

Randomization of Approximate Bilinear Computation for Matrix Multiplication

Osman Asif Malik and Stephen Becker

Department of Applied Mathematics, University of Colorado Boulder, USA

ARTICLE HISTORY

Compiled January 11, 2022

ABSTRACT

We present a method for randomizing formulas for bilinear computation of matrix products. We consider the implications of such randomization when there are two sources of error: One due to the formula itself only being approximately correct, and one due to using floating point arithmetic. Our theoretical results and numerical experiments indicate that our method can improve performance when each of these error sources are present individually, as well as when they are present at the same time.

KEYWORDS

randomized algorithms; approximate algorithms; matrix multiplication; Strassen's algorithm

1. Introduction

Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. In this paper, we are concerned with formulas for computing $\mathbf{C} = \mathbf{AB}$ that take the form

$$c_{ij} = \sum_{r=1}^R w_{ijr} \left(\sum_{k,l=1}^n u_{klr} a_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'r} b_{k'l'} \right), \quad (1)$$

$i, j \in [n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$, where each $\mathbf{u} = (u_{klr})$, $\mathbf{v} = (v_{k'l'r})$, and $\mathbf{w} = (w_{ijr})$ is a tensor containing real numbers, and c_{ij} is the element at position (i, j) in \mathbf{C} , with similar notation for elements of \mathbf{A} and \mathbf{B} . Such a formula is called a *bilinear computation* (BC) [7]. We can rewrite the expression above as

$$c_{ij} = \sum_{k,l=1}^n \sum_{k',l'=1}^n a_{kl} b_{k'l'} \sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr}.$$

Comparing this to the standard algorithm for matrix multiplication, we can see that for the computation (1) to be exact, \mathbf{U} , \mathbf{V} and \mathbf{W} must satisfy

$$\sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} = \delta_{ki} \delta_{l'j} \delta_{lk'} \quad \text{for all } (k, l, k', l', i, j) \in [n]^6, \quad (2)$$

where δ is the Kronecker delta (i.e., $\delta_{ki} = 1$ if $k = i$ and 0 otherwise) [8]. If (2) is satisfied, we say that (1) is an *exact bilinear computation* (EBC). The smallest positive integer R for which there exist \mathbf{U} , \mathbf{V} and \mathbf{W} such that (2) holds is referred to as the *rank* of the computation. We refer to any algorithm of the form (1) for which $R < n^3$ as *fast*, since the asymptotic complexity is smaller than that of standard matrix multiplication, which has complexity $O(n^3)$. Examples of fast computations of the form (1) include Strassen's algorithm for 2×2 matrices [30], and Laderman's algorithm for 3×3 matrices [25]. The formula (1) is also valid if \mathbf{A} , \mathbf{B} and \mathbf{C} are of size $mn \times mn$ and a_{kl} , $b_{k'l'}$ and c_{ij} are replaced by submatrices \mathbf{A}_{kl} , $\mathbf{B}_{k'l'}$ and \mathbf{C}_{ij} of size $m \times m$, so that e.g.

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1n} \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2n} \\ \vdots & \vdots & & \vdots \\ \mathbf{A}_{n1} & \mathbf{A}_{n2} & \cdots & \mathbf{A}_{nn} \end{bmatrix}.$$

This is why fast algorithms of the form (1) can be used recursively to compute the product of larger matrices.

Define 6-way tensors

$$\mathbf{X} \stackrel{\text{def}}{=} (\delta_{ki} \delta_{l'j} \delta_{lk'})_{(k,l,k',l',i,j) \in [n]^6}, \quad (3)$$

$$\mathbf{Y} \stackrel{\text{def}}{=} \left(\sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} \right)_{(k,l,k',l',i,j) \in [n]^6}. \quad (4)$$

The condition in (2) can be written succinctly as $\mathbf{Y} = \mathbf{X}$. For both matrix and tensor inputs, let $\|\cdot\|$ denote the Frobenius norm, i.e., the square root of the sum of the square of all elements. We are interested in BCs that are only approximately correct, which motivates the following definition.

