

# WORKING DRAFT

## Bit Centering: A Technical Report

Christopher R. Aberger<sup>2</sup>, Christopher De Sa<sup>1</sup>, and Christopher Ré<sup>2</sup>

<sup>1</sup> Cornell University, Ithaca, NY  
<sup>2</sup> Stanford University, Stanford, CA

July 7, 2018

### Abstract

The motivation for this work is the rise of accelerators for analytics, in which a host machine with large memory but limited compute uses a high-performance accelerator that is capable of performing many low-precision operations cheaply. This technical report describes techniques to cope with low-precision inference and learning based on a novel *bit centering technique*. The key idea is to gracefully tradeoff representation using host-memory and an accelerator. Applications to learning are described including novel convergence rates. This allows one to train high-precision models by off-loading as much as possible to the low-precision accelerator. **Please note this draft is a working draft and may be revised without notice.**

## 1 Introduction

The technical setting for this technical report note is the recent resurgence of hardware accelerators for machine learning and other applications.<sup>1</sup> In such a scenario, there is a host device with a larger memory and general purpose processor (e.g., a CPU). The accelerator typically has many more computational units of lower precision and smaller memory. In turn, this motivates the use of low-precision arithmetic. We think of the accelerator as working on the inner loop of the computation and the host computation as infrequently (orders of magnitude less frequently) performing some high-precision computation to help direct (or center) the low-precision computation on the accelerator. This technical report describes a technique that uses limited amounts of host memory and computation to improve the low-precision computation of an accelerator.

As we will see, the conventional wisdom is that training algorithms using low-precision iterations are limited to low accuracy. Our goal is to formalize situations in which a modest amount of host memory allows us to achieve dramatically better error with minimal overhead.

**The Core Idea: Bit Centering** Bit centering is a technique in which we represent each number as the sum of a pre-computed full-precision *offset* and a low-precision *delta*:

$$x = o_x + \delta_x.$$

The benefit of bit centering, as compared with simply storing the number  $x$  itself in low-precision, comes from the fact that floating point computation has *relative* error. That is, if we do some computation in low precision that is intended to approximate some number  $x \in \mathbb{R}$ , and the result of that computation is a floating point number  $\tilde{x}$ , then

$$|\tilde{x} - x| \leq \epsilon|x|,$$

---

<sup>1</sup>The authors have a financial interest in SambaNova Systems, which produces machine learning and big data platform including hardware accelerators and an accelerated software stack.

where  $\epsilon$  is some fixed-but-small number that is a function of the number of bits used. In comparison, if we compute only the delta in low-precision, using bit centering, the error will look like

$$|\tilde{x} - x| = |\tilde{\delta}_x - \delta_x| \leq \epsilon |\delta_x| = \epsilon |x - o_x|.$$

So if we can choose  $o_x$  so that it will be close to  $x$  (equivalently, so that  $\delta_x$  will be small), then bit centering can produce a much more accurate estimate of  $x$  using the same number of bits.

In this document, we will describe bit centering, and compare it with other quantization strategies, when used to compute gradients for learning algorithms. The quantization strategies are:

- **High-precision.** Compute the whole algorithm in high-precision. This is the null quantization strategy.
- **Quantize last.** Compute the gradient in high-precision floating point, then quantize to low-precision floating point at the end. This is sometimes called *gradient quantization*.
- **Quantize first.** Quantize all the inputs to the gradient procedure into low-precision floating point, and compute the gradient in low-precision. This quantizes all operations, but can cause accuracy to suffer.
- **Bit centering.** Compute the gradient using bit-centered numbers with a pre-computed-on-the-host high-precision offset and a computed-on-the-accelerator low-precision floating-point delta.
- **Block floating point bit centering.** The same as bit centering, but using block-floating-point numbers (which share an exponent) for the delta rather than floating-point numbers. This lets us save memory and compute that would otherwise be used to store and process each number's exponent individually.

## 1.1 Notation

We will use  $\tilde{x}$  to denote an approximation. We will follow the notation of Trefethen and Bau III [3] and their error model. Namely, given a number  $x$  and a floating point number system with machine error  $\epsilon_{\text{machine}}$ , then if  $\text{fl}(x)$  is the floating point representation of  $x$ , there will exist an  $\epsilon$  such that  $|\epsilon| \leq \epsilon_{\text{machine}}$  and:

$$\text{fl}(x) = (1 + \epsilon)x$$

Moreover, if  $*$  denotes any basic operation in scalar arithmetic (but also fused-multiply and add on two floating point representations), and  $\otimes$  denotes its floating point analogue, then for some  $\epsilon$  such that  $|\epsilon| \leq \epsilon_{\text{machine}}$

$$x \otimes y = (1 + \epsilon)(x * y).$$

This is sometimes called the fundamental theorem of floating point arithmetic. In this note, we will deal with just two floating point number systems to make our life easier: (1) a **low precision** number system with machine precision  $\epsilon_{\text{machine-lo}}$  and (2) a **high precision** number system with machine precision  $\epsilon_{\text{machine-hi}}$ . We let the low-precision representation of a number  $x$  be denoted  $\text{fl}_{\text{lo}}(x)$  and the high-precision representation be denoted  $\text{fl}_{\text{hi}}(x)$ . In the analysis that follows, we will use  $\epsilon_{\text{lo}}$  to generally denote a value that satisfies  $\epsilon_{\text{lo}} \leq \epsilon_{\text{machine-lo}}$  and  $\epsilon_{\text{hi}}$  to generally denote a value that satisfies  $\epsilon_{\text{hi}} \leq \epsilon_{\text{machine-hi}}$ . We also define the precision ration  $k$  as

$$k = \frac{\epsilon_{\text{machine-lo}}}{\epsilon_{\text{machine-hi}}}.$$

Intuitively, we can imagine that the high-precision representation is a 32-bit floating point number, and the low-precision representation is an 8-bit float.

## 2 Case Study: Dot Product

To better understand the error improvement tradeoff of splitting versus full precision and low precision in host memory, we consider computing a dot product  $x^T y$  for  $x, y \in \mathbb{R}^n$ .

To simplify the discussion, we will consider bit recentering on  $x$  to start. That is, we write  $x = o_x + \delta_x$ . The idea is that  $o_x$  will be stored in high precision on the host and we can do limited operations on it, while  $\delta_x$  is stored in low precision on the accelerator.

Consider computing  $x^T y$  where  $x = o_x + \delta_x$  with  $x$  and  $y$  in  $\mathbb{R}_+^n$  (both stored in high precision) and storing the result in low precision (i.e. the floating point representation with machine error  $\varepsilon_{\text{machine-lo}}$ ). We are going to evaluate three methods of computing this expression: *quantize last*, which computes the dot product in high precision and quantizes at the end; *quantize first*, which quantizes  $x$  and  $y$  into low precision and then computes the dot product in low precision; and *bit centering*. Our results in this section are summarized in the following table, which also compares to *high-precision*, a strategy which just computes the dot product in high precision and stores the result in high precision as well.

In the derivations that follow, we will use the well-known fact [1] that for a nonnegative floating-point dot product with machine error  $\varepsilon_{\text{machine}}$ , the result of the computation will be, for some  $|\varepsilon| \leq \varepsilon_{\text{machine}}$ ,

$$\text{fl}(x)^T \text{fl}(y) = x^T y + |x|^T |y| n \varepsilon$$

where  $|x|$  denotes the absolute value of  $x$  applies entrywise, and where we are ignoring terms proportional to  $\varepsilon^2$ . (Note that this is for using naive summation; if we used pairwise summation the  $n$  would be replaced by a  $\log n$  factor.) For non-negative  $x$  and  $y$ , this is equivalent to

$$\text{fl}(x)^T \text{fl}(y) = x^T y (1 + n \varepsilon),$$

a fact we use in the subsequent derivation.

Quantization strategy	# LP Ops	# HP Ops	Result
High-precision	—	$n$ multiplies $n - 1$ adds	$(x^T y) (1 + \varepsilon_{\text{hi}} \cdot n)$
Quantize last	—	$n$ multiplies $n - 1$ adds	$(x^T y) (1 + \varepsilon_{\text{hi}} \cdot (n + k))$
Quantize first	$n$ multiplies $n - 1$ adds	—	$(x^T y) (1 + \varepsilon_{\text{hi}} \cdot kn)$
Bit centering	$n$ multiplies $n$ adds	$n$ multiplies $n - 1$ adds precomputed on host	$(x^T y) (1 + \varepsilon_{\text{hi}} \cdot (n + k))$ $+ (\delta_x^T y) (\varepsilon_{\text{hi}} \cdot (kn + n + 2k + 1))$ $+ (x^T y) \cdot \varepsilon_{\text{hi}} \cdot k$

**Quantize last** We compute the dot product first in high precision, producing an intermediate value which will satisfy

$$\text{fl}_{\text{hi}}(x)^T \text{fl}_{\text{hi}}(y) = x^T y \cdot (1 + \varepsilon_{\text{hi}} \cdot n).$$

If we then quantize that high-precision intermediate value to low precision, we will get

$$\begin{aligned} \text{fl}_{\text{lo}}(\text{fl}_{\text{hi}}(x)^T \text{fl}_{\text{hi}}(y)) &= (x^T y) (1 + \varepsilon_{\text{hi}} \cdot n) (1 + \varepsilon_{\text{lo}}) \\ &= (x^T y) \left( 1 + \varepsilon_{\text{hi}} \cdot \left( n + \frac{\varepsilon_{\text{machine-lo}}}{\varepsilon_{\text{machine-hi}}} \right) \right) \\ &= (x^T y) (1 + \varepsilon_{\text{hi}} \cdot (n + k)). \end{aligned}$$

Note again that we drop terms that are  $O(\varepsilon^2)$ , as is standard in Trefethen and Bau III [3]. While this is a small relative error, all of the computation we did here was in high precision. Can we do better?

**Quantize first** If we first quantize to low precision and then do the multiplication, then

$$\text{fl}_{\text{lo}}(x)^T \text{fl}_{\text{lo}}(y) = (x^T y)(1 + \varepsilon_{\text{lo}} \cdot n) = (x^T y)(1 + \varepsilon_{\text{hi}} \cdot kn)$$

That is, we suffer the extra error *multiplicatively*, rather than additively, as we did in the quantize-last case. Can we do better? Our goal is to smoothly interpolate between the two approaches.

**Hybrid method: bit centering with an offset** Let’s compare this with bit centering, which uses an offset. Here, we decompose  $x = o_x + \delta_x$  and represent  $o_x$  in high precision and  $\delta_x$  in low-precision. So part of the dot product is computed in high precision and rounded (the offset), and part is computed in low precision (the delta). We expect that if the delta is small, this algorithm will have lower error. Let’s upper bound the error of this algorithm:

$$\begin{aligned} & \text{fl}_{\text{lo}}(\text{fl}_{\text{hi}}(o_x)^T \text{fl}_{\text{hi}}(y)) + \text{fl}_{\text{lo}}(\delta_x)^T \text{fl}_{\text{lo}}(y) \\ &= ((o_x^T y)(1 + \varepsilon_{\text{hi}} \cdot (n + k)) + (\delta_x^T y)(1 + \varepsilon_{\text{hi}} \cdot kn))(1 + \varepsilon_{\text{lo}}) \\ &= (o_x^T y)(1 + \varepsilon_{\text{hi}} \cdot (n + 2k)) + (\delta_x^T y)(1 + \varepsilon_{\text{hi}} \cdot k(n + 1)) \\ &= ((x - \delta_x)^T y)(1 + \varepsilon_{\text{hi}} \cdot (n + 2k)) + (\delta_x^T y)(1 + \varepsilon_{\text{hi}} \cdot k(n + 1)) \\ &= (x^T y)(1 + \varepsilon_{\text{hi}} \cdot (n + k)) + (\delta_x^T y)(\varepsilon_{\text{hi}} \cdot (kn + n + 2k + 1)) + (x^T y) \cdot \varepsilon_{\text{hi}} \cdot k \end{aligned}$$

This formula smoothly interpolates between the two above formulae. We’ve written this so that the first term is identical to quantize last, the second term is the quantize first error, and the third term is a small extra error from doing the extra addition to add the offset and delta components together. Since the quantize last term is the best we could hope for, we can view the second and third terms as additional error. When  $\delta_x$  is small enough this approach improves error—as we suspected. In contrast, if  $\delta_x$  is large and the offset is 0 we suffer like quantize first, and suggests that when  $\delta_x$  is small the hybrid strategy is likely to provide an improvement.

A difference of this approach from quantize last is that we pre-compute  $\text{fl}_{\text{lo}}(\text{fl}_{\text{hi}}(o_x)^T \text{fl}_{\text{hi}}(y))$  on the host and then load it into memory. In terms of accelerator compute, it costs at worst a  $2/n$  additional relative computation time over quantize first (since we need to do an extra add and an extra load from the host, compared with quantize first). This may save compute overall if we are evaluating dot products with  $y$  for multiple different values of  $\delta_x$ , but the same offset  $o_x$ .

