer version of the most vital components of deep learning such as linear layer, convolutional
 310 layer, batch-norm, layer-norm, and stochastic gradient descent (SGD). Our experimental results show
 311 the effectiveness of the proposed method, where the accuracy loss is negligible in a wide variety
 312 of the training tasks. Furthermore, in our classification tests all training components are in integer
 313 arithmetic. We theoretically studied the effect of our proposed method on SGD, and demonstrated
 314 why the optimality gap of convergence is shifted by a negligible amount.

³Refer to : <https://github.com/kazuto1011/deeplab-pytorch>

References

- 315
- 316 [1] Ron Banner, Itay Hubara, Elad Hoffer, and Daniel Soudry. Scalable methods for 8-bit training of neural
317 networks. *Advances in neural information processing systems*, 31, 2018.
- 318 [2] Xishan Zhang, Shaoli Liu, Rui Zhang, Chang Liu, Di Huang, Shiyi Zhou, Jiaming Guo, Qi Guo, Zidong
319 Du, Tian Zhi, et al. Fixed-point back-propagation training. In *Proceedings of the IEEE/CVF Conference*
320 *on Computer Vision and Pattern Recognition*, pages 2330–2338, 2020.
- 321 [3] Kang Zhao, Sida Huang, Pan Pan, Yinghan Li, Yingya Zhang, Zhenyu Gu, and Yinghui Xu. Distribution
322 adaptive int8 quantization for training cnns. In *Proceedings of the Thirty-Fifth AAAI Conference on*
323 *Artificial Intelligence*, 2021.
- 324 [4] Feng Zhu, Ruihao Gong, Fengwei Yu, Xianglong Liu, Yanfei Wang, Zhelong Li, Xiuqi Yang, and Junjie
325 Yan. Towards unified int8 training for convolutional neural network. In *Proceedings of the IEEE/CVF*
326 *Conference on Computer Vision and Pattern Recognition*, pages 1969–1979, 2020.
- 327 [5] Charbel Sakr and Naresh Shanbhag. Per-tensor fixed-point quantization of the back-propagation algorithm.
328 *arXiv preprint arXiv:1812.11732*, 2018.
- 329 [6] Qing Jin, Jian Ren, Richard Zhuang, Sumant Hanumante, Zhengang Li, Zhiyu Chen, Yanzhi Wang,
330 Kaiyuan Yang, and Sergey Tulyakov. F8net: Fixed-point 8-bit only multiplication for network quantization.
331 *arXiv preprint arXiv:2202.05239*, 2022.
- 332 [7] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Training deep neural networks with low
333 precision multiplications. arxiv 2014. *arXiv preprint arXiv:1412.7024*, 2014.
- 334 [8] Mario Drumond, Tao Lin, Martin Jaggi, and Babak Falsafi. Training dnns with hybrid block floating point.
335 *Advances in Neural Information Processing Systems*, 31, 2018.
- 336 [9] Bitar Rouhani, Daniel Lo, Ritchie Zhao, Ming Liu, Jeremy Fowers, Kalin Ovtcharov, Anna Vinogradsky,
337 Sarah Massengill, Lita Yang, Ray Bittner, et al. Pushing the limits of narrow precision inferencing at cloud
338 scale with microsoft floating point. *Advances in Neural Information Processing Systems*, 33:10271–10281,
339 2020.
- 340 [10] Dan Zuras, Mike Cowlshaw, Alex Aiken, Matthew Applegate, David Bailey, Steve Bass, Dileep Bhan-
341 darkar, Mahesh Bhat, David Bindel, Sylvie Boldo, et al. Ieee standard for floating-point arithmetic. *IEEE*
342 *Std*, 754(2008):1–70, 2008.
- 343 [11] Martin S Schmoockler and Kevin J Nowka. Leading zero anticipation and detection—a comparison of
344 methods. In *Proceedings 15th IEEE Symposium on Computer Arithmetic. ARITH-15 2001*, pages 7–12.
345 IEEE, 2001.
- 346 [12] Michael P Connolly, Nicholas J Higham, and Theo Mary. Stochastic rounding and its probabilistic
347 backward error analysis. *SIAM Journal on Scientific Computing*, 43(1):A566–A585, 2021.
- 348 [13] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning.
349 *Siam Review*, 60(2):223–311, 2018.
- 350 [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
351 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth
352 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 353 [15] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric
354 Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art
355 natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 356 [16] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The
357 PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. [http://www.pascal-](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html)
358 [network.org/challenges/VOC/voc2012/workshop/index.html](http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html), .
- 359 [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
360 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on*
361 *computer vision*, pages 740–755. Springer, 2014.
- 362 [18] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab:
363 Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.
364 *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.