Definition 1 (Approximate bilinear computation). Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{mn \times mn}$, where m is a positive integer, and let f be defined blockwise via

$$f(\mathbf{A}, \mathbf{B})_{ij} \stackrel{\text{def}}{=} \sum_{r=1}^R w_{ijr} \left(\sum_{k,l=1}^n u_{klr} \mathbf{A}_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'r} \mathbf{B}_{k'l'} \right) \in \mathbb{R}^{m \times m} \quad \text{for } i, j \in [n], \quad (5)$$

where $f(\mathbf{A}, \mathbf{B})_{ij}$, \mathbf{A}_{kl} and $\mathbf{B}_{k'l'}$ are submatrices of $f(\mathbf{A}, \mathbf{B})$, \mathbf{A} and \mathbf{B} , respectively, of size $m \times m$. We say that f is an *approximate bilinear computation* (ABC) with parameters n and τ , or (n, τ) -ABC for short, if $\|\mathbf{Y} - \mathbf{X}\| \leq \tau$.

We present a method for randomizing the computation (5) and consider the implications of this approach when there are two kinds of error present.

- **Error due to an approximate algorithm.** Using an (n, τ) -ABC with $\tau > 0$ will introduce

error even if exact arithmetic is used. ABCs are interesting since it is sometimes hard or impossible to convert an approximate algorithm found numerically to an exact algorithm. Some examples of such approximate algorithms found numerically appear in [29]. Other papers that search for fast algorithms numerically include [2, 8, 17, 23].

- **Numerical error due to using floating point arithmetic.** Although the standard algorithm for matrix multiplication also incurs numerical error, it is more severe in fast algorithms since the computation of an element c_{ij} can involve elements from \mathbf{A} and \mathbf{B} other than the vectors \mathbf{a}_i and \mathbf{b}_j , which are the i th row of \mathbf{A} and j th column of \mathbf{B} , respectively. In exact arithmetic, these additional terms cancel out for exact fast algorithms, but in finite precision those cancellations are typically not exact which can lead to substantial numerical error. These issues are exacerbated when the algorithm itself is only approximately correct. These considerations are especially important when using low precision environments, such as when computing on a GPU [21]. Low precision computation—using 32-bit, 16-bit, and even 8-bit precision numbers—is popular in, e.g., machine learning [19, 20, 31].

We make the following contributions in this paper:

- We propose a method for randomizing BCs for matrix multiplication which does not increase the leading order computational complexity of the algorithm.
- When exact arithmetic is used, we show that our randomized ABCs compute the correct matrix product in expectation. We also provide some performance guarantees.
- We show that these exact arithmetic results largely carry over to a setting when all computations are done in floating point arithmetic.
- When floating point arithmetic is used, we provide numerical evidence that randomizing EBCs using our scheme can reduce numerical error and improve robustness to adversarial examples.

1.1. Related work

Bini et al. [4, 6] introduce a concept similar to ABC called Arbitrary Precision Approximating (APA) algorithms for matrix multiplication. For an APA algorithm, the tensor \mathbf{X} and \mathbf{Y} defined in (3) and (4) satisfy the relationship

$$\mathbf{Y}(\varepsilon) + \mathbf{E}(\varepsilon) = \mathbf{X},$$

where

$$\mathbf{Y}(\varepsilon) \stackrel{\text{def}}{=} \left(\sum_{r=1}^R u_{klr}(\varepsilon) v_{k'l'r}(\varepsilon) w_{ijr}(\varepsilon) \right)_{(k,l,k',l',i,j) \in [n]^6} \quad (6)$$

is a function of ε , and \mathbf{E} represents the error in the approximation $\mathbf{Y} \approx \mathbf{X}$. Each entry in \mathbf{E} is assumed to be a polynomial with zero constant coefficient, hence for every entry $\mathbf{E}_{klk'l'ij}(0) = 0$. Consequently, an APA algorithm can be made arbitrarily accurate by making ε small enough. Moreover, as shown in [3], an EBC can be derived from an APA algorithm (see also the related discussion in Section 15.2 of [10]). This EBC takes the form