### 3 Learning Using Bit Centering

We want to use bit centering for SGD and SGD-like algorithms. To do this, we need to compute gradients using bit centering. That is, we need to apply bit centering to compute things like  $\nabla f_i(w)$ . First, we need to decide how to assign the offset. One natural way to do this is to choose some  $o_w$  that we believe is close to  $w$ , and let  $o_w$  be the offset for storing  $w$ . That is, we store  $w$  using bit centering as

$$w = o_w + \delta_w.$$

If we do this, the natural way to assign the offset for  $\nabla f_i(w)$  is to assign  $\nabla f_i(o_w)$ . Then the gradient, represented using bit centering, will look like

$$\nabla f_i(w) = \nabla f_i(o_w) + (\nabla f_i(o_w + \delta_w) - \nabla f_i(o_w)),$$

where the term on the left is the offset and the term on the right is the delta. To compute gradients using bit centering, we would pre-compute  $\nabla f_i(o_w)$  in high precision for some  $o_w$ , and then compute  $\nabla f_i(o_w + \delta_w) - \nabla f_i(o_w)$  using low-precision arithmetic. For bit centering to be effective we need to compute this in such a way that our approximation satisfies a good relative error bound similar to the one we had for the dot product example. However, if we just compute this naively using low-precision arithmetic by computing  $\nabla f_i(o_w + \delta_w)$  and  $\nabla f_i(o_w)$  and then subtracting them, we will observe *catastrophic cancellation* resulting in a loss of precision for the delta component, because  $\nabla f_i(o_w + \delta_w)$  and  $\nabla f_i(o_w)$  are very close to each other and cancel out to a value near 0 when  $\delta$  is small. In order to prevent this

loss of precision, we need to compute  $\nabla f_i(o_w + \delta_w) - \nabla f_i(o_w)$  in another way, one that does not involve computing the terms independently and subtracting them. This “other way” must depend on the computational structure of the functions  $f_i$ . Next, we will present ways that we can do this, first for simple linear-model objectives, and then for more general objectives, such as deep neural nets.

### 3.1 Gradient Estimation for Linear Models

A regularized linear model is a model with loss of the form

$$f(w) = \frac{1}{n} \sum_{i=1}^n l_i(w^T x_i) + \frac{1}{2} \sigma \|w\|^2$$

where  $x_i \in \mathbb{R}^d$  is the training example,  $w \in \mathbb{R}^d$  are the weights,  $l_i : \mathbb{R} \rightarrow \mathbb{R}$  is the component loss function,  $d$  is the number of dimensions, and  $\sigma$  is the regularization parameter. Examples of linear model problems include linear regression and logistic regression. For this class of problems, the individual component functions can be written simply as

$$f_i(w) = l_i(w^T x_i) + \frac{1}{2} \sigma \|w\|^2,$$

and the gradient of this is

$$g = f'_i(w) = l'_i(w^T x_i) x_i + \sigma w.$$

It is straightforward to imagine computing this using the quantize-first and quantize-last strategies, by simply doing all the computations in low-precision or high-precision, respectively. To be concrete, in the quantize-last strategy, we do

$$\tilde{g} = \text{fl}_{\text{lo}}(l'_i(\text{fl}_{\text{hi}}(w)^T \text{fl}_{\text{hi}}(x_i)) \text{fl}_{\text{hi}}(x_i) + \text{fl}_{\text{hi}}(\sigma) \text{fl}_{\text{hi}}(w))$$

and in the quantize-first strategy, we do

$$\tilde{g} = l'_i(\text{fl}_{\text{lo}}(w)^T \text{fl}_{\text{lo}}(x_i)) \text{fl}_{\text{lo}}(x_i) + \text{fl}_{\text{lo}}(\sigma) \text{fl}_{\text{lo}}(w).$$

On the other hand, if we compute this gradient using bit centering as described above for  $w = o_w + \delta_w$ , the offset will be

$$o_g = l'_i(o_w^T x_i) x_i + \sigma o_w$$

and the delta will be

$$\begin{aligned} \delta_g &= (l'_i(w^T x_i) x_i + \sigma w) - (l'_i(o_w^T x_i) x_i + \sigma o_w) \\ &= (l'_i(o_w^T x_i + \delta_w^T x_i) - l'_i(o_w^T x_i)) x_i + \sigma \delta_w. \end{aligned}$$

We can compute this in low-precision without loss of precision by using the following algorithm. To do this, we make the following assumptions. First, we assume that our floating point computation is accurate enough to compute  $l'_i(u)$  for any  $u$  up to machine precision. That is, for any  $u$ , and for either high or low precision,

$$l'_i(\text{fl}(u)) = (1 + \varepsilon) l'_i(u).$$

Second, we assume that we have the ability to approximate scalar functions accurately using low-precision arithmetic. Given any “sufficiently nice” scalar function  $h : \mathbb{R} \rightarrow \mathbb{R}$  with  $h(0) = 0$  and any range  $[-a, a]$  about zero, we suppose that we can (via some amount of high-precision computation) produce a low-precision-computable routine  $\tilde{h}$  that can compute  $h(x)$  for any low-precision number  $x \in [-a, a]$  to within relative accuracy

$$\left| \tilde{h}(x) - h(x) \right| \leq \zeta \cdot |x| \cdot \max_{u \in [-a, a]} \left| \frac{h(u)}{u} \right| \cdot \varepsilon_{\text{machine-lo}}$$

for some constant  $\zeta \geq 1$ . Methods for doing this include low-precision Chebyshev approximation and linear interpolation: any desired method may be used, as is appropriate for the accelerator hardware. We further assume that  $l'_i$  is part of the class of functions that are sufficiently nice for this approximation method to work. Third, we assume that

$l'_i$  is Lipschitz continuous with parameter  $L$ . Finally, we assume that we know some bound on  $\delta_w$ , and can guarantee that  $\|\delta_w\| \leq R$ . Now the algorithm.

---

**Algorithm 1** Bit Centering for Linear Models

---

- 1: **given:** loss function  $l$ , training example  $x$ , and delta maximum range  $R$
  
  - 2: **high-precision precomputation on host:**
  - 3: **given:** offset  $o_w$
  - 4:  $\phi \leftarrow \text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)$
  - 5:  $a \leftarrow R \cdot \|\text{fl}_{\text{hi}}(x)\|$
  - 6: **define:**  $h(\psi) = l'(\phi + \psi) - l'(\phi)$
  - 7: **compute LP function approximation:**  $\tilde{h} \approx h$  in range  $[-a, a]$
  - 8: **store:**  $\tilde{h}$  and  $\text{fl}_{\text{lo}}(x)$  and  $\text{fl}_{\text{lo}}(\sigma)$
  - 9: **output:**  $o_g = l'(\phi) \cdot \text{fl}_{\text{hi}}(x) + \text{fl}_{\text{hi}}(\sigma) \cdot \text{fl}_{\text{hi}}(o_w)$
  
  - 10: **low-precision computation on accelerator:**
  - 11: **given:** delta  $\delta_w$
  - 12: **load:**  $\tilde{h}$  and  $\text{fl}_{\text{lo}}(x)$  and  $\text{fl}_{\text{lo}}(\sigma)$
  - 13:  $\psi \leftarrow \text{fl}_{\text{lo}}(\delta_w)^T \text{fl}_{\text{lo}}(x)$
  - 14: **output:**  $\delta_g = \tilde{h}(\psi) \cdot \text{fl}_{\text{lo}}(x) + \text{fl}_{\text{lo}}(\sigma) \cdot \text{fl}_{\text{lo}}(\delta_w)$
- 

To compute each of these low-precision gradients requires  $3d$  low-precision multiplies on the accelerator ( $d$  to compute  $\text{fl}_{\text{lo}}(\delta_w)^T \text{fl}_{\text{lo}}(x)$ ,  $d$  to compute  $\tilde{h}(\psi) \cdot \text{fl}_{\text{lo}}(x)$ , and  $d$  to compute  $\text{fl}_{\text{lo}}(\sigma) \cdot \text{fl}_{\text{lo}}(\delta_w)$ ), plus an additional  $O(1)$  low-precision compute to compute the approximation  $\tilde{h}(\psi)$ . It also requires  $2d$  low-precision loads (for  $\delta_w$  and  $\text{fl}_{\text{lo}}(x)$ ) plus an addition  $O(1)$  low-precision loads for  $h_i$  and  $\text{fl}_{\text{lo}}(\sigma)$ . Thus, the total on-accelerator cost of bit centering for linear models is  $3d + O(1)$  low-precision multiplies plus  $2d + O(1)$  low-precision loads. This can be compared with the cost of the quantize-last strategy, which would require  $3d$  high-precision multiplies plus  $2d + 1$  high-precision loads. Thus, through the addition of only  $O(1)$  extra low-precision loads and compute, compared with quantize-last, we have converted all the computation to low-precision. Note that this comes at the cost of some high-precision computation, which may be precomputed and done on the host as it only depends on  $o_w$ . The amount of this pre-computation is also  $4d + O(1)$  high-precision multiplies plus  $2d + 1$  high-precision loads and  $d + O(1)$  low-precision stores.

Next, we will analyze the errors produced by this method. But first, we summarize our results in a table.

Quant. strat.	# LP Ops	# HP Ops	Result error ( $\ell_2$ norm)
High-precision	—	$2d$ loads $3d$ multiplies $d$ adds 1 compute of $l'$	$\left(2 \cdot  l'(w^T x)  \cdot \ x\  + \ w\  \cdot \ x\ ^2 \cdot Ld +  \sigma  \ w\ \right) \cdot \varepsilon_{\text{hi}}$
Quantize last	—	$2d$ loads $3d$ multiplies $d$ adds 1 compute of $l'$	$\left((2+k) \cdot  l'(w^T x)  \cdot \ x\  + \ w\  \cdot \ x\ ^2 \cdot Ld + (1+k) \cdot  \sigma  \ w\ \right) \cdot \varepsilon_{\text{hi}}$
Quantize first	$3d$ multiplies $d$ adds 1 compute of $l'$	$2d$ loads	$\left(2 l'(w^T x)  \cdot \ x\  + \ w\  \cdot \ x\ ^2 \cdot Ld +  \sigma  \ w\ \right) \cdot k \cdot \varepsilon_{\text{hi}}$
Bit centering	$2d + O(1)$ loads $3d$ multiplies $d$ adds 1 eval of $\tilde{h}$	$2d$ loads $4d$ multiplies $2d - 1$ adds 1 approx of $\tilde{h}$ $2d + O(1)$ lp stores precomputed on host	error for $\delta_g$ : $\left(2L \cdot \ x\ ^2 \cdot \zeta + L \cdot \ x\ ^2 \cdot d \cdot k +  \sigma  \cdot k\right) \cdot \ \delta_w\  \cdot \varepsilon_{\text{hi}}$

**Error analysis.** The error analysis for  $o_g$  is straightforward. From the algorithm statement, we have that

$$o_g = l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \cdot \text{fl}_{\text{hi}}(x) + \text{fl}_{\text{hi}}(\sigma) \cdot \text{fl}_{\text{hi}}(o_w).$$

First, notice that by the dot product error analysis,

$$|\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x) - o_w^T x| \leq |o_w|^T |x| \cdot d \cdot \varepsilon_{\text{machine-hi}}.$$

This means that if  $l'$  is Lipschitz continuous, then

$$|l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) - l'(o_w^T x)| \leq |o_w|^T |x| \cdot Ld \cdot \varepsilon_{\text{machine-hi}} + |l'(o_w^T x)| \cdot \varepsilon_{\text{machine-hi}}$$

where  $L$  is the Lipschitz constant of the function  $l'$ , and the last term results from our assumption that we can compute the function  $l'$  to machine precision. It also follows that, if  $x_j$  is the  $j$ th component of the vector  $x$ ,

$$\begin{aligned}
& |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \text{fl}_{\text{hi}}(x_j) - l'(o_w^T x) x_j| \\
& \leq |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \text{fl}_{\text{hi}}(x_j) - l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) x_j| + |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) x_j - l'(o_w^T x) x_j| \\
& \leq |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \text{fl}_{\text{hi}}(x_j) - l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) x_j| \cdot \varepsilon_{\text{machine-hi}} + |o_w|^T |x| \cdot |x_j| \cdot Ld \cdot \varepsilon_{\text{machine-hi}} + |l'(o_w^T x)| \cdot |x_j| \cdot \varepsilon_{\text{machine-hi}} \\
& \leq 2|l'(o_w^T x) x_j| \cdot \varepsilon_{\text{machine-hi}} + |o_w|^T |x| \cdot |x_j| \cdot Ld \cdot \varepsilon_{\text{machine-hi}}
\end{aligned}$$

where in the last line we implicitly dropped terms in  $\varepsilon^2$ . It follows that the error for the  $j$ th component of  $o_g$ , which we denote  $o_{g,j}$ , is

$$\begin{aligned}
|o_{g,j} - (l'(o_w^T x) x_j + \sigma o_{w,j})| &= |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \text{fl}_{\text{hi}}(x_j) + \text{fl}_{\text{hi}}(\sigma) \text{fl}_{\text{hi}}(o_{w,j}) - (l'(o_w^T x) x_j + \sigma o_{w,j})| \\
&\leq |l'(\text{fl}_{\text{hi}}(o_w)^T \text{fl}_{\text{hi}}(x)) \text{fl}_{\text{hi}}(x_j) - l'(o_w^T x) x_j| + |\text{fl}_{\text{hi}}(\sigma) \text{fl}_{\text{hi}}(o_{w,j}) - \sigma o_{w,j}| \\
&\leq 2|l'(o_w^T x) x_j| \cdot \varepsilon_{\text{machine-hi}} + |o_w|^T |x| \cdot |x_j| \cdot Ld \cdot \varepsilon_{\text{machine-hi}} + |\sigma o_{w,j}| \varepsilon_{\text{machine-hi}}.
\end{aligned}$$

Or, in terms of infinity norms, the error will be

$$\|o_g - (l'(o_w^T x) x + \sigma o_w)\|_\infty \leq \left(2|l'(o_w^T x)| \|x\|_\infty + |o_w|^T |x| \cdot \|x\|_\infty \cdot Ld + |\sigma| \|o_w\|_\infty\right) \cdot \varepsilon_{\text{machine-hi}}$$