- 365 [19] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection
366 with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- 367 [20] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and
368 Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages
369 21–37. Springer, 2016.
- 370 [21] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng,
371 Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu
372 Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change
373 Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint*
374 *arXiv:1906.07155*, 2019.
- 375 [22] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The
376 PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. [http://www.pascal-](http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html)
377 [network.org/challenges/VOC/voc2007/workshop/index.html](http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html), .
- 378 [23] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson,
379 Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding.
380 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223,
381 2016.

382 Checklist

- 383 1. For all authors...
- 384 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
385 tions and scope? [Yes]
- 386 (b) Did you describe the limitations of your work? [Yes] See *Theorem 1* and *Remark 3* in Section
387 4.2 reflecting the effect of representation mapping noise. Moreover, the accuracy deviation from
388 baseline floating-point training experiments are reported in Table 1, Table 2, and Table 3
- 389 (c) Did you discuss any potential negative societal impacts of your work? [N/A] The work is purely
390 methodological and related to number format.
- 391 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 392 2. If you are including theoretical results...
- 393 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 4
- 394 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.1, Appendix
395 A.2, Appendix A.3, and Appendix A.4
- 396 3. If you ran experiments...
- 397 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
398 results (either in the supplemental material or as a URL)? [No] The code is proprietary, however
399 we will provide them upon publication. We ran the experiments on publicly available datasets.
- 400 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
401 [Yes] See Appendix A.5
- 402 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
403 multiple times)? [No] . It was too costly to train multiple times, we just ran the experiments
404 once. Also note the type of training experiments we reported are resilient to random initialization
405 changes.
- 406 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
407 internal cluster, or cloud provider)? [Yes] See Appendix A.5
- 408 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 409 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 410 (b) Did you mention the license of the assets? [N/A]
- 411 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 412 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-
413 ing/curating? [N/A]
- 414 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
415 tion or offensive content? [N/A]
- 416 5. If you used crowdsourcing or conducted research with human subjects...

- 417 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
418 [N/A]
419 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
420 approvals, if applicable? [N/A]
421 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
422 participant compensation? [N/A]

423 A Appendix

424 A.1 Stochastic Rounding

425 In this paper, stochastic rounding for a variable x where $x_1 < x < x_2$ is defined as

$$\hat{x} = \begin{cases} x_1 & \text{with probability } \frac{x_2 - x}{x_2 - x_1} \\ x_2 & \text{with probability } \frac{x - x_1}{x_2 - x_1} \end{cases} \quad (13)$$

426 it is easy to see that \hat{x} is an unbiased estimator of x :

$$\mathbb{E}\{\hat{x}\} = x_1 \frac{x_2 - x}{x_2 - x_1} + x_2 \frac{x - x_1}{x_2 - x_1} = x \quad (14)$$

427 and if we relate x and \hat{x} using an error term δ (i.e. $\hat{x} = x + \delta$), then $\mathbb{E}\{\delta\} = 0$.

428 A.2 Quantization increases the gradients variance: Linear layer example

429 A linear layer is essentially a matrix multiplication. Let us denote X as inputs, W as weights and Y as the
430 output of a linear layer where $Y = XW$. Following the notation of this paper, we also denote our fixed-point
431 version of this layer as $\hat{Y} = \hat{X}\hat{W}$.

432 In our proposed integer back-propagation, the layer receives the upstream gradient $\hat{G} := \frac{\partial \hat{L}}{\partial \hat{Y}}$. Using the
433 chain-rule, the gradient with respect to weights is

$$\frac{\partial \hat{L}}{\partial \hat{W}} = \frac{\partial \hat{Y}}{\partial \hat{W}} \frac{\partial \hat{L}}{\partial \hat{Y}} = \hat{X}^\top \frac{\partial \hat{L}}{\partial \hat{Y}} = \hat{X}^\top \hat{G}. \quad (15)$$