$$\sum_{i=1}^{d+1} \alpha_i \mathbf{Y}(\varepsilon_i) = \mathbf{X}, \quad (7)$$

where all ε_i are distinct, $d := \max_{klk'lj \in [n]^6} \deg(\mathbf{E}_{klk'lj}(\varepsilon))$ is the maximum degree of the polynomials describing the entries in \mathbf{E} , and where $[\alpha_1, \dots, \alpha_{d+1}]$ is chosen as the solution to a certain linear system. Bini [3] uses this approach to convert the APA scheme for 12×12 matrix multiplication in [4] to an exact one, which can be used recursively to get a matrix multiplication algorithm with complexity $O(n^{2.7799})$. Fixing the error in an APA algorithm yields an ABC with some parameter τ . In some of our numerical experiments, we use an ABC which we get by fixing ε in the 12×12 APA scheme of [4]. We also consider an associated EBC derived via (7) in other experiments. Our paper focuses on ABCs instead of APA algorithms since, in practice, it may be difficult or impossible to express a pair (\mathbf{Y}, \mathbf{E}) found numerically in terms of polynomials in ε .

To the best of our knowledge, our paper is the first to consider randomization as a tool for improving BCs for matrix multiplication which are only approximately correct. Various randomized algorithms for the standard matrix multiplication algorithm have been considered in other works; see e.g. [15, 26]. For EBCs in floating point arithmetic, a patent by Castrapel and Gustafson [11] describes a randomized version of Strassen’s algorithm which they claim reduces numerical error. They provide empirical support for this, but no mathematical proof. Our method generalizes their approach by randomly choosing from a wider range of equivalent algorithms. Additionally, our method can be applied to any formula of the form (1).

Early works that analyze the stability of fast algorithms for matrix multiplication include [5, 6, 9]. Other works, such as [1, 12, 14, 16, 24], attempt to improve the stability of fast algorithms using other approaches that do not rely on randomization. In Section 3.2, we compare our proposed method to the rescaling method in [1]. In these experiments, we also consider a restricted version of our method which corresponds to the method in [11].

2. Randomization of bilinear computation for matrix multiplication

In Section 2.1, we first consider a setting in which exact arithmetic is used. In Section 2.2, we then consider a setting in which floating point arithmetic is used.

2.1. In exact arithmetic

We present our randomization scheme in the setting when the input matrices are $mn \times mn$ with $m \geq 1$ ¹. Accordingly, suppose $f : \mathbb{R}^{mn \times mn} \times \mathbb{R}^{mn \times mn} \rightarrow \mathbb{R}^{mn \times mn}$ is an (n, τ) -ABC. We will now define a randomized version of f , denoted by \hat{f} , which has the following property: For all $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{mn \times mn}$, $\mathbb{E}[\hat{f}(\mathbf{A}, \mathbf{B})] = \mathbf{A}\mathbf{B}$. To that end, let $\{s_i(j)\}_{(i,j) \in [3] \times [n]}$ be a collection of i.i.d. Rademacher random variables, i.e., each satisfying $\mathbb{P}[s_i(j) = +1] = \mathbb{P}[s_i(j) = -1] = 1/2$. Moreover, let $\pi_i : [n] \rightarrow [n]$, $i \in [3]$, be independent random permutation functions, each satisfying $(\forall (j, k) \in [n]^2) \mathbb{P}[\pi_i(j) = k] = 1/n$. Let $\mathbf{S}_i \in \mathbb{R}^{mn \times mn}$, $i \in [3]$, be block diagonal matrices with the j th nonzero block equal to $s_i(j)\mathbf{I}_m$, where \mathbf{I}_m is the $m \times m$ identity matrix. Also, let $\mathbf{P}_i \in \mathbb{R}^{mn \times mn}$, $i \in [3]$, be permutation matrices divided into $m \times m$ blocks, with blocks on position $(\pi_i(j), j)$, $j \in [n]$, equal to \mathbf{I}_m and all other blocks equal to zero. Define $\mathbf{M}_i = \mathbf{P}_i\mathbf{S}_i$, $i \in [3]$. Note that each \mathbf{M}_i is orthogonal (i.e., $\mathbf{M}_i^{-1} = \mathbf{M}_i^\top$). We propose the following definition of \hat{f} .