Note that basically the same analysis would apply to give us a bound on the error of the quantize-first and quantize-last strategies. Specifically, for quantize-first, everything will be the same, except in low-precision, so we get

$$\begin{aligned} \|g_{\text{quantize-first}} - (l'(o_w^T x) x + \sigma w)\|_\infty &\leq \left(2|l'(w^T x)|\|x\|_\infty + |w|^T|x| \cdot \|x\|_\infty \cdot Ld + |\sigma|\|w\|_\infty\right) \cdot \varepsilon_{\text{machine-lo}} \\ &\leq \left(2|l'(w^T x)|\|x\|_\infty + |w|^T|x| \cdot \|x\|_\infty \cdot Ld + |\sigma|\|w\|_\infty\right) \cdot k \cdot \varepsilon_{\text{machine-hi}}. \end{aligned}$$

For quantize last, everything will again be the same, except that we will get an extra error term from quantizing into low-precision at the end, so

$$\begin{aligned} \|g_{\text{quantize-last}} - (l'(w^T x) x + \sigma w)\|_\infty &\leq \left(2|l'(w^T x)|\|x\|_\infty + |w|^T|x| \cdot \|x\|_\infty \cdot Ld + |\sigma|\|w\|_\infty + k \cdot \|l'(w^T x) x + \sigma w\|_\infty\right) \cdot \varepsilon_{\text{machine-hi}} \\ &\leq \left((2+k) \cdot |l'(w^T x)|\|x\|_\infty + |w|^T|x| \cdot \|x\|_\infty \cdot Ld + (1+k) \cdot |\sigma|\|w\|_\infty\right) \cdot \varepsilon_{\text{machine-hi}}. \end{aligned}$$

The analysis for the low-precision section is similar. First, notice that

$$|\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x) - \delta_w^T x| \leq |\delta_w|^T|x| \cdot d \cdot \varepsilon_{\text{machine-lo}}.$$

It follows that, since  $h(\psi) = l'(\phi + \psi) - l'(\phi)$ ,  $h$  has the same Lipschitz constant as  $l'$ , and so

$$\begin{aligned} |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) - h(\delta_w^T x)| &\leq |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) - h(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x))| + |h(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) - h(\delta_w^T x)| \\ &\leq \zeta \cdot |\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)| \cdot \max_{u \in [-a, a]} \left| \frac{h(u)}{u} \right| \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x) - \delta_w^T x| \\ &\leq \zeta \cdot |\delta_w^T x| \cdot \max_{u \in [-a, a]} \left| \frac{h(u)}{u} \right| \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x) - \delta_w^T x|, \end{aligned}$$

where in the last line we are implicitly also dropping terms that are  $O(\varepsilon^2)$ . Next, we need to bound this max term. To do so, notice that for any  $u$ ,

$$\left| \frac{h(u)}{u} \right| = \left| \frac{l'(\phi + u) - l'(\phi)}{u} \right| \leq L.$$

So,

$$\begin{aligned} |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) - h(\delta_w^T x)| &\leq L \cdot |\delta_w^T x| \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x) - \delta_w^T x| \\ &\leq L \cdot |\delta_w^T x| \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\delta_w|^T|x| \cdot d \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

It also follows that, if  $x_j$  is the  $j$ th component of the vector  $x$ ,

$$\begin{aligned} &|\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) \text{fl}_{10}(x_j) - h(\delta_w^T x) x_j| \\ &\leq |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) \text{fl}_{10}(x_j) - \tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) x_j| + |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) x_j - h(\delta_w^T x) x_j| \\ &\leq |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) x_j| \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\delta_w^T x| \cdot |x_j| \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\delta_w|^T|x| \cdot |x_j| \cdot d \cdot \varepsilon_{\text{machine-lo}} \\ &\leq 2L \cdot |\delta_w^T x| \cdot |x_j| \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\delta_w|^T|x| \cdot |x_j| \cdot d \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

where in the second-to last line we implicitly dropped terms that are  $O(\varepsilon^2)$ , and in the last line we leveraged the fact that  $h$  is Lipschitz continuous with parameter  $L$ . It follows that the error for the  $j$ th component of  $\delta_g$ , which we denote  $\delta_{g,j}$ , is

$$\begin{aligned} |\delta_{g,j} - (h(\delta_w^T x) x_j + \sigma \delta_{w,j})| &= |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) \text{fl}_{10}(x_j) + \text{fl}_{10}(\sigma) \text{fl}_{10}(\delta_{w,j}) - (h(\delta_w^T x) x_j + \sigma \delta_{w,j})| \\ &\leq |\tilde{h}(\text{fl}_{10}(\delta_w)^T \text{fl}_{10}(x)) \text{fl}_{10}(x_j) - h(\delta_w^T x) x_j| + |\text{fl}_{10}(\sigma) \text{fl}_{10}(\delta_{w,j}) - \sigma \delta_{w,j}| \\ &\leq \left(2L \cdot |\delta_w^T x| \cdot |x_j| \cdot \zeta + L \cdot |\delta_w|^T|x| \cdot |x_j| \cdot d + |\sigma \delta_{w,j}|\right) \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$



It follows that we can bound this in terms of the infinity norm as

$$\|\delta_g - (h(\delta_w^T x) x + \sigma \delta_w)\|_\infty \leq \left( 2L \cdot |\delta_w^T x| \cdot \|x\|_\infty \cdot \zeta + L \cdot |\delta_w|^T |x| \cdot \|x\|_\infty \cdot d + |\sigma| \cdot \|\delta_w\|_\infty \right) \cdot \varepsilon_{\text{machine-lo}}$$

or in terms of the  $\ell_2$ -norm as

$$\begin{aligned} \|\delta_g - (h(\delta_w^T x) x + \sigma \delta_w)\| &\leq 2L \cdot |\delta_w^T x| \cdot \|x\| \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot |\delta_w|^T |x| \cdot \|x\| \cdot d \cdot \varepsilon_{\text{machine-lo}} + |\sigma| \cdot \|\delta_w\| \cdot \varepsilon_{\text{machine-lo}} \\ &\leq 2L \cdot \|\delta_w\| \cdot \|x\|^2 \cdot \zeta \cdot \varepsilon_{\text{machine-lo}} + L \cdot \|\delta_w\| \|x\|^2 \cdot d \cdot \varepsilon_{\text{machine-lo}} + |\sigma| \cdot \|\delta_w\| \cdot \varepsilon_{\text{machine-lo}} \\ &= \left( L \cdot \|x\|^2 \cdot (2\zeta + d) + |\sigma| \right) \cdot \|\delta_w\| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

Notice a very important property of this error bound: it is proportional to  $\delta_w$ . That is, as our offset guess becomes closer to the true value of  $w$  (equivalently, as  $\delta_w$  gets smaller), the error of our  $\delta_g$  computation becomes smaller at the same rate.

### 3.2 Full Learning Algorithms

Combining everything together produces the following full algorithm for learning on linear models. We call this method *Bit-Centered SGD*.

---

#### Algorithm 2 BC-SGD (Bit-Centered SGD) for Linear Models

---

```

1: given:  $N$  loss functions  $l_i$  and training examples  $x_i$  in high precision, number of epochs  $K$ , epoch length  $T$ , step size  $\alpha$ 
2: given: initial iterate  $w_1 = o_{w,1}$  in high precision
3: for  $k = 1$  to  $K$  do
4:   for  $i = 1$  to  $N$  do
5:      $\phi_{i,k} \leftarrow \text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})$ 
6:     define:  $h_i(\psi) = \alpha l'_i(\phi_i + \psi) - \alpha l'_i(\phi_i)$ 
7:     compute approximation:  $\tilde{h}_i \approx h_i$ 
8:     quantize:  $\tilde{x}_i \leftarrow \text{fl}_{\text{lo}}(x_i)$ 
9:     quantize:  $\tilde{g}_{k,i} \leftarrow \text{fl}_{\text{lo}}(\text{fl}_{\text{hi}}(\alpha) \cdot l'_i(\text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})) \cdot \text{fl}_{\text{hi}}(x_i))$ 
10:   end for
11:    $\delta_{w,k,0} \leftarrow \text{fl}_{\text{lo}}(0)$ 
12:   for  $t = 1$  to  $T$  do
13:     sample  $i$  uniformly from  $\{1, \dots, N\}$ 
14:      $\psi_{k,t} \leftarrow \text{fl}_{\text{lo}}(\tilde{x}_i)^T \text{fl}_{\text{lo}}(\delta_{w,k,t-1})$ 
15:      $u_{k,t} \leftarrow \tilde{h}_i(\text{fl}_{\text{lo}}(\psi_{k,t})) \cdot \text{fl}_{\text{lo}}(\tilde{x}_i)$ 
16:     update model:  $\delta_{w,k,t} \leftarrow \text{fl}_{\text{lo}}(\delta_{w,k,t-1}) - \text{fl}_{\text{lo}}(u_{k,t}) - \text{fl}_{\text{lo}}(\tilde{g}_{k,i})$ 
17:   end for
18:    $o_{w,k+1} \leftarrow \text{fl}_{\text{hi}}(o_{w,k}) + \text{fl}_{\text{hi}}(\delta_{w,k,T})$ 
19: end for
20: return  $o_{w,K+1}$ 

```

---

Notice that the entire inner loop (here outlined in red) is computed only using low-precision arithmetic. Furthermore, the bulk of the outer-loop computation (the part outlined in blue) is embarrassingly parallel across training examples. The bit centering method also naturally lends itself to a variant of *stochastic variance reduced gradient* (SVRG) Johnson and Zhang [2]: all that we need to do is replace the individual stored gradients  $\tilde{g}_{k,i}$  with the average gradient.

---

**Algorithm 3** BC-SVRG (Bit-Centered SVRG) for Linear Models
 

---

```

1: given:  $N$  loss functions  $l_i$  and training examples  $x_i$  in high precision, number of epochs  $K$ , epoch length  $T$ , step size  $\alpha$ 
2: given: initial iterate  $w_1 = o_{w,1}$  in high precision
3: for  $k = 1$  to  $K$  do
4:   for  $i = 1$  to  $N$  do
5:      $\phi_{i,k} \leftarrow \text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})$ 
6:     define:  $h_i(\psi) = \alpha l'_i(\phi_i + \psi) - \alpha l'_i(\phi_i)$ 
7:     compute approximation:  $\tilde{h}_i \approx h_i$ 
8:     quantize:  $\tilde{x}_i \leftarrow \text{fl}_{\text{lo}}(x_i)$ 
9:   end for
10:   $g_k \leftarrow \frac{\text{fl}_{\text{hi}}(\alpha)}{N} \sum_{i=1}^N l'_i(\text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})) \cdot \text{fl}_{\text{hi}}(x_i)$ 
11:  quantize:  $\tilde{g}_k \leftarrow \text{fl}_{\text{lo}}(g_k)$ 
12:   $\delta_{w,k,0} \leftarrow \text{fl}_{\text{lo}}(0)$ 
13:  for  $t = 1$  to  $T$  do
14:    sample  $i$  uniformly from  $\{1, \dots, N\}$ 
15:     $\psi_{k,t} \leftarrow \text{fl}_{\text{lo}}(\tilde{x}_i)^T \text{fl}_{\text{lo}}(\delta_{w,k,t-1})$ 
16:     $u_{k,t} \leftarrow \tilde{h}_i(\text{fl}_{\text{lo}}(\psi_{k,t})) \cdot \text{fl}_{\text{lo}}(\tilde{x}_i)$ 
17:    update model:  $\delta_{w,k,t} \leftarrow \text{fl}_{\text{lo}}(\delta_{w,k,t-1}) - \text{fl}_{\text{lo}}(u_{k,t}) - \text{fl}_{\text{lo}}(\tilde{g}_k)$ 
18:  end for
19:   $o_{w,k+1} \leftarrow \text{fl}_{\text{hi}}(o_{w,k}) + \text{fl}_{\text{hi}}(\delta_{w,k,T})$ 
20: end for
21: return  $o_{w,K+1}$ 

```

---

Notice again that the entire inner loop of the algorithm (the part outlined in red) will be done *entirely* using low-precision loads and stores, since it only reads and writes the variables  $v$ ,  $z$ , and  $\tilde{h}$ , all of which are low-precision values. Additionally, the entire inner loop of the algorithm will also be computed using only low-precision arithmetic. Furthermore, the bulk of the outer-loop computation (the part outlined in blue) is embarrassingly parallel across training examples. Figure 1 illustrates the convergence of BC-SGD and BC-SVRG on a toy logistic regression problem, with model dimension  $d = 128$ , number of training examples  $N = 1024$ . This problem was generated by first sampling a random generative model  $w_{\text{gen}} \sim \mathcal{N}(0, I_d)$ , where  $\mathcal{N}$  denotes the normal distribution and  $I_d \in \mathbb{R}^{d \times d}$  is the identity matrix, and then sampling  $n$  training examples  $(x_i, y_i)$  for  $i \in \{1, \dots, N\}$  as

$$x_i \sim \mathcal{N}(0, I_d), \quad /du_i \sim \mathcal{N}(0, 1), \quad y_i = \text{sign}(x_i^T w_{\text{gen}} + u_i).$$

This process results in a logistic regression problem with objective

$$f(w) = \frac{1}{N} \sum_{i=1}^n f_i(w) = \frac{1}{N} \sum_{i=1}^n \log(1.0 + \exp(-y_i x_i^T w)).$$

This problem will have condition number  $\kappa = 1150$  near the global optimum (the condition number  $\kappa$  in this context is the ratio between the maximum Lipschitz constant  $L$  of  $\nabla f_i$  and the strong convexity constant  $\mu$  of  $f$ ). For the experiments, we ran each algorithm for  $K = 200$  epochs, each of which had  $T = 1000$  inner iterations. For SGD and BC-SGD, we used a learning rate  $\alpha = 0.03$ , and for BC-SVRG we used  $\alpha = 1.0$ ; these learning rates were hand-tuned.