434 Thus, computation of gradient with respect to weights is essentially a matrix multiplication too. We are interested
435 in understanding the relation between variance of $\hat{C} = \hat{X}^\top \hat{G}$ in the integer back-propagation and the true
436 gradients $C = X^\top G$. For each element \hat{C}_{ij} we have

$$\begin{aligned} \mathbb{V}\{\hat{C}_{ij}\} &= \mathbb{V}\left\{\sum_{k=1}^K \hat{X}_{ik}^\top \hat{G}_{kj}\right\} = \sum_{k=1}^K \mathbb{V}\left\{\hat{X}_{ik}^\top \hat{G}_{kj}\right\} + \sum_{\substack{k=1 \\ k \neq k'}}^K \sum_{k'=1}^K \text{COV}\left(\hat{X}_{ik}^\top \hat{G}_{kj}, \hat{X}_{ik'}^\top \hat{G}_{k'j}\right) \\ &= \sum_{k=1}^K \left\{\mathbb{E}\{(\hat{X}_{ik}^\top)^2\} \mathbb{E}\{(\hat{G}_{kj})^2\} - \mathbb{E}^2\{(\hat{X}_{ik}^\top)\} \mathbb{E}^2\{(\hat{G}_{kj})\}\right\} + \sum_{\substack{k=1 \\ k \neq k'}}^K \sum_{k'=1}^K \text{COV}\left(\hat{X}_{ik}^\top \hat{G}_{kj}, \hat{X}_{ik'}^\top \hat{G}_{k'j}\right) \\ &= \sum_{k=1}^K \left\{\mathbb{E}\{(X_{ik}^\top + \delta_{ik}^X)^2\} \mathbb{E}\{(G_{kj} + \delta_{kj}^G)^2\} - \mathbb{E}^2\{(X_{ik}^\top + \delta_{ik}^X)\} \mathbb{E}^2\{(G_{kj} + \delta_{kj}^G)\}\right\} \\ &\quad + \sum_{\substack{k=1 \\ k \neq k'}}^K \sum_{k'=1}^K \text{COV}\left(X_{ik}^\top G_{kj}, X_{ik'}^\top G_{k'j}\right) \\ &\leq \sum_{k=1}^K \left\{\mathbb{V}\{X_{ik}^\top G_{kj}\} + \sigma_X^2 \mathbb{E}\{G_{kj}\} + \sigma_G^2 \mathbb{E}\{X_{ik}^\top\} + \sigma_X^2 \sigma_G^2\right\} + \sum_{\substack{k=1 \\ k \neq k'}}^K \sum_{k'=1}^K \text{COV}\left(X_{ik}^\top G_{kj}, X_{ik'}^\top G_{k'j}\right) \\ &= \mathbb{V}\left\{\sum_{k=1}^K X_{ik}^\top G_{kj}\right\} + \sigma_G^2 \mathbb{E}\{\|X_i^\top\|_2^2\} + \sigma_X^2 \mathbb{E}\{\|G_j^\top\|_2^2\} + K \sigma_X^2 \sigma_G^2 \\ &= \mathbb{V}\{C_{ij}\} + \sigma_G^2 \mathbb{E}\{\|X_i^\top\|_2^2\} + \sigma_X^2 \mathbb{E}\{\|G_j^\top\|_2^2\} + K \sigma_X^2 \sigma_G^2. \end{aligned} \quad (16)$$

437 Note that in inequality (16), $\sigma_G^2 = \max(\mathbb{V}\{\delta_{i,j}^G\})$ knowing that error terms $\delta_{i,j}^G$ have essentially similar
438 distributions. Likewise, $\sigma_X^2 = \max(\mathbb{V}\{\delta_{i,j}^X\})$. Also note that X and G are matrices of random variables in the
439 back-propagation and $\|X_{i,\cdot}^\top\|_2^2$ denotes the norm-2 of the i^{th} row of X^\top and $\|G_{\cdot,j}^\top\|_2^2$ denotes the norm-2 of the
440 j^{th} column of G . Also, \hat{X}_{ik}^\top denotes ik^{th} element of the matrix \hat{X}^\top .

441 By having the following definitions

$$\begin{cases} M^q := \sigma_G^2 \mathbb{E}\{\|X_{i,\cdot}^\top\|_2^2\} + K \sigma_X^2 \sigma_G^2 \\ M_V^q := \sigma_X^2 \end{cases} \quad (17)$$

442 we can re-organize the inequality (16) as

$$\begin{aligned} \mathbb{V}\{\hat{C}_{ij}\} &\leq \mathbb{V}\{C_{ij}\} + M_V^q \mathbb{E}\{\|G_{\cdot,j}^\top\|_2^2\} + M^q \\ \xrightarrow{\text{Assumption 2.(iii.a)}} \mathbb{V}\{\hat{C}_{ij}\} &\leq (M_V + M_V^q) \mathbb{E}\{\|G_{\cdot,j}^\top\|_2^2\} + (M + M^q). \end{aligned} \quad (18)$$

Remark. Inequality (18) supports our *Assumption 2 (iii,b)* i.e.

$$\mathbb{V}_{\xi_k} \{\hat{g}(w_k, \xi_k)\} \leq M + M^q + (M_V + M_V^q) \|\nabla \mathcal{L}(w_k)\|_2^2$$

443 for considering the effect of our representation mapping method on gradients variance.