Definition 2 (Randomized approximate bilinear computation). Let f be an (n, τ) -ABC. We define a corresponding *randomized approximate bilinear computation* with parameters n, τ and κ , or (n, τ, κ) -

¹We present results for square matrices, but we believe the results can be extended to rectangular matrices.

RandABC for short, via

$$\hat{f}(\mathbf{A}, \mathbf{B}) \stackrel{\text{def}}{=} (1 - \kappa)^{-1} \mathbf{M}_1^\top f(\mathbf{M}_1 \mathbf{A} \mathbf{M}_2^\top, \mathbf{M}_2 \mathbf{B} \mathbf{M}_3^\top) \mathbf{M}_3, \quad (8)$$

where

$$\kappa \stackrel{\text{def}}{=} \frac{1}{n^3} \sum_{(i,j,l) \in [n]^3} \left(1 - \sum_{r=1}^R u_{ilr} v_{ljr} w_{ijr} \right) \quad (9)$$

is assumed to satisfy $\kappa \neq 1$.

Observe that if f was an exact algorithm, then $\kappa = 0$ and $\hat{f}(\mathbf{A}, \mathbf{B}) = \mathbf{A}\mathbf{B}$ since the \mathbf{M}_i 's would cancel out due to orthogonality. Since the cost of applying \mathbf{M}_i to an $mn \times mn$ matrix is $O(m^2n^2)$, computing $\hat{f}(\mathbf{A}, \mathbf{B})$ has the same leading order complexity as computing $f(\mathbf{A}, \mathbf{B})$.

Proposition 3. *Let \hat{f} be an (n, τ, κ) -RandABC with $\tau \geq 0$ and $\kappa \neq 1$. For all $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{mn \times mn}$ we have $\mathbb{E}[\hat{f}(\mathbf{A}, \mathbf{B})] = \mathbf{A}\mathbf{B}$.*

Proof. Let $\hat{\mathbf{C}} \stackrel{\text{def}}{=} \hat{f}(\mathbf{A}, \mathbf{B})$. Considering the (i, j) th block of $\hat{\mathbf{C}}$, and going through some tedious but straightforward algebra, we get

$$\hat{\mathbf{C}}_{ij} = \frac{s_1(i)s_3(j)}{1 - \kappa} \sum_{k,l=1}^n \sum_{k',l'=1}^n s_1(k)s_2(l)s_2(k')s_3(l') \mathbf{A}_{kl} \mathbf{B}_{k'l'} \sum_{r=1}^R u_{\pi_1(k)\pi_2(l)r} v_{\pi_2(k')\pi_3(l')r} w_{\pi_1(i)\pi_3(j)r}. \quad (10)$$

If we take the expectation of this equation with respect to the random variables $\{s_i(j)\}$, most terms will vanish: If $i = k$, $l = k'$ and $j = l'$ for a given term, then the product of the s_i 's will be 1; otherwise, the expectation of that term will be zero due to independence and the fact that each $\mathbb{E}[s_i(j)] = 0$. Consequently, we have

$$\mathbb{E}[\hat{\mathbf{C}}_{ij}] = (1 - \kappa)^{-1} \sum_{l=1}^n \mathbf{A}_{il} \mathbf{B}_{lj} \mathbb{E} \left[\sum_{r=1}^R u_{\pi_1(i)\pi_2(l)r} v_{\pi_2(l)\pi_3(j)r} w_{\pi_1(i)\pi_3(j)r} \right] = \sum_{l=1}^n \mathbf{A}_{il} \mathbf{B}_{lj} = (\mathbf{A}\mathbf{B})_{ij},$$

where the second equality is true since (9) implies that

$$\mathbb{E} \left[\sum_{r=1}^R u_{\pi_1(i)\pi_2(l)r} v_{\pi_2(l)\pi_3(j)r} w_{\pi_1(i)\pi_3(j)r} \right] = 1 - \kappa.$$

□

It may seem surprising that Proposition 3 holds for *any* $\tau \geq 0$, since this means that even an ABC with an arbitrarily large error will be correct in expectation once randomized as in Definition 2. However, as we will see in Proposition 4, in order to be able to guarantee a small error for *any realization* of $\hat{f}(\mathbf{A}, \mathbf{B})$, τ also needs to be small.