Notice that the SVRG implementation has the very interesting property of having *high accuracy*. That is, it is able to converge to solutions of accuracy determined by the machine epsilon of the high-precision numbers, even though it uses low-precision numbers for its inner loop computations (in this case, sometimes very low precision numbers). The figure also illustrates an effect that has been previously observed: that randomized rounding, which rounds up or down at random when quantizing such that the quantized output is an unbiased estimator of the original value, performs better than standard round-to-nearest quantization.

Another interesting observation in these figures is that when the number of exponent bits is small, performance degrades for BC-SVRG, and the algorithm converges to a higher error than comparable runs with lower numbers of bits. This occurs because eventually, the low-precision  $\delta$  values used in bit centering start underflowing, which limits convergence to a level at which  $\delta$  values are still large enough to be representable in low precision. One way to address

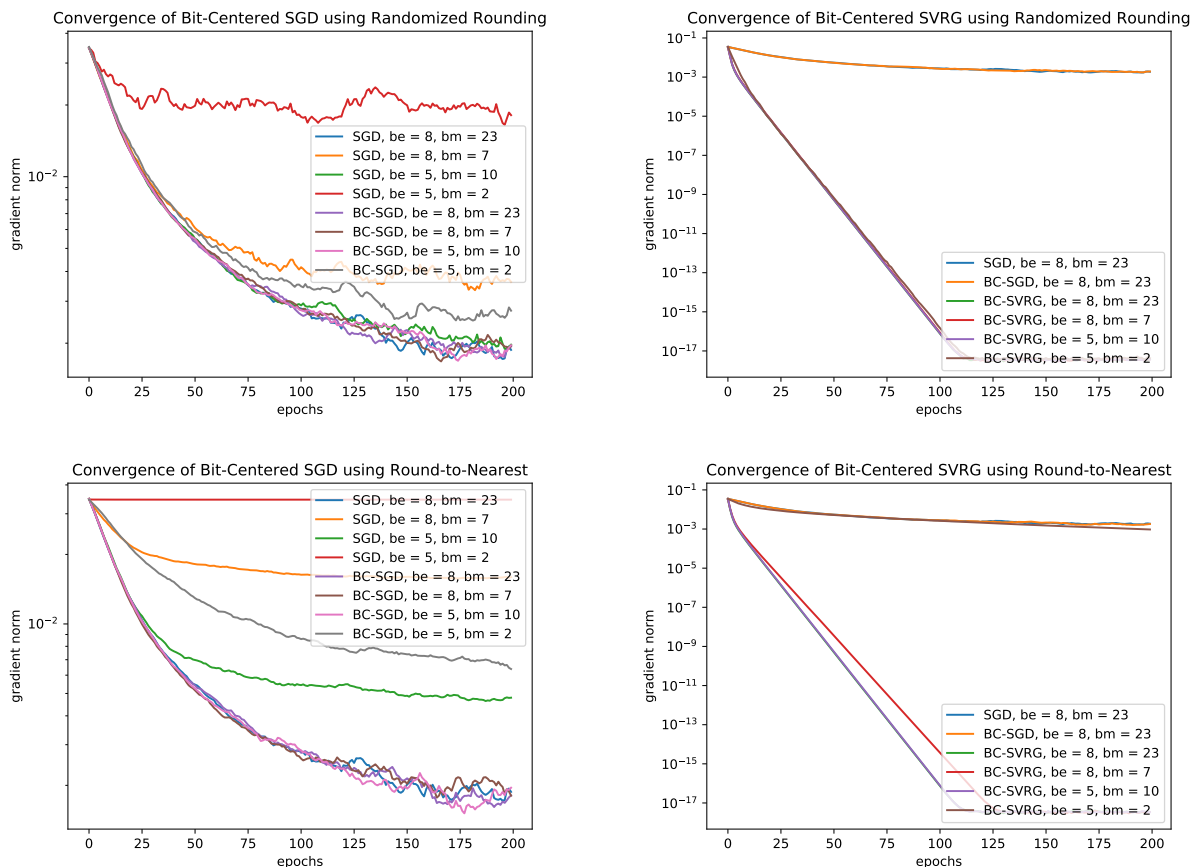


Figure 1: Convergence of BC-SGD (left) and BC-SVRG (right) compared with ordinary quantize-first low-precision SGD on a toy synthetic logistic regression problem (with model dimension  $d = 128$ , number of training examples  $n = 1024$ , condition number  $\kappa = 1150$  near the global optimum, and no regularization). Single-precision floating point numbers (8-bit exponent, 23-bit mantissa) were used for the high-precision numbers in all experiments, and for low-precision numbers the number of exponent bits  $be$  and the number of mantissa bits  $bm$  are reported (the number of exponent bits was fixed at 11, as is the case for double-precision floats). Randomized rounding was used for all quantization in the top row, while nearest-neighbor rounding was used for all quantization in the bottom row.

this is by dynamically adjusting the *exponent bias*. Exponent bias is used in floating point numbers to control the range of the exponents that are supported. For example, a single precision number with exponent  $e$  (an unsigned 8-bit number) and mantissa  $m$  (a 23-bit digit string) has value

$$(-1)^{\text{sign}} \cdot 2^{e-127} \cdot (1.m);$$

this 127 number is the exponent bias. Instead of using a fixed bias for representing  $\delta$  for bit-centering, we can use a bias that decreases over time proportional to the magnitude of the  $\delta$  values we know we are going to compute. That is, we define a new low-precision floating point representation in which the value is

$$(-1)^{\text{sign}} \cdot 2^{e-127+\text{extra bias}} \cdot (1.m);$$

We let  $\text{fl}_{\text{lo-bias}}()$  denote this bias-adjusted low-precision representation. Note that we would still want to use standard-bias low-precision numbers for storing intermediate values that do not decrease in magnitude with  $\delta$ , such as the  $\text{fl}_{\text{lo}}(o_x)$  used in a bit-centered multiply. Doing this results in the following algorithm.

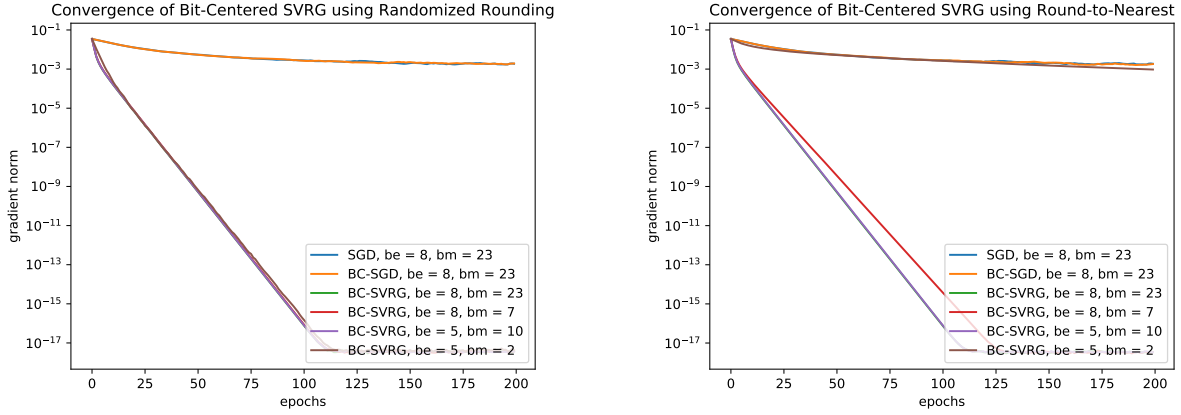


Figure 2: Convergence of dynamic-bias-adjusted BC-SVRG with randomized rounding (left) and nearest-neighbor rounding (right) on the same toy synthetic logistic regression problem used in Figure 1. For the bias control parameter, we used  $\chi = 0.01$ .

---

#### Algorithm 4 Dynamic-Bias-Adjusted BC-SVRG for Linear Models

---

- 1: **given:**  $N$  loss functions  $l_i$  and training examples  $x_i$  in high precision, number of epochs  $K$ , epoch length  $T$ , step size  $\alpha$
  - 2: **given:** initial iterate  $w_1 = o_{w,1}$  in high precision
  - 3: **given:** bias control parameter  $\chi$
  - 4: **for**  $k = 1$  to  $K$  **do**
  - 5:   **for**  $i = 1$  to  $N$  **do**
  - 6:      $\phi_{i,k} \leftarrow \text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})$
  - 7:     **define:**  $h_i(\psi) = \alpha l'_i(\phi_i + \psi) - \alpha l'_i(\phi_i)$
  - 8:     **compute approximation:**  $\tilde{h}_i \approx h_i$
  - 9:     **quantize:**  $\tilde{x}_i \leftarrow \text{fl}_{\text{lo}}(x_i)$
  - 10:   **end for**
  - 11:    $g_k \leftarrow \frac{\text{fl}_{\text{hi}}(\alpha)}{N} \sum_{i=1}^N \cdot l'_i(\text{fl}_{\text{hi}}(x_i)^T \text{fl}_{\text{hi}}(o_{w,k})) \cdot \text{fl}_{\text{hi}}(x_i)$
  - 12:   **extra bias**  $\leftarrow \lfloor \log_2(\chi \cdot \|g_k\|_\infty) \rfloor$
  - 13:   **quantize:**  $\tilde{g}_k \leftarrow \text{fl}_{\text{lo-bias}}(g_k)$
  - 14:    $\delta_{w,k,0} \leftarrow \text{fl}_{\text{lo-bias}}(0)$
  - 15:   **for**  $t = 1$  to  $T$  **do**
  - 16:     **sample**  $i$  uniformly from  $\{1, \dots, N\}$
  - 17:      $\psi_{k,t} \leftarrow \text{fl}_{\text{lo}}(\tilde{x}_i)^T \text{fl}_{\text{lo-bias}}(\delta_{w,k,t-1})$
  - 18:      $u_{k,t} \leftarrow \tilde{h}_i(\text{fl}_{\text{lo-bias}}(\psi_{k,t})) \cdot \text{fl}_{\text{lo}}(\tilde{x}_i)$
  - 19:     **update model:**  $\delta_{w,k,t} \leftarrow \text{fl}_{\text{lo-bias}}(\delta_{w,k,t-1}) - \text{fl}_{\text{lo-bias}}(u_{k,t}) - \text{fl}_{\text{lo-bias}}(\tilde{g}_k)$
  - 20:   **end for**
  - 21:    $o_{w,k+1} \leftarrow \text{fl}_{\text{hi}}(o_{w,k}) + \text{fl}_{\text{hi}}(\delta_{w,k,T})$
  - 22: **end for**
  - 23: **return**  $o_{w,K+1}$
- 

Note that the extra bias value is global across all numbers, and so little extra memory is used to compute and store it. This technique can have a dramatic effect on the convergence of BC-SVRG, as shown in Figure 2.

## 4 Bit Centering for General Loss Functions

In the previous section, we described how to do bit centering for linear loss functions, which had a particular form. We noticed that we were able to produce a bound on the error of the delta of the gradients of the form

$$\text{error}(\delta_g) \leq C \cdot \|\delta_w\| \cdot \varepsilon_{\text{machine-lo}}$$

for a constant  $C$  that was a function of the problem parameters. Next, we will show how we can get a bound of the same form on the delta error of the gradients for *arbitrary* loss functions which can be computed via a computation graph. We will do this with a particular focus on deep neural networks as an application.

To start, recall an interesting property of our error analysis for the linear models case. When we did the analysis, we derived bounds on  $o_g$  using a particular approach, and then repeated *almost the same calculation* to get a bound on  $\delta_g$ . This suggests that a more sophisticated analysis of a bit centering computation can get a bound on the error of computing the offset *and* the error of computing the delta at the same time.

Let  $\text{bc}(x)$  denote the value  $x$  represented using bit centering. That is, we represent it as  $x = o_x + \delta_x$  using a high-precision-floating-point offset  $o_x$  and a low-precision-floating-point delta  $\delta_x$ . Then we can compute the error of bit-centered operations just like we would compute the error of floating point operations. We will start this section by doing this. But first, we present a table summarizing our results. For matrix operands, we assume dimension of  $m \times n$ , and for matrix product the first operand is  $m \times n$  and the second is  $n \times p$ .