444 A.3 Proof of Theorem 1

445 The proof goes along the proof of Bottou et al. [13, Theorem 4.6] in the case that the gradient variance bound
446 increased as stated in *Assumption 2 (iii, b)*.

447 A convex function satisfying the inequality (11) given $\bar{w}, w \in \mathbb{R}^d$ represents a quadratic model

$$q(\bar{w}) := \mathcal{L}(w) + \nabla \mathcal{L}(w)^\top (\bar{w} - w) + \frac{1}{2} c \|\bar{w} - w\|_2^2, \quad (19)$$

448 and has a unique minimizer at w_*

$$\begin{aligned} w_* &:= w - \frac{1}{c} \nabla \mathcal{L}(w) \\ q(w_*) &= \mathcal{L}(w) - \frac{1}{2c} \|\nabla \mathcal{L}(w)\|_2^2 \\ \rightarrow 2c(\mathcal{L}(w) - \mathcal{L}_*) &\leq \|\nabla \mathcal{L}(w)\|_2^2; \quad \forall w \in \mathbb{R}^d. \end{aligned} \quad (20)$$

449 Also remember that the fixed learning rate in our integer back-propagation has the following constraint

$$0 < \bar{\alpha} \leq \frac{1}{L(M_G + M_G^q)}. \quad (21)$$

450 Starting from inequality (10) and using inequalities (20) and (21) we can write

$$\begin{aligned} \mathbb{E}_{\xi_k} \{\mathcal{L}(w_{k+1})\} - \mathcal{L}(w_k) &\leq -\left(1 - \frac{1}{2} \bar{\alpha} L(M_G + M_G^q)\right) \bar{\alpha} \|\nabla \mathcal{L}(w_k)\|_2^2 + \frac{1}{2} \bar{\alpha}^2 L(M + M^q) \\ &\leq -\frac{1}{2} \bar{\alpha} \|\nabla \mathcal{L}(w_k)\|_2^2 + \frac{1}{2} \bar{\alpha}^2 L(M + M^q) \\ &\leq -\bar{\alpha} c (\mathcal{L}(w_k) - \mathcal{L}_*) + \frac{1}{2} \bar{\alpha}^2 L(M + M^q). \end{aligned} \quad (22)$$

451 By subtracting \mathcal{L}_* , rearrange, and taking total expectation from both sides of inequality (22) we have

$$\mathbb{E}\{\mathcal{L}(w_{k+1}) - \mathcal{L}_*\} \leq (1 - \bar{\alpha} c) \mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\} + \frac{1}{2} \bar{\alpha}^2 L(M + M^q). \quad (23)$$

452 Then by subtracting $\frac{\bar{\alpha} L(M + M^q)}{2c}$ from both sides

$$\begin{aligned} \mathbb{E}\{\mathcal{L}(w_{k+1}) - \mathcal{L}_*\} - \frac{\bar{\alpha} L(M + M^q)}{2c} &\leq (1 - \bar{\alpha} c) \mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\} + \frac{1}{2} \bar{\alpha}^2 L(M + M^q) - \frac{\bar{\alpha} L(M + M^q)}{2c} \\ &= (1 - \bar{\alpha} c) \left(\mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\} - \frac{\bar{\alpha} L(M + M^q)}{2c} \right). \end{aligned} \quad (24)$$

453 Thus, *Theorem 1* can be proven by applying inequality (24) repeatedly for $k \in \mathbb{N}$. Also note that using inequality
 454 (21) it is easy to derive that $0 < (1 - \bar{\alpha}c) < 1$ because

$$0 < \bar{\alpha}c \leq \frac{c}{L(M_G + M_G^q)} \leq \frac{c}{L} \leq 1, \quad (25)$$

455 hence, in *Theorem 1*, if $k \rightarrow \infty$, then $(1 - \bar{\alpha}c)^k \rightarrow 0$ and the optimality gap of integer training algorithm using
 456 inequality (12) is

$$\mathbb{E}\{\mathcal{L}(w_k) - \mathcal{L}_*\} \leq \frac{\bar{\alpha}L(M + M^q)}{2c} \quad (26)$$

s.t. $k \rightarrow \infty$.