Define $B_\mu^{(n)} \stackrel{\text{def}}{=} \{\mathbf{M} \in \mathbb{R}^{n \times n} : \|\mathbf{M}\| \leq \mu\}$ and $\eta \stackrel{\text{def}}{=} (1 - \kappa)^{-1} - 1 = O(\kappa)$. The following proposition provides performance guarantees for f and \hat{f} .

Proposition 4. Consider matrices $\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}$.

(i) If f is an (n, τ) -ABC, then $\|f(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| \leq \mu^2 \tau$.

(ii) If \hat{f} is an (n, τ, κ) -RandABC, then $\|\hat{f}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| \leq \mu^2 |\eta| \|\mathbf{Y}\| + \mu^2 \tau \lesssim \mu^2 |\kappa| \|\mathbf{Y}\| + \mu^2 \tau$.

(iii) Moreover, $\sup_{\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}} \|\hat{f}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| \leq \sup_{\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}} \|f(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| + \mu^2 |\eta| \|\mathbf{Y}\|$.

Proof. We have

$$\begin{aligned}
\|f(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\|^2 &= \sum_{i,j=1}^n \left\| \sum_{k,l=1}^n \sum_{k',l'=1}^n \mathbf{A}_{kl} \mathbf{B}_{k'l'} (y_{klk'l'ij} - x_{klk'l'ij}) \right\|^2 \\
&\leq \sum_{i,j=1}^n \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n \|\mathbf{A}_{kl}\| \|\mathbf{B}_{k'l'}\| |y_{klk'l'ij} - x_{klk'l'ij}| \right)^2 \\
&\leq \sum_{i,j=1}^n \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n (\|\mathbf{A}_{kl}\| \|\mathbf{B}_{k'l'}\|)^2 \right) \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n |y_{klk'l'ij} - x_{klk'l'ij}|^2 \right) \\
&= \left(\sum_{k,l=1}^n \|\mathbf{A}_{kl}\|^2 \sum_{k',l'=1}^n \|\mathbf{B}_{k'l'}\|^2 \right) \left(\sum_{i,j=1}^n \sum_{k,l=1}^n \sum_{k',l'=1}^n |y_{klk'l'ij} - x_{klk'l'ij}|^2 \right) \\
&= \|\mathbf{A}\|^2 \|\mathbf{B}\|^2 \|\mathbf{Y} - \mathbf{X}\|^2 \leq \mu^4 \tau^2,
\end{aligned}$$

where the first inequality follows from first applying the triangle inequality and then using the sub-multiplicativity of the Frobenius norm, and the second inequality follows from applying the Cauchy–Schwarz inequality. This proves (i). Since the Frobenius norm is invariant under unitary transformations,

$$\|\hat{f}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| = \|(1 - \kappa)^{-1} f(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) - \tilde{\mathbf{A}}\tilde{\mathbf{B}}\|, \quad (11)$$

where $\tilde{\mathbf{A}} = \mathbf{M}_1 \mathbf{A} \mathbf{M}_2^\top$, $\tilde{\mathbf{B}} = \mathbf{M}_2 \mathbf{B} \mathbf{M}_3^\top$. Going through the same computations as in the proof of (i), but with the extra $(1 - \kappa)^{-1}$ term, we therefore get

$$\|\hat{f}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| \leq \|\tilde{\mathbf{A}}\| \|\tilde{\mathbf{B}}\| \|(1 - \kappa)^{-1} \mathbf{Y} - \mathbf{X}\| = \|\mathbf{A}\| \|\mathbf{B}\| \|(1 - \kappa)^{-1} \mathbf{Y} - \mathbf{X}\| \leq \mu^2 (|\eta| \|\mathbf{Y}\| + \tau)$$

where the equality once again uses the unitary invariance of the Frobenius norm, and the last inequality follows from the definitions of η and the triangle inequality. This proves (ii). Fix $\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}$. Applying the triangle inequality to (11) and using the definition of η , we get

$$\|\hat{f}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| \leq \|f(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) - \tilde{\mathbf{A}}\tilde{\mathbf{B}}\| + |\eta| \|f(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})\|. \quad (12)$$