Operation	# LP Ops	# HP Ops	Result error ( $\cdot \varepsilon_{\text{machine-hi}}$ )
Scalar sum	1 add	1 add	$ x + y  +  \delta_x + \delta_y  \cdot (k + 1)$
Scalar difference	1 sub	1 sub	$ x - y  +  \delta_x - \delta_y  \cdot (k + 1)$
Scalar product	2 extra loads 2 multiplies 2 adds	1 multiply 2 extra stores	$ x \cdot y  + ( o_x \cdot \delta_y  +  \delta_x \cdot o_y  +  \delta_x \cdot \delta_y ) \cdot (4k + 1)$
Scalar function eval	$O(1)$ extra loads 1 eval of $\tilde{h}$	1 eval of $f$ 1 function approx $O(1)$ extra stores	$ f(x)  + (k \cdot \zeta + 1) \cdot L_f \cdot  \delta_x $
Matrix sum	$mn$ adds	$mn$ adds	$\ X + Y\  + \ \delta_X + \delta_Y\  \cdot (k + 1)$
Matrix difference	$mn$ subs	$mn$ subs	$\ X - Y\  + \ \delta_X - \delta_Y\  \cdot (k + 1)$
Hadamard product	$2mn$ extra loads $2mn$ multiplies $2mn$ adds	$mn$ multiplies $2mn$ extra stores	$\ X \cdot Y\  +$ $\   o_X \cdot \delta_Y  +  \delta_X \cdot o_Y  +  \delta_X \cdot \delta_Y  \  \cdot (4k + 1)$
Entrywise f eval	$O(mn)$ extra loads $mn$ evals of $\tilde{h}$	$mn$ evals of $f$ $mn$ function approx $O(mn)$ extra stores	$\ f(x)\  + (k \cdot \zeta + 1) \cdot L_f \cdot \ \delta_x\ $
Matrix product	$mn + np$ extra loads $2mnp$ multiplies $2mnp + mp + np$ adds	$mnp$ multiplies $mnp$ adds $mn + np$ extra stores	$\   X  \cdot  Y  \  \cdot n +$ $\   o_X  \cdot  \delta_Y  +  \delta_X  \cdot  o_Y  +  \delta_X  \cdot  \delta_Y  \  \cdot$ $(nk + 3k + n)$
Const matrix prod	$np$ extra loads $mnp$ multiplies $mnp$ adds	$mnp$ multiplies $mnp$ adds $np$ extra stores	$\   X  \cdot  o_Y  \  \cdot n + \   \delta_X  \cdot  o_Y  \  \cdot (nk + k + 1)$

## 4.1 Computing the Error of Bit-Centered Operations

**Addition** Suppose that we want to do addition using bit centering, to compute  $\text{bc}(x) + \text{bc}(y)$  where  $x = o_x + \delta_x$  and  $y = o_y + \delta_y$ . The natural way to do this is to let

$$\text{bc}(x) + \text{bc}(y) \stackrel{\text{def}}{=} (\text{fl}_{\text{hi}}(o_x) + \text{fl}_{\text{hi}}(o_y)) + (\text{fl}_{\text{lo}}(\delta_x) + \text{fl}_{\text{lo}}(\delta_y)),$$

where more explicitly the offset is the first term  $\text{fl}_{\text{hi}}(o_x) + \text{fl}_{\text{hi}}(o_y)$  and the delta is the second term  $\text{fl}_{\text{lo}}(\delta_x) + \text{fl}_{\text{lo}}(\delta_y)$ . This will have error

$$\begin{aligned} |\text{bc}(x) + \text{bc}(y) - (x + y)| &= |(\text{fl}_{\text{hi}}(o_x) + \text{fl}_{\text{hi}}(o_y)) + (\text{fl}_{\text{lo}}(\delta_x) + \text{fl}_{\text{lo}}(\delta_y)) - (o_x + o_y + \delta_x + \delta_y)| \\ &\leq |(\text{fl}_{\text{hi}}(o_x) + \text{fl}_{\text{hi}}(o_y)) - (o_x + o_y)| + |(\text{fl}_{\text{lo}}(\delta_x) + \text{fl}_{\text{lo}}(\delta_y)) - (\delta_x + \delta_y)| \\ &\leq |o_x + o_y| \cdot \varepsilon_{\text{machine-hi}} + |\delta_x + \delta_y| \cdot \varepsilon_{\text{machine-lo}} \\ &\leq |x + y| \cdot \varepsilon_{\text{machine-hi}} + |\delta_x + \delta_y| \cdot (k + 1) \cdot \varepsilon_{\text{machine-hi}}. \end{aligned}$$

**Subtraction** The analysis is the same. We will have

$$|\text{bc}(x) - \text{bc}(y) - (x - y)| \leq |x - y| \cdot \varepsilon_{\text{machine-hi}} + |\delta_x - \delta_y| \cdot (k + 1) \cdot \varepsilon_{\text{machine-hi}}.$$

**Multiplication** Suppose that we want to do multiplication using bit centering, to compute  $\text{bc}(x) \cdot \text{bc}(y)$  where  $x = o_x + \delta_x$  and  $y = o_y + \delta_y$ . The natural way to do this is to let the offset of the result be computed as

$$o_{\text{bc}(x) \cdot \text{bc}(y)} \stackrel{\text{def}}{=} \text{fl}_{\text{hi}}(o_x) \cdot \text{fl}_{\text{hi}}(o_y)$$

and the delta be computed as

$$\delta_{\text{bc}(x) \cdot \text{bc}(y)} \stackrel{\text{def}}{=} \text{fl}_{\text{lo}}(o_x) \cdot \text{fl}_{\text{lo}}(\delta_y) + \text{fl}_{\text{lo}}(\delta_x) \cdot (\text{fl}_{\text{lo}}(o_y) + \text{fl}_{\text{lo}}(\delta_y)).$$

Note that in order to compute this delta, we will need to additionally quantize and load low-precision version of  $o_x$  and  $o_y$ ; this will require additional loads on the accelerator. Using the  $1 + \varepsilon$  model, this will look like

$$\begin{aligned} \text{bc}(x) \cdot \text{bc}(y) &= o_{\text{bc}(x) \cdot \text{bc}(y)} + \delta_{\text{bc}(x) \cdot \text{bc}(y)} \\ &= \text{fl}_{\text{hi}}(o_x) \cdot \text{fl}_{\text{hi}}(o_y) + \text{fl}_{\text{lo}}(o_x) \cdot \text{fl}_{\text{lo}}(\delta_y) + \text{fl}_{\text{lo}}(\delta_x) \cdot (\text{fl}_{\text{lo}}(o_y) + \text{fl}_{\text{lo}}(\delta_y)) \\ &= (o_x \cdot o_y) \cdot (1 + \varepsilon_{\text{hi}}) \\ &\quad + \left( ((o_x \cdot (1 + \varepsilon_{\text{lo}})) \cdot \delta_y) \cdot (1 + \varepsilon_{\text{lo}}) \right. \\ &\quad \left. + (\delta_x \cdot ((o_y \cdot (1 + \varepsilon_{\text{lo}})) + \delta_y) \cdot (1 + \varepsilon_{\text{lo}})) \cdot (1 + \varepsilon_{\text{lo}}) \right) \cdot (1 + \varepsilon_{\text{lo}}) \end{aligned}$$

where as usual in the  $1+\varepsilon$  model, all the  $\varepsilon$  in this expression can take on different values and are bounded as appropriate by the machine epsilon. By dropping terms in  $\varepsilon^2$ , we can simplify this to

$$\begin{aligned}
\text{bc}(x) \cdot \text{bc}(y) &= (o_x \cdot o_y) \cdot (1 + \varepsilon_{\text{hi}}) \\
&\quad + (o_x \cdot \delta_y) \cdot 2\varepsilon_{\text{lo}} \\
&\quad + (\delta_x \cdot o_y) \cdot \varepsilon_{\text{lo}} \\
&\quad + (\delta_x \cdot (o_y + \delta_y)) \cdot 2\varepsilon_{\text{lo}} \\
&\quad + (o_x \cdot \delta_y + \delta_y \cdot (o_y + \delta_y)) \cdot (1 + \varepsilon_{\text{lo}}) \\
&= (x \cdot y) \cdot (1 + \varepsilon_{\text{hi}}) \\
&\quad + (o_x \cdot \delta_y) \cdot 2\varepsilon_{\text{lo}} \\
&\quad + (\delta_x \cdot o_y) \cdot \varepsilon_{\text{lo}} \\
&\quad + (\delta_x \cdot (o_y + \delta_y)) \cdot 2\varepsilon_{\text{lo}} \\
&\quad + (o_x \cdot \delta_y + \delta_y \cdot (o_y + \delta_y)) \cdot (\varepsilon_{\text{lo}} + \varepsilon_{\text{hi}}) \\
&= (x \cdot y) \cdot (1 + \varepsilon_{\text{hi}}) \\
&\quad + (o_x \cdot \delta_y) \cdot (3k + 1) \cdot \varepsilon_{\text{hi}} \\
&\quad + (\delta_x \cdot o_y) \cdot (4k + 1) \cdot \varepsilon_{\text{hi}} \\
&\quad + (\delta_x \cdot \delta_y) \cdot (3k + 1) \cdot \varepsilon_{\text{hi}}.
\end{aligned}$$

It follows that the error will be bounded by

$$|\text{bc}(x) \cdot \text{bc}(y) - (x \cdot y)| \leq |x \cdot y| \cdot \varepsilon_{\text{hi}} + (|o_x \cdot \delta_y| + |\delta_x \cdot o_y| + |\delta_x \cdot \delta_y|) \cdot (4k + 1) \cdot \varepsilon_{\text{hi}}.$$

**Function application** Suppose that we want to apply a function to a value stored using bit centering, to compute  $f(\text{bc}(x))$  where  $x = o_x + \delta_x$ . Assume that, as we did above, we are able to compute  $f$  to machine precision in high-precision arithmetic, and assume that  $f$  is Lipschitz continuous with parameter  $L_f$ . The natural way to do bit centering for this, following our approach in the linear model case, is to let the offset of the result be computed as

$$o_{f(\text{bc}(x))} \stackrel{\text{def}}{=} f(\text{fl}_{\text{hi}}(o_x)).$$

Next, define the function

$$h(\psi) = f(o_x + \psi) - f(o_x),$$

assume that it is sufficiently nice to approximate, and using our function approximation capabilities, compute a low-precision approximation  $\tilde{h} \approx h$  over an appropriate range (either one computed by looking at  $\delta_x$  or one computed from some *a priori* bound on  $\delta_x$ ). Recall from before that we assumed that given any “sufficiently nice” scalar function  $h : \mathbb{R} \rightarrow \mathbb{R}$  with  $h(0) = 0$  and any range  $[-a, a]$  about zero, we can (via some amount of high-precision computation) produce a low-precision-computable routine  $\tilde{h}$  that can compute  $h(x)$  for any low-precision number  $x \in [-a, a]$  to within relative accuracy

$$\left| \tilde{h}(x) - h(x) \right| \leq \zeta \cdot |x| \cdot \max_{u \in [-a, a]} \left| \frac{h(u)}{u} \right| \cdot \varepsilon_{\text{machine-lo}}$$

Once we’ve computed the approximation  $\tilde{h}$ , we can use it to compute the delta

$$\delta_{f(\text{bc}(x))} \stackrel{\text{def}}{=} \tilde{h}(\text{fl}_{\text{lo}}(\delta_x)).$$

Note that in order to compute this, we will need to load the value of  $\tilde{h}$ . This will require additional loads on the accelerator. The resulting computed bit centered value will have error

$$\begin{aligned}
|f(\text{bc}(x)) - f(x)| &= \left| f(\text{fl}_{\text{hi}}(o_x)) + \tilde{h}(\text{fl}_{\text{lo}}(\delta_x)) - (f(o_x) + h(\delta_x)) \right| \\
&\leq |f(\text{fl}_{\text{hi}}(o_x)) - f(o_x)| + \left| \tilde{h}(\text{fl}_{\text{lo}}(\delta_x)) - h(\delta_x) \right| \\
&\leq |f(o_x)| \cdot \varepsilon_{\text{machine-hi}} + \zeta \cdot |\delta_x| \cdot \max_{u \in [-a, a]} \left| \frac{h(u)}{u} \right| \cdot \varepsilon_{\text{machine-lo}} \\
&\leq f(o_x) \cdot \varepsilon_{\text{machine-hi}} + \zeta \cdot |\delta_x| \cdot \max_{u \in [-a, a]} \left| \frac{f(o_x + u) - f(o_x)}{u} \right| \cdot \varepsilon_{\text{machine-lo}} \\
&\leq |f(x) + (f(o_x) - f(x))| \cdot \varepsilon_{\text{machine-hi}} + \zeta \cdot |\delta_x| \cdot L_f \cdot \varepsilon_{\text{machine-lo}} \\
&\leq |f(x)| \cdot \varepsilon_{\text{machine-hi}} + L_f \cdot |\delta_x| \cdot \varepsilon_{\text{machine-hi}} + \zeta \cdot |\delta_x| \cdot L_f \cdot \varepsilon_{\text{machine-lo}} \\
&\leq |f(x)| \cdot \varepsilon_{\text{machine-hi}} + (k \cdot \zeta + 1) \cdot L_f \cdot |\delta_x| \cdot \varepsilon_{\text{machine-hi}}.
\end{aligned}$$

**The fundamental axiom of bit centered arithmetic** Now that we've done this, we can state an analogue of the fundamental axiom of floating point arithmetic, but applied to bit centering. Specifically, for any primitive operator  $\star$  (except division), and any bit-centered values  $x$  and  $y$

$$|\text{bc}(x) \star \text{bc}(y) - x| \leq |x \star y| \cdot \varepsilon_{\text{machine-hi}} + (O(|\delta_x|) + O(|\delta_y|) + O(|\delta_x \delta_y|)) \cdot \varepsilon_{\text{machine-lo}}$$

where the big-O terms here hide parameters that depend on the operator and possibly on the offsets  $o_x$  and  $o_y$ . But the important thing to notice here is that the low-precision error *is proportional to the input delta*. So as long as we keep these deltas small, we can get very accurate computations, for *arbitrary* compute graphs.

## 4.2 Vector and Matrix Operations

While the above analysis does give us a way of bounding the error of any deep neural network gradient computed with bit centering, this bound may be annoying to compute in practice because so far, we have only bounded the error of *scalar* computations, and neural networks in both their forward and backward passes are written in terms of vector and matrix computations. Here, we derive bounds on the error of common vector and matrix computations done using bit centering.