457 A.4 Integer weight update

458 In the integer weight update, the equation (6) transforms to its fixed-point version with integer-only arithmetic

$$\hat{w}_{k+1} = \hat{w}_k + \hat{\alpha}_k \hat{g}(w_k, \xi_k). \quad (27)$$

459 By expanding the error terms for the representation mapping and taking expectation on both sides we can see the
 460 weights are on the average updated equivalent to the true weights.

$$\begin{aligned} \mathbb{E}\{\hat{w}_{k+1}\} &= \mathbb{E}\{\hat{w}_k + \hat{\alpha}_k \hat{g}(w_k, \xi_k)\} \\ &= \mathbb{E}\{w_k + \delta_k^w + (\alpha_k + \delta^\alpha)(g(w_k, \xi_k) + \delta^g)\} \\ &= w_k + \alpha_k g(w_k, \xi_k) \\ &= w_{k+1} \end{aligned} \quad (28)$$

461 We used stochastic rounding for our weight update operation, thus $\mathbb{E}\{\delta\} = 0$ and using our proposed method, \hat{g}
 462 is unbiased estimator of g .

463 A.5 Experimental setup

464 **Computing resources:** We ran our experiments using an in-house developed integer emulator to avoid common
 465 floating-point quantization techniques. In our emulator, we perform the *representation mapping* within the GPU
 466 memory. We further developed the integer deep learning modules using Pytorch autograd functionality on top of
 467 our integer emulator. Experimental results of this paper are run using the following number of GPUs.

- 468 • ResNet18 on ImageNet requires $4 \times V100$ GPUs when batch size is 512 and $2 \times V100$ GPUs when
 469 batch size is 256.
- 470 • MobileNetV2 on ImageNet requires $8 \times V100$ GPUs when batch size is 512.
- 471 • ResNet18 on CIFAR10 runs on $1 \times V100$ GPUs when batch size is 128.
- 472 • Semantic segmentation experiments require $2 \times V100$ GPUs.
- 473 • Object detection experiments require $8 \times V100$ GPUs.
- 474 • ViT-B fine-tuning experiment requires $8 \times V100$ GPUs.

475 **Classification:** The hyper-parameters of our classification experiments are reported in Table 4. We used SGD
 476 with momentum of 0.9 for conventional classification experiments and AdamW for ViT fine-tuning.

Table 4: Hyper-parameters for classification experiments

Dataset	Model	Training epochs	Learning rate	LR scheduling	Weight decay	Batch size
ImageNet	ResNet18	90 or 100	0.1	$\times 0.1$ every 30 epochs	1e-4	512
	MobileNetV2	70	0.1	Cosine with $T_{max} = 70$	4e-5	512
CIFAR10	ResNet18	200	0.1	Cosine with $T_{max} = 90$	5e-4	128
	ViT-B fine-tuning	100	5e-5	Cosine with $T_{max} = 100$	0.01	512

477 **Semantic segmentation:** The hyper-parameters for the semantic segmentation experiments are provided in
 478 Table5. For DeepLabV1 and DeepLabV2 one-scale and multi-scale loss was used respectively. Conditional
 479 random field post-processing was used at the network’s output as proposed in Chen et al. [18]. In all the
 480 experiments we use SGD with momentum of 0.9 and weight decay of 5×10^{-4} .

Table 5: Hyper-parameters for semantic segmentation experiments

Model	Dataset	Training iterations	Learning rate	Batch size	Data Augmentation	CRF
DeepLabV1/V2	VOC	20000	2.5e-4	16	✓	✓
	COCO	30000	2.5e-4	16	-	✓

Table 6: Hyper-parameters for object detection experiments

Model	Dataset	Training epochs	Learning rate	Per GPU Batch size	LR reduction epochs	LR Warmup
Faster R-CNN	COCO	12	0.2	1	8 and 11	✓
	VOC	12	0.1	2	9	-
	Cityscapes	64	0.1	1	56	✓
SSD	COCO	24	0.2	1	16 and 22	✓

481 **Object detection:** The hyper-parameters are provided in Table 6. All experiments use an SGD optimizer with a
482 momentum of 0.9 and a weight decay parameter of 10^{-5} . The experiments indicated with *LR Warmup* use a
483 linear warm-up function with a ratio of 10^{-3} for the first 500 iterations. For the rest of training, the learning rate
484 is reduced by 0.1 at the epochs indicated in the *LR Reduction Epochs* column.