By doing computations almost identical to those in the proof of (i), we get the bound

$$\|f(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})\| \leq \|\tilde{\mathbf{A}}\| \|\tilde{\mathbf{B}}\| \|\mathbf{Y}\| \leq \mu^2 \|\mathbf{Y}\|, \quad (13)$$

since $\tilde{\mathbf{A}}, \tilde{\mathbf{B}} \in B_\mu^{(mn)}$. Combining (12) and (13) and taking supremums appropriately proves (iii). \square

Points (i) and (ii) in Proposition 4 provide absolute performance guarantees for f and \hat{f} , respectively. Note that a tighter version of (i) holds, with $\|\mathbf{A}\|$ replaced by $\sqrt{\sum_{kl} \|\mathbf{A}_{kl}\|_2^2}$, where $\|\cdot\|_2$ is the

spectral norm. We keep (i) in its current looser form to make it easier to compare to (ii) and (iii). Point (iii) shows that the worst case performance of \hat{f} is no worse than that of f plus a constant. In fact, that constant, which also appears in (ii), can be bounded as follows.

Proposition 5. *If $|\kappa| \leq 1/2$, then*

$$|\eta|\mu^2\|\mathbf{Y}\| \leq 2\mu^2(n^{-5/2}\tau + n^{-1})\tau.$$

Proof. Let $\mathbf{z} \in \mathbb{R}^{n^3}$ be the vector with elements $(\sum_{r=1}^R u_{ilr}v_{ljr}w_{ijr} - 1)$ for $(i, j, l) \in [n]^3$, i.e., containing the elements of $(\mathbf{Y} - \mathbf{X})$ in positions for which \mathbf{X} has an entry 1 (the element order in \mathbf{z} is irrelevant). Note that

$$|\kappa| = \left| n^{-3} \sum_{i=1}^{n^3} z_i \right| \leq n^{-3} \|\mathbf{z}\|_1 \leq n^{-3} \sqrt{n} \|\mathbf{z}\| \leq n^{-5/2} \|\mathbf{Y} - \mathbf{X}\|, \quad (14)$$

where the first equality follows from (9), the first inequality follows from the triangle inequality and the definition of the 1-norm, and the second inequality is a well known relation (see e.g. Equation (2.2.5) in [18]). Now, note that

$$\|\mathbf{Y}\| \leq \|\mathbf{Y} - \mathbf{X}\| + \|\mathbf{X}\| = \|\mathbf{Y} - \mathbf{X}\| + n^{3/2}. \quad (15)$$

Combining (14), (15), the fact that $|\eta| \leq 2|\kappa|$ when $|\kappa| \leq 1/2$, and the definition of τ gives us the desired bound. \square

The upper bound $\mu^2\tau$ in Proposition 4 (i) is also an upper bound to the sup term on the right hand side of the inequality in (iii) of the same proposition. Proposition 5 shows that the size of the additional constant $|\eta|\mu^2\|\mathbf{Y}\|$ is not much larger than the bound that we already have on this sup term, and that it will be smaller than that bound if e.g. $n \geq 3$ and $\tau < 3^{3/2}/2$.

As is clear from the proof of Proposition 3, the constant κ in (9) is used to rescale the computation in (8) so that the output of \hat{f} is correct in expectation. It is important to note that κ , which can be both positive and negative, is *not* a measure of error in the algorithm. Indeed, setting $R = 1$ and all elements of \mathbf{U} , \mathbf{V} and \mathbf{W} to 1 would result in $\kappa = 0$, but this would clearly be a very poor algorithm. Although the corresponding randomized algorithm \hat{f} would be correct in expectation, the error guarantees in Proposition 4 (i) and (ii) would be very poor, since τ would be large.