**Addition** As a consequence of our scalar bounds, the following will hold for any norm for which  $|X_{i,j}| \leq Y_{i,j}$  for all indices  $i$  and  $j$  implies  $\|X\| \leq \|Y\|$ . Note that this condition will hold for all induced norms, as well as the Frobenius norm.

$$\|\text{bc}(x) + \text{bc}(y) - (x + y)\| \leq \|x + y\| \cdot \varepsilon_{\text{machine-hi}} + \|\delta_x + \delta_y\| \cdot (k + 1) \cdot \varepsilon_{\text{machine-hi}}.$$

**Subtraction** An analogous result will also hold for subtraction.

$$\|\text{bc}(x) - \text{bc}(y) - (x - y)\| \leq \|x - y\| \cdot \varepsilon_{\text{machine-hi}} + \|\delta_x - \delta_y\| \cdot (k + 1) \cdot \varepsilon_{\text{machine-hi}}.$$

**Hadamard product** An analogous result will also hold for the Hadamard product (or entrywise product).

$$\|\text{bc}(x) \odot \text{bc}(y) - (x \odot y)\| \leq \|x \odot y\| \cdot \varepsilon_{\text{machine-hi}} + \||o_x \odot \delta_y| + |\delta_x \odot o_y| + |\delta_x \odot \delta_y|\| \cdot (4k + 1) \cdot \varepsilon_{\text{machine-hi}}.$$

**Entrywise function application** An analogous result will also hold for entrywise function application.

$$\|f(\text{bc}(x)) - f(x)\| \leq \|f(x)\| \cdot \varepsilon_{\text{machine-hi}} + (k \cdot \zeta + 1) \cdot L_f \cdot \|\delta_x\| \cdot \varepsilon_{\text{machine-hi}}.$$



**Matrix product** So far all our bounds have been trivial applications of our previous scalar bounds in matrix form. The bound for the matrix product is somewhat more complicated. Suppose that we want to do matrix multiplication using bit centering, to compute  $\text{bc}(X) \cdot \text{bc}(Y)$  where  $X \in \mathbb{R}^{m \times n}$  and  $Y \in \mathbb{R}^{n \times p}$  and  $X = o_X + \delta_X$  and  $Y = o_Y + \delta_Y$ . The natural way to do this is to let the offset of the result be computed as

$$o_{\text{bc}(X) \cdot \text{bc}(Y)} \stackrel{\text{def}}{=} \text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y)$$

and the delta be computed as

$$\delta_{\text{bc}(X) \cdot \text{bc}(Y)} \stackrel{\text{def}}{=} \text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y) + \text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)).$$

Note that as in the scalar product case, in order to compute this delta, we will need to additionally quantize and load low-precision versions of  $o_X$  and  $o_Y$ ; this will require additional loads on the accelerator. In index  $i, j$ , this will have error

$$\begin{aligned} & |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\ &= \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} + (\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - ((o_X + \delta_X) \cdot (o_Y + \delta_Y))_{i,j} \right| \\ &\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} - (o_X \cdot \delta_Y)_{i,j} \right| \\ &\quad + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (o_Y + \delta_Y))_{i,j} \right| \\ &\quad + |(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}| \cdot \varepsilon_{\text{machine-lo}} \end{aligned}$$

where this last term comes from bounding the error of the low-precision sum used to add  $(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j}$  and  $(\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}$  when computing the delta term. We can further bound this with

$$\begin{aligned} & |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\ &\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} - (o_X \cdot \delta_Y)_{i,j} \right| \\ &\quad + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} \right| \\ &\quad + \left| (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (o_Y + \delta_Y))_{i,j} \right| \\ &\quad + |(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}| \cdot \varepsilon_{\text{machine-lo}} \\ &\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} - (\text{fl}_{\text{lo}}(o_X) \cdot \delta_Y)_{i,j} \right| \\ &\quad + \left| (\text{fl}_{\text{lo}}(o_X) \cdot \delta_Y)_{i,j} - (o_X \cdot \delta_Y)_{i,j} \right| \\ &\quad + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} \right| \\ &\quad + \left| (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \delta_Y))_{i,j} \right| \\ &\quad + \left| (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \delta_Y))_{i,j} - (\delta_X \cdot (o_Y + \delta_Y))_{i,j} \right| \\ &\quad + |(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}| \cdot \varepsilon_{\text{machine-lo}} \\ &\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} - (\text{fl}_{\text{lo}}(o_X) \cdot \delta_Y)_{i,j} \right| \\ &\quad + (|o_X| \cdot |\delta_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ &\quad + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} - (\delta_X \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j} \right| \\ &\quad + (|\delta_X| \cdot |o_Y + \delta_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ &\quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ &\quad + |(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

All the remaining terms that we need to bound are effectively dot products, since computing a single entry of a matrix product is a dot product. Recall that the error of a  $n$ -dimensional floating point dot product is

$$|\text{fl}(x)^T \text{fl}(y) - x^T y| \leq |x|^T |y| \cdot n \cdot \varepsilon_{\text{machine}},$$

so for a matrix product with inner dimension  $n$ , it follows that

$$|(\text{fl}(X)\text{fl}(Y) - XY)_{i,j}| \leq (|X||Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine}}.$$

Using this, we can bound our error with

$$\begin{aligned} & |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\ & \leq (|o_X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|\text{fl}_{\text{lo}}(o_X)| \cdot |\delta_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|o_X| \cdot |\delta_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|\delta_X| \cdot |\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|\delta_X| \cdot |o_Y + \delta_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + |(\text{fl}_{\text{lo}}(o_X) \cdot \text{fl}_{\text{lo}}(\delta_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot (\text{fl}_{\text{lo}}(o_Y) + \text{fl}_{\text{lo}}(\delta_Y)))_{i,j}| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

Now dropping terms in  $\varepsilon^2$  gives us

$$\begin{aligned} & |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\ & \leq (|o_X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|o_X| \cdot |\delta_Y|)_{i,j} \cdot (n+1) \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|\delta_X| \cdot |o_Y + \delta_Y|)_{i,j} \cdot (n+1) \cdot \varepsilon_{\text{machine-lo}} + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + |(o_X \cdot \delta_Y)_{i,j} + (\delta_X \cdot (o_Y + \delta_Y))_{i,j}| \cdot \varepsilon_{\text{machine-lo}} \\ & \leq (|o_X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|o_X| \cdot |\delta_Y|)_{i,j} \cdot (n+2) \cdot \varepsilon_{\text{machine-lo}} \\ & \quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot (n+3) \cdot \varepsilon_{\text{machine-lo}} + (|\delta_X| \cdot |\delta_Y|)_{i,j} \cdot (n+2) \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

And so,

$$\begin{aligned} & |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\ & \leq (|X| \cdot |Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|o_X| \cdot |\delta_Y|)_{i,j} \cdot (nk + 2k + n) \cdot \varepsilon_{\text{machine-hi}} \\ & \quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot (nk + 3k + n) \cdot \varepsilon_{\text{machine-hi}} + (|\delta_X| \cdot |\delta_Y|)_{i,j} \cdot (nk + 2k + n) \cdot \varepsilon_{\text{machine-hi}} \\ & \leq (|X| \cdot |Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|o_X| \cdot |\delta_Y| + |\delta_X| \cdot |o_Y| + |\delta_X| \cdot |\delta_Y|)_{i,j} \cdot (nk + 3k + n) \cdot \varepsilon_{\text{machine-hi}}. \end{aligned}$$

Finally, for any norm of the type described above, we can bound the error with

$$\begin{aligned} & \|\text{bc}(X) \cdot \text{bc}(Y) - X \cdot Y\| \\ & \leq \| |X| \cdot |Y| \| \cdot n \cdot \varepsilon_{\text{machine-hi}} + \| |o_X| \cdot |\delta_Y| + |\delta_X| \cdot |o_Y| + |\delta_X| \cdot |\delta_Y| \| \cdot (nk + 3k + n) \cdot \varepsilon_{\text{machine-hi}}. \end{aligned}$$

**Special case: constant matrix product** One simpler special case of the matrix product is a product where one of the terms (say,  $Y$ ) is known *a priori*. In this case, we can set  $o_Y = Y$  and  $\delta_Y = 0$ , and substantially simplify the formula for computing the matrix product in bit centering to

$$o_{\text{bc}(X) \cdot \text{bc}(Y)} \stackrel{\text{def}}{=} \text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y)$$

and the delta be computed as

$$\delta_{\text{bc}(X) \cdot \text{bc}(Y)} \stackrel{\text{def}}{=} \text{fl}_{\text{lo}}(\delta_X) \cdot \text{fl}_{\text{lo}}(o_Y).$$

In index  $i, j$ , this will have error

$$\begin{aligned}
& |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\
&= \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} + (\text{fl}_{\text{lo}}(\delta_X) \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} - ((o_X + \delta_X) \cdot o_Y)_{i,j} \right| \\
&\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} - (\delta_X \cdot o_Y)_{i,j} \right| \\
&\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} - (\delta_X \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} \right| \\
&\quad + \left| (\delta_X \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} - (\delta_X \cdot o_Y)_{i,j} \right| \\
&\leq \left| (\text{fl}_{\text{hi}}(o_X) \cdot \text{fl}_{\text{hi}}(o_Y))_{i,j} - (o_X \cdot o_Y)_{i,j} \right| + \left| (\text{fl}_{\text{lo}}(\delta_X) \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} - (\delta_X \cdot \text{fl}_{\text{lo}}(o_Y))_{i,j} \right| \\
&\quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}}
\end{aligned}$$

Each of these terms is a dot product, so we can bound them with

$$\begin{aligned}
& |(\text{bc}(X) \cdot \text{bc}(Y))_{i,j} - (X \cdot Y)_{i,j}| \\
&\leq (|o_X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-lo}} \\
&\quad + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot \varepsilon_{\text{machine-lo}} \\
&\leq (|X| \cdot |o_Y|)_{i,j} \cdot n \cdot \varepsilon_{\text{machine-hi}} + (|\delta_X| \cdot |o_Y|)_{i,j} \cdot (nk + k + 1) \cdot \varepsilon_{\text{machine-hi}}.
\end{aligned}$$

And so, in terms of any norm of the type described above,

$$\begin{aligned}
& \|\text{bc}(X) \cdot \text{bc}(Y) - X \cdot Y\| \\
&\leq \| |X| \cdot |o_Y| \| \cdot n \cdot \varepsilon_{\text{machine-hi}} + \| |\delta_X| \cdot |o_Y| \| \cdot (nk + k + 1) \cdot \varepsilon_{\text{machine-hi}}.
\end{aligned}$$

## 5 Proof of Convergence for BC-SVRG

In this section, we will prove convergence for one algorithm using bit centering: bit centered SVRG. Here, describe the general BC-SVRG algorithm, and prove its convergence for strongly convex objectives.

---

### Algorithm 5 BC-SVRG: Bit-Centered Stochastic Variance-Reduced Gradient

---

- 1: **given:**  $N$  loss functions  $f_i$ , number of epochs  $K$ , epoch length  $T$ , step size  $\alpha$
  - 2: **given:** initial iterate  $w_1 = o_{w,1} \in \mathbb{R}^d$  in high precision
  - 3: **for**  $k = 1$  **to**  $K$  **do**
  - 4:   **for**  $i = 1$  **to**  $N$  **do**
  - 5:     **compute gradient in HP:**  $g_{k,i} \leftarrow \nabla f_i(\text{fl}_{\text{hi}}(o_{w,k}))$
  - 6:     **pre-compute and store** all necessary data for computing bit-centered delta of  $\nabla f_i$  later
  - 7:   **end for**
  - 8:    $\bar{g}_k \leftarrow \frac{\text{fl}_{\text{hi}}(\alpha)}{\text{fl}_{\text{hi}}(N)} \cdot \sum_{i=1}^N \text{fl}_{\text{hi}}(g_{k,i})$
  - 9:   **quantize and store:**  $\tilde{g}_k \leftarrow \text{fl}_{\text{lo}}(\bar{g}_k)$
  - 10:    $\delta_{w,k,0} \leftarrow \text{fl}_{\text{lo}}(0) \in \mathbb{R}^d$
  - 11:   **for**  $t = 1$  **to**  $T$  **do**
  - 12:     **sample**  $i$  uniformly from  $\{1, \dots, N\}$
  - 13:     **load** data necessary for computing bit-centered delta of  $\nabla f_i$
  - 14:     **compute**  $u_{k,t} \leftarrow \delta_{\nabla f_i(w_{k,t-1})} \approx \nabla f_i(o_{w,k} + \delta_{w,k,t-1}) - \nabla f_i(o_{w,k})$  **in low precision**
  - 15:     **update model:**  $\delta_{w,k,t} \leftarrow \text{fl}_{\text{lo}}(\delta_{w,k,t-1}) - \text{fl}_{\text{lo}}(\alpha) \cdot \text{fl}_{\text{lo}}(u_{k,t}) - \text{fl}_{\text{lo}}(\tilde{g}_k)$
  - 16:   **end for**
  - 17:   **sample**  $t$  uniformly from  $\{0, \dots, T-1\}$ , then set  $o_{w,k+1} \leftarrow \text{fl}_{\text{hi}}(o_{w,k}) + \text{fl}_{\text{hi}}(\delta_{w,k,t})$
  - 18: **end for**
  - 19: **return**  $o_{w,K+1}$
-

We will prove this algorithm converges at a linear rate, subject to the following assumptions. Note that as before, we are going to ignore terms in  $O(\varepsilon^2)$  in this section. First, we assume strong convexity with parameter  $\mu$  of the objective function  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$ , and Lipschitz continuity of the gradients of each component function  $\nabla f_i(w)$  with parameter  $L$ . These are standard assumptions used in the analysis of SVRG. Second, we assume that the quantization error of the computation of the delta on Line 5 is always bounded by

$$\|u_{k,t} - \nabla f_i(o_{w,k} + \delta_{w,k,t-1}) - \nabla f_i(o_{w,k})\| \leq C \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}}$$

for some constant parameter  $C$ . This will hold (over any bounded region, for some parameter  $C$ ) for any function  $\nabla f_i$  that is a composition of the operators we described in the above sections, as a consequence of the fundamental axiom of bit centered arithmetic. Third, we make the simplifying assumption that  $\varepsilon_{\text{machine-hi}} = 0$ : that the high-precision numbers used are effectively full-precision real-number arithmetic. This is equivalent to ignoring terms in  $\varepsilon_{\text{machine-hi}} = 0$ , and this is a reasonable assumption when we the high-precision format has enough precision that the accuracy of the output of Algorithm 5

Under these conditions, we can prove the following theorem.