2.1.1. A recursive algorithm for approximate bilinear computation

We can extend the result in Proposition 3 to a recursive version of \hat{f} . We denote the recursive algorithm with Q recursions for multiplication of $mn^Q \times mn^Q$ matrices by $\hat{F}^{(Q)}$. Let $\{s_i^{(q)}(j)\}_{(i,j,q) \in [3] \times [n] \times [Q]}$ be a collection of i.i.d. Rademacher random variables, and let $\pi_i^{(q)} : [n] \rightarrow [n]$, $(i, q) \in [3] \times [Q]$, be independent random permutation functions, each satisfying $(\forall (j, k) \in [n]^2) \mathbb{P}[\pi_i^{(q)}(j) = k] = 1/n$. Let $\hat{F}^{(1)}$ be defined exactly as \hat{f} in (8) but based on the random variables $\{s_i^{(1)}(j)\}_{(i,j) \in [3] \times [n]}$ and $\{\pi_i^{(1)}\}_{i \in [3]}$ and define $\hat{F}^{(q)} : \mathbb{R}^{mn^q \times mn^q} \times \mathbb{R}^{mn^q \times mn^q} \rightarrow \mathbb{R}^{mn^q \times mn^q}$, $q \in \{2, 3, \dots, Q\}$, recursively via

$$(\hat{F}^{(q)}(\mathbf{A}, \mathbf{B}))_{ij} \stackrel{\text{def}}{=} (1 - \kappa)^{-1} s_1^{(q)}(i) s_3^{(q)}(j) \sum_{r=1}^R w_{\pi_1^{(q)}(i) \pi_3^{(q)}(j)r} \hat{F}^{(q-1)}(\mathbf{A}^{(q)}, \mathbf{B}^{(q)}), \quad (16)$$

where

$$\mathbf{A}^{(q)} \stackrel{\text{def}}{=} \sum_{k,l=1}^n u_{\pi_1^{(q)}(k)\pi_2^{(q)}(l)r} s_1^{(q)}(k) s_2^{(q)}(l) \mathbf{A}_{kl},$$

$$\mathbf{B}^{(q)} \stackrel{\text{def}}{=} \sum_{k',l'=1}^n v_{\pi_2^{(q)}(k')\pi_3^{(q)}(l')r} s_2^{(q)}(k') s_3^{(q)}(l') \mathbf{B}_{k'l'},$$

and $(\hat{F}^{(q)}(\mathbf{A}, \mathbf{B}))_{ij}$ is the subblock of size $mn^{q-1} \times mn^{q-1}$ on position (i, j) . If \mathbf{A} and \mathbf{B} are of size $p \times p$ but there is no integer m such that $p = mn^Q$, e.g. if p is prime, then we can simply pad the matrices appropriately [22]. If the recursive formula (16) is used Q times to compute the product of two $N \times N$ matrices, where $N \stackrel{\text{def}}{=} mn^Q$, and $\hat{F}^{(0)}(\mathbf{A}^{(0)}, \mathbf{B}^{(0)}) = \mathbf{A}^{(0)}\mathbf{B}^{(0)}$ is computed via the standard matrix multiplication formula, the asymptotic cost of the recursive algorithm is $O(m^{3-\log_n R} N^{\log_n R})$. This is the same leading order complexity as the corresponding computation without any randomization. For example, if we insert $m = 1$, $n = 2$ and $R = 7$, we recover the asymptotic cost of Strassen's algorithm: $O(N^{\log_2 7}) \approx O(N^{2.81})$.