**Theorem 1.** *Suppose that we run BC-SVRG under the above conditions. Then for any parameter  $0 < \gamma < 1$ , if we set our step size  $\alpha$  and epoch length  $T$  to be*

$$\alpha = \frac{\gamma}{4L(1+\gamma)} \quad T = \frac{8\kappa(1+\gamma)}{\gamma^2},$$

then BC-SVRG converges in expected value at a linear rate given by

$$\mathbf{E} [f(o_{w,k+1}) - f(w^*)] \leq \left( \gamma + 9 \cdot (23 + 2C) \cdot \kappa \cdot \sqrt{\frac{2}{\gamma}} \cdot \varepsilon_{\text{machine-lo}} + O(\varepsilon_{\text{machine-lo}}^2) \right) \cdot \mathbf{E} [f(o_{w,k}) - f(w^*)].$$

Note that in order for this rate to be useful, we need the low-precision format to have epsilon small enough that

$$9 \cdot (23 + 2C) \cdot \kappa \cdot \sqrt{\frac{2}{\gamma}} \cdot \varepsilon_{\text{machine-lo}} < 1.$$

To do this, we need

$$\varepsilon_{\text{machine-lo}} = O(\kappa^{-1}).$$

This theory exposes a relationship between the amount of precision we need and the condition number of the problem. However, as long as this is satisfied, this theorem shows that BC-SVRG can produce solutions of arbitrarily small objective gap—that is, its performance is not bound by the accuracy of the low-precision representation.

Now we will present the proof of the theorem.

*Proof of Theorem 1.* First, we analyze the error caused by the low-precision updates in the inner loop. First, we notice that, using the  $1 + \varepsilon$  model,

$$\begin{aligned} \delta_{w,k,t,j} &= ((\delta_{w,k,t-1,j} - \alpha \cdot u_{k,t,j} \cdot (1 + \varepsilon_{\text{lo}})) (1 + \varepsilon_{\text{lo}}) - \tilde{g}_{k,j}) (1 + \varepsilon_{\text{lo}}) \\ &= \delta_{w,k,t-1,j} \cdot (1 + 2\varepsilon_{\text{lo}}) - \alpha \cdot u_{k,t,j} (1 + 3\varepsilon_{\text{lo}}) - \tilde{g}_{k,j} (1 + \varepsilon_{\text{lo}}) \end{aligned}$$

where the extra subscript  $j$  denotes the  $j$ th element of the vector. It follows that

$$\begin{aligned} &\|\delta_{w,k,t} - (\delta_{w,k,t-1} - \alpha \cdot u_{k,t} - \tilde{g}_k)\| \\ &\leq ((\delta_{w,k,t-1} - \alpha \cdot u_{k,t} \cdot (1 + \varepsilon_{\text{lo}})) (1 + \varepsilon_{\text{machine-lo}}) - \tilde{g}_k) (1 + \varepsilon_{\text{machine-lo}}) \\ &= 2\|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 3\alpha \cdot \|u_{k,t}\| \cdot \varepsilon_{\text{machine-lo}} + 3\|\tilde{g}_k\| \varepsilon_{\text{machine-lo}} \\ &= 2\|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 3\alpha \cdot \|\nabla f_i(o_{w,k} + \delta_{w,k,t-1}) - \nabla f_i(o_{w,k})\| \cdot \varepsilon_{\text{machine-lo}} + 3\|\tilde{g}_k\| \cdot \varepsilon_{\text{machine-lo}} \end{aligned}$$

where in the last line we dropped terms in  $\varepsilon^2$ , as usual. Next, by Lipschitz continuity,

$$\begin{aligned} &\|\delta_{w,k,t} - (\delta_{w,k,t-1} - \alpha \cdot u_{k,t} - \tilde{g}_k)\| \\ &\leq 2\|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 3\alpha L \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 3\|\tilde{g}_k\| \cdot \varepsilon_{\text{machine-lo}} \\ &= (2 + 3\alpha L) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 3\|\tilde{g}_k\| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

Now, we also know that, since  $\tilde{g}$  is just a quantization of  $\bar{g}$ ,

$$\|\tilde{g}_k - \bar{g}_k\| \leq \|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}},$$

and so adding this to our bound by the triangle inequality gives us

$$\begin{aligned} & \|\delta_{w,k,t} - (\delta_{w,k,t-1} - \alpha \cdot u_{k,t} - \bar{g}_k)\| \\ & \leq (2 + 3\alpha L) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 4\|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

Finally, adding in our bound on the quantization error for the delta on Line 5 gives us

$$\begin{aligned} & \|\delta_{w,k,t} - (\delta_{w,k,t-1} - \alpha \cdot \nabla f_i(o_{w,k} + \delta_{w,k,t-1}) + \alpha \nabla f_i(o_{w,k}) - \bar{g}_k)\| \\ & \leq (2 + 3\alpha L) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 4\|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}} + C \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} \\ & \leq (2 + 3\alpha L + C) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 4\|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}}. \end{aligned}$$

Now we have this bound on the error of the update step, we will use it to analyze the expected distance to the optimum of the iterates as the algorithm progresses. At each step, define

$$w_{k,t} = o_{w,k} + \delta_{w,k,t}$$

and define

$$v_{k,t} = o_{w,k} + \delta_{w,k,t-1} - \alpha \nabla f_i(o_{w,k} + \delta_{w,k,t-1}) + \alpha \nabla f_i(o_{w,k}) - \bar{g}_k.$$

Notice that our bound above implies that

$$\|w_{k,t} - v_{k,t}\| \leq (2 + 3\alpha L + C) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 4\|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}}.$$

At the next timestep,

$$\begin{aligned} \mathbf{E} \left[ \|w_{k,t} - w^*\|^2 \right] &= \mathbf{E} \left[ \|v_{k,t} - w^* + (w_{k,t} - v_{k,t})\|^2 \right] \\ &= \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + 2(v_{k,t} - w^*)^T (w_{k,t} - v_{k,t}) + \|w_{k,t} - v_{k,t}\|^2 \right] \\ &= \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + 2(v_{k,t} - w^*)^T (w_{k,t} - v_{k,t}) \right], \end{aligned}$$

where we dropped the term  $\|w_{k,t} - v_{k,t}\|^2$  because our inequality above shows that this is a  $O(\varepsilon^2)$  term. Next, by Cauchy-Schwarz, we can further bound this with

$$\begin{aligned} \mathbf{E} \left[ \|w_{k,t} - w^*\|^2 \right] &\leq \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + 2\|v_{k,t} - w^*\| \|w_{k,t} - v_{k,t}\| \right] \\ &= \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + 2\|w_{k,t} - w^*\| \|w_{k,t} - v_{k,t}\| \right] \end{aligned}$$

where in the last line we again dropped  $O(\varepsilon^2)$  terms, since  $w_{k,t}$  differs from  $v_{k,t}$  by only at most a factor of  $\varepsilon$ . Since  $2xy \leq x^2 + y^2$ , it follows that for any  $\eta > 0$  (a parameter to be set later),

$$\mathbf{E} \left[ \|w_{k,t} - w^*\|^2 \right] \leq \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + \eta \|w_{k,t} - w^*\|^2 + \frac{1}{\eta} \|w_{k,t} - v_{k,t}\|^2 \right].$$

Note that we do not drop the  $\|w_{k,t} - v_{k,t}\|^2$  term because later, we will choose  $\eta$  to be proportional to  $\varepsilon$  so that it cancels out one of the  $\varepsilon$  terms in that norm. Collecting terms and dividing gives us

$$\begin{aligned} \mathbf{E} \left[ \|w_{k,t} - w^*\|^2 \right] &\leq \frac{1}{1 - \eta} \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 + \frac{1}{\eta} \|w_{k,t} - v_{k,t}\|^2 \right] \\ &= \mathbf{E} \left[ (1 + \eta) \|v_{k,t} - w^*\|^2 + \frac{1}{\eta} \|w_{k,t} - v_{k,t}\|^2 \right], \end{aligned}$$

where the transformation in the last line applies because we will set  $\eta = O(\varepsilon)$ , and so up to terms in  $O(\varepsilon^2)$ , it will hold that

$$\frac{1}{1-\eta} = 1 + \eta;$$

this factor of  $1 + \eta$  disappears from the last term above because that term was already  $O(\varepsilon)$ .

Now, the term  $v_{k,t}$  is exactly what vanilla (non-quantized) SVRG would update its model to. As a result, we can bound  $\|v_{k,t} - w^*\|^2$  using exactly the argument used on page 4 of the original SVRG paper Johnson and Zhang [2]. This straightforward derivation will give

$$\mathbf{E} \left[ \|v_{k,t} - w^*\|^2 \middle| w_{k,t-1} \right] \leq \|w_{k,t-1} - w^*\|^2 - 2\alpha(1 - 2L\alpha) (f(w_{k,t-1}) - f(w^*)) + 4\alpha^2 L (f(o_{w,k}) - f(w^*)).$$

As a result,

$$\begin{aligned} & (1 + \eta) \mathbf{E} \left[ \|v_{k,t} - w^*\|^2 \middle| w_{k,t-1} \right] \\ & \leq (1 + \eta) \|w_{k,t-1} - w^*\|^2 - 2(1 + \eta)\alpha(1 - 2L\alpha) (f(w_{k,t-1}) - f(w^*)) \\ & \quad + 4(1 + \eta)\alpha^2 L (f(o_{w,k}) - f(w^*)) \\ & \leq \|w_{k,t-1} - w^*\|^2 + \eta \|w_{k,t-1} - w^*\|^2 \\ & \quad - 2(1 + \eta)\alpha(1 - 2L\alpha) (f(w_{k,t-1}) - f(w^*)) + 4(1 + \eta)\alpha^2 L (f(o_{w,k}) - f(w^*)) \\ & \leq \|w_{k,t-1} - w^*\|^2 + \eta \frac{2}{\mu} (f(w_{k,t-1}) - f(w^*)) \\ & \quad - 2(1 + \eta)\alpha(1 - 2L\alpha) (f(w_{k,t-1}) - f(w^*)) + 4(1 + \eta)\alpha^2 L (f(o_{w,k}) - f(w^*)) \\ & = \|w_{k,t-1} - w^*\|^2 - 2 \left( (1 + \eta)\alpha(1 - 2L\alpha) - \eta\mu^{-1} \right) (f(w_{k,t-1}) - f(w^*)) \\ & \quad + 4(1 + \eta)\alpha^2 L (f(o_{w,k}) - f(w^*)). \end{aligned}$$

On the other hand, since  $(x + y)^2 \leq 2x^2 + 2y^2$ ,

$$\begin{aligned} \|w_{k,t} - v_{k,t}\|^2 & \leq ((2 + 3\alpha L + C) \cdot \|\delta_{w,k,t-1}\| \cdot \varepsilon_{\text{machine-lo}} + 4\|\bar{g}_k\| \cdot \varepsilon_{\text{machine-lo}})^2 \\ & \leq ((2 + 3\alpha L + C) \cdot \|w_{k,t-1} - o_{w,k}\| \cdot \varepsilon_{\text{machine-lo}} + 4\alpha L \|o_{w,k} - w^*\| \cdot \varepsilon_{\text{machine-lo}})^2 \\ & \leq ((2 + 3\alpha L + C) \cdot (\|w_{k,t-1} - w^*\| + \|o_{w,k} - w^*\|) \cdot \varepsilon_{\text{machine-lo}} + 4\alpha L \|o_{w,k} - w^*\| \cdot \varepsilon_{\text{machine-lo}})^2 \\ & \leq ((2 + 3\alpha L + C) \cdot \|w_{k,t-1} - w^*\| \cdot \varepsilon_{\text{machine-lo}} + (2 + 7\alpha L + C) \cdot \|o_{w,k} - w^*\| \cdot \varepsilon_{\text{machine-lo}})^2 \\ & \leq (2 + 7\alpha L + C)^2 \cdot \varepsilon_{\text{machine-lo}}^2 \cdot \left( 2\|w_{k,t-1} - w^*\|^2 + 2\|o_{w,k} - w^*\|^2 \right) \\ & \leq 4(2 + 7\alpha L + C)^2 \mu^{-1} \cdot \varepsilon_{\text{machine-lo}}^2 \cdot (f(w_{k,t-1}) - f(w^*) + f(o_{w,k}) - f(w^*)), \end{aligned}$$

where the last line follows from strong convexity of  $f$ . Combining our bounds together gives us

$$\begin{aligned} \mathbf{E} \left[ \|w_{k,t} - w^*\|^2 \right] & \leq \mathbf{E} \left[ (1 + \eta) \|v_{k,t} - w^*\|^2 + \frac{1}{\eta} \|w_{k,t} - v_{k,t}\|^2 \right] \\ & \leq \mathbf{E} \left[ \|w_{k,t-1} - w^*\|^2 - 2 \left( (1 + \eta)\alpha(1 - 2L\alpha) - \eta\mu^{-1} \right) (f(w_{k,t-1}) - f(w^*)) \right. \\ & \quad + 4(1 + \eta)\alpha^2 L (f(o_{w,k}) - f(w^*)) \\ & \quad \left. + \frac{4(2 + 7\alpha L + C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2 \cdot (f(w_{k,t-1}) - f(w^*) + f(o_{w,k}) - f(w^*)) \right] \\ & = \mathbf{E} \left[ \|w_{k,t-1} - w^*\|^2 - \Phi \cdot (f(w_{k,t-1}) - f(w^*)) + \Psi \cdot (f(o_{w,k}) - f(w^*)) \right] \end{aligned}$$

where

$$\Phi = 2\alpha(1 - 2L\alpha) + 2\eta (\alpha(1 - 2L\alpha) - \mu^{-1}) - \frac{4(2 + 7\alpha L + C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2$$

and

$$\Psi = 4\alpha^2 L + 4\eta\alpha^2 L + \frac{4(2 + 7\alpha L + C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2.$$