Proposition 6. *For any positive integer Q and for all $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{mn^Q \times mn^Q}$ we have $\mathbb{E}[\hat{F}^{(Q)}(\mathbf{A}, \mathbf{B})] = \mathbf{AB}$.*

Proof. Note that the claim is true for $Q = 1$ due to Proposition 3. Now assume it is true for some $Q \geq 1$. We will show that it is also true for $Q + 1$, i.e., that $\mathbb{E}[\hat{F}^{(Q+1)}(\mathbf{A}, \mathbf{B})] = \mathbf{AB}$ for $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{mn^{Q+1} \times mn^{Q+1}}$. Let \mathcal{S} denote the σ -algebra generated by the random variables $\{s_i^{(Q+1)}(j)\}_{(i,j) \in [3] \times [n]}$ and $\{\pi_i^{(Q+1)}(j)\}_{(i,j) \in [3] \times [n]}$, and let $\mathbb{E}_{\mathcal{S}}[\cdot] \stackrel{\text{def}}{=} \mathbb{E}[\cdot | \mathcal{S}]$. Then

$$\mathbb{E}[(\hat{F}^{(Q+1)}(\mathbf{A}, \mathbf{B}))_{ij}] = \mathbb{E} \left[\mathbb{E}_{\mathcal{S}}[(\hat{F}^{(Q+1)}(\mathbf{A}, \mathbf{B}))_{ij}] \right]$$

due to the smoothing property of expectation (property 10 in [27, p. 348]),

$$\begin{aligned} &= \mathbb{E} \left[(1 - \kappa)^{-1} s_1^{(Q+1)}(i) s_3^{(Q+1)}(j) \sum_{r=1}^R w_{\pi_1^{(Q+1)}(i)\pi_3^{(Q+1)}(j)r} \right. \\ &\quad \times \mathbb{E}_{\mathcal{S}} \left[\hat{F}^{(Q)} \left(\sum_{k,l=1}^n u_{\pi_1^{(Q+1)}(k)\pi_2^{(Q+1)}(l)r} s_1^{(Q+1)}(k) s_2^{(Q+1)}(l) \mathbf{A}_{kl}, \right. \right. \\ &\quad \left. \left. \sum_{k',l'=1}^n v_{\pi_2^{(Q+1)}(k')\pi_3^{(Q+1)}(l')r} s_2^{(Q+1)}(k') s_3^{(Q+1)}(l') \mathbf{B}_{k'l'} \right) \right] \end{aligned}$$

since each $s_i^{(Q+1)}(j)$ and $\pi_i^{(Q+1)}(j)$ is \mathcal{S} -measurable,

$$\begin{aligned}
&= \mathbb{E} \left[(1 - \kappa)^{-1} s_1^{(Q+1)}(i) s_3^{(Q+1)}(j) \sum_{r=1}^R w_{\pi_1^{(Q+1)}(i) \pi_3^{(Q+1)}(j) r} \right. \\
&\quad \times \left(\sum_{k,l=1}^n u_{\pi_1^{(Q+1)}(k) \pi_2^{(Q+1)}(l) r} s_1^{(Q+1)}(k) s_2^{(Q+1)}(l) \mathbf{A}_{kl} \right) \\
&\quad \left. \times \left(\sum_{k',l'=1}^n v_{\pi_2^{(Q+1)}(k') \pi_3^{(Q+1)}(l') r} s_2^{(Q+1)}(k') s_3^{(Q+1)}(l') \mathbf{B}_{k'l'} \right) \right]
\end{aligned}$$

due to the induction hypothesis and since all random variables $\{s_i^{(q)}(j)\}_{(i,j,q) \in [3] \times [n] \times [Q]}$ and $\{\pi_i^{(q)}(j)\}_{(i,j,q) \in [3] \times [n] \times [Q]}$ are independent of \mathcal{S} (and using property 12 in [27, pp. 349–350]),

$$\begin{aligned}
&= \mathbb{E} \left[(1 - \kappa)^{-1} s_1^{(Q+1)}(i) s_3^{(Q+1)}(j) \sum_{k,l=1}^n \sum_{k',l'=1}^n s_1^{(Q+1)}(k) s_2^{(Q+1)}(l) s_2^{(Q+1)}(k') s_3^{(Q+1)}(l') \mathbf{A}_{kl} \mathbf{B}_{k'l'} \right