Now summing up over all the iterations of an epoch gives us

$$\mathbf{E} \left[ \|w_{k,T} - w^*\|^2 \right] \leq \mathbf{E} \left[ \|w_{k,0} - w^*\|^2 - \Phi \cdot \sum_{t=1}^T (f(w_{k,t-1}) - f(w^*)) + \Psi \cdot T \cdot (f(o_{w,k}) - f(w^*)) \right],$$

and by the definition of  $o_{w,k+1}$ ,

$$\mathbf{E} \left[ \frac{1}{T} \sum_{t=1}^T (f(w_{k,t-1}) - f(w^*)) \right] = \mathbf{E} [(f(o_{w,k+1}) - f(w^*))],$$

so

$$\mathbf{E} \left[ \|w_{k,T} - w^*\|^2 \right] \leq \mathbf{E} \left[ \|w_{k,0} - w^*\|^2 - \Phi \cdot T \cdot (f(o_{w,k+1}) - f(w^*)) + \Psi \cdot T \cdot (f(o_{w,k}) - f(w^*)) \right].$$

It follows that, since

$$\|w_{k,0} - w^*\|^2 = \|o_{w,k} - w^*\|^2 \leq 2\mu^{-1}(f(o_{w,k}) - f(w^*)),$$

we can conclude that

$$\Phi \cdot T \cdot \mathbf{E} [f(o_{w,k+1}) - f(w^*)] \leq (\Psi T + 2\mu^{-1}) \cdot \mathbf{E} [f(o_{w,k}) - f(w^*)],$$

and so

$$\mathbf{E} [f(o_{w,k+1}) - f(w^*)] \leq \left( \frac{\Psi}{\Phi} + \frac{2}{\mu\Phi T} \right) \cdot \mathbf{E} [f(o_{w,k}) - f(w^*)].$$

Now, ignoring terms of  $\varepsilon^2$ , we can show that

$$\begin{aligned} \frac{\Psi}{\Phi} &= \frac{4\alpha^2 L + 4\eta\alpha^2 L + \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{2\alpha(1-2L\alpha) + 2\eta(\alpha(1-2L\alpha) - \mu^{-1}) - \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2} \\ &= \frac{4\alpha^2 L}{2\alpha(1-2L\alpha)} \left( 1 + \frac{4\eta\alpha^2 L + \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{4\alpha^2 L} - \frac{2\eta(\alpha(1-2L\alpha) - \mu^{-1}) - \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{2\alpha(1-2L\alpha)} \right) \\ &= \frac{4\alpha^2 L}{2\alpha(1-2L\alpha)} \left( 1 + \frac{\frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{4\alpha^2 L} + \frac{2\eta\mu^{-1} + \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{2\alpha(1-2L\alpha)} \right) \\ &= \frac{4\alpha^2 L}{2\alpha(1-2L\alpha)} \left( 1 + \left( \frac{1}{2\alpha L} + \frac{1}{1-2L\alpha} \right) \cdot \frac{2(2+7\alpha L+C)^2}{\alpha\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2 + \frac{\eta}{\alpha\mu(1-2L\alpha)} \right) \\ &= \frac{4\alpha^2 L}{2\alpha(1-2L\alpha)} \left( 1 + \left( \frac{1}{2\alpha L(1-2L\alpha)} \right) \cdot \frac{2(2+7\alpha L+C)^2}{\alpha\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2 + \frac{\eta}{\alpha\mu(1-2L\alpha)} \right) \\ &= \frac{2\alpha L}{1-2L\alpha} \left( 1 + \frac{(2+7\alpha L+C)^2}{\alpha^2\mu L(1-2L\alpha)\eta} \cdot \varepsilon_{\text{machine-lo}}^2 + \frac{\eta}{\alpha\mu(1-2L\alpha)} \right) \\ &= \frac{2\alpha L}{1-2L\alpha} + \frac{2 \cdot (2+7\alpha L+C)^2}{\alpha\mu(1-2L\alpha)^2\eta} \cdot \varepsilon_{\text{machine-lo}}^2 + \frac{2\eta L}{\mu(1-2L\alpha)^2} \end{aligned}$$

Now, suppose we set

$$\eta = \frac{(2+7\alpha L+C)}{\sqrt{\alpha L}} \cdot \varepsilon_{\text{machine-lo}}.$$

Then,

$$\frac{\Psi}{\Phi} = \frac{2\alpha L}{1-2L\alpha} + \frac{4 \cdot (2+7\alpha L+C)}{\mu(1-2L\alpha)^2} \sqrt{\frac{L}{\alpha}} \cdot \varepsilon_{\text{machine-lo}}.$$

Similarly,

$$\begin{aligned}
\frac{2}{\mu\Phi T} &= \frac{2}{\mu T \left( 2\alpha(1-2L\alpha) + 2\eta(\alpha(1-2L\alpha) - \mu^{-1}) - \frac{4(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2 \right)} \\
&= \frac{1}{\mu\alpha(1-2L\alpha)T} \left( 1 - \frac{\eta(\alpha(1-2L\alpha) - \mu^{-1}) - \frac{2(2+7\alpha L+C)^2}{\mu\eta} \cdot \varepsilon_{\text{machine-lo}}^2}{\alpha(1-2L\alpha)} \right) \\
&= \frac{1}{\mu\alpha(1-2L\alpha)T} \left( 1 - \eta + \frac{\eta}{\alpha\mu(1-2L\alpha)} + \frac{2(2+7\alpha L+C)^2 \cdot \varepsilon_{\text{machine-lo}}^2}{\alpha\mu\eta(1-2L\alpha)} \right) \\
&\leq \frac{1}{\mu\alpha(1-2L\alpha)T} + \frac{\eta}{\alpha^2\mu^2(1-2L\alpha)^2T} + \frac{2 \cdot (2+7\alpha L+C)^2 \cdot \varepsilon_{\text{machine-lo}}^2}{\alpha^2\mu^2\eta(1-2L\alpha)^2T} \\
&= \frac{1}{\mu\alpha(1-2L\alpha)T} + \frac{(2+7\alpha L+C)}{\alpha^2\mu^2(1-2L\alpha)^2T\sqrt{\alpha L}} \cdot \varepsilon_{\text{machine-lo}} + \frac{2 \cdot (2+7\alpha L+C)}{\alpha\mu^2(1-2L\alpha)^2T} \sqrt{\frac{L}{\alpha}} \cdot \varepsilon_{\text{machine-lo}} \\
&= \frac{1}{\mu\alpha(1-2L\alpha)T} + \frac{2+7\alpha L+C}{\alpha\mu^2(1-2L\alpha)^2T} \sqrt{\frac{L}{\alpha}} \left( \frac{1}{\alpha L} + 2 \right) \cdot \varepsilon_{\text{machine-lo}}.
\end{aligned}$$

Putting these bounds together gives us

$$\begin{aligned}
\frac{\Psi}{\Phi} + \frac{2}{\mu\Phi T} &\leq \frac{2\alpha L}{1-2L\alpha} + \frac{1}{\mu\alpha(1-2L\alpha)T} + \frac{4 \cdot (2+7\alpha L+C)}{\mu(1-2L\alpha)^2} \sqrt{\frac{L}{\alpha}} \cdot \varepsilon_{\text{machine-lo}} \\
&\quad + \frac{2+7\alpha L+C}{\alpha\mu^2(1-2L\alpha)^2T} \sqrt{\frac{L}{\alpha}} \left( \frac{1}{\alpha L} + 2 \right) \cdot \varepsilon_{\text{machine-lo}} \\
&\leq \frac{2\alpha L}{1-2L\alpha} + \frac{1}{\mu\alpha(1-2L\alpha)T} + \frac{2+7\alpha L+C}{\mu(1-2L\alpha)^2} \sqrt{\frac{L}{\alpha}} \left( 4 + \frac{1}{\alpha^2\mu L T} + \frac{2}{\alpha\mu T} \right) \cdot \varepsilon_{\text{machine-lo}}.
\end{aligned}$$

Next, we substitute the assigned values for  $\alpha$  and  $T$ ,

$$\alpha = \frac{\gamma}{4L(1+\gamma)} \quad T = \frac{8\kappa(1+\gamma)}{\gamma^2}.$$

First, notice that

$$\begin{aligned}
\frac{2\alpha L}{1-2L\alpha} + \frac{1}{\mu\alpha(1-2L\alpha)T} &= \frac{2L \frac{\gamma}{4L(1+\gamma)}}{1-2L \frac{\gamma}{4L(1+\gamma)}} + \frac{1}{\mu \frac{\gamma}{4L(1+\gamma)} (1-2L \frac{\gamma}{4L(1+\gamma)}) \cdot \frac{8\kappa(1+\gamma)}{\gamma^2}} \\
&= \frac{\gamma}{2+\gamma} + \frac{\gamma(1+\gamma)}{(2+\gamma)^2} \\
&= \gamma.
\end{aligned}$$

For the first part of the error term,

$$\begin{aligned}
\frac{2+7\alpha L+C}{\mu(1-2L\alpha)^2} \sqrt{\frac{L}{\alpha}} &= \frac{2+7L \frac{\gamma}{4L(1+\gamma)} + C}{\mu(1-2L \frac{\gamma}{4L(1+\gamma)})^2} \sqrt{\frac{L}{\frac{\gamma}{4L(1+\gamma)}}} \\
&= \frac{2+7 \frac{\gamma}{4(1+\gamma)} + C}{\mu(\frac{2+\gamma}{2(1+\gamma)})^2} \sqrt{\frac{4L^2(1+\gamma)}{\gamma}} \\
&= \frac{8(1+\gamma)^2 + 7\gamma(1+\gamma) + C(1+\gamma)^2}{\mu(2+\gamma)^2} \cdot L \sqrt{\frac{4(1+\gamma)}{\gamma}} \\
&\leq \frac{16+7+2C}{2} \cdot \kappa \sqrt{\frac{8}{\gamma}},
\end{aligned}$$



where the last line follows from the fact that  $0 < \gamma < 1$ . Simplifying further,

$$\frac{2 + 7\alpha L + C}{\mu(1 - 2L\alpha)^2} \sqrt{\frac{L}{\alpha}} \leq (23 + 2C) \cdot \kappa \sqrt{\frac{2}{\gamma}}.$$

For the second part of the error term,

$$\begin{aligned} 4 + \frac{1}{\alpha^2 \mu L T} + \frac{2}{\alpha \mu T} &= 4 + \frac{1}{\left(\frac{\gamma}{4L(1+\gamma)}\right)^2 \mu L \frac{8\kappa(1+\gamma)}{\gamma^2}} + \frac{2}{\frac{\gamma}{4L(1+\gamma)} \mu \frac{8\kappa(1+\gamma)}{\gamma^2}} \\ &= 4 + 2(1 + \gamma) + \gamma \\ &\leq 9, \end{aligned}$$

where again this last line follows from the fact that  $0 < \gamma < 1$ . So, our whole error term is bounded by

$$\frac{2 + 7\alpha L + C}{\mu(1 - 2L\alpha)^2} \sqrt{\frac{L}{\alpha}} \left(4 + \frac{1}{\alpha^2 \mu L T} + \frac{2}{\alpha \mu T}\right) \cdot \varepsilon_{\text{machine-lo}} \leq 9 \cdot (23 + 2C) \cdot \kappa \cdot \sqrt{\frac{2}{\gamma}} \cdot \varepsilon_{\text{machine-lo}}.$$

Substituting this back into our original bound on the convergence of BC-SVRG,

$$\mathbf{E} [f(o_{w,k+1}) - f(w^*)] \leq \left(\gamma + 9 \cdot (23 + 2C) \cdot \kappa \cdot \sqrt{\frac{2}{\gamma}} \cdot \varepsilon_{\text{machine-lo}}\right) \cdot \mathbf{E} [f(o_{w,k}) - f(w^*)].$$

Which is what we wanted to prove. □

## References

- [1] David Bindel. CS 6210 course notes: Dot products in floating point. <http://www.cs.cornell.edu/~bindel/class/cs6210-f12/notes/lec05.pdf>, 2012.
- [2] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- [3] Lloyd N Trefethen and David Bau III. *Numerical linear algebra*, volume 50. Siam, 1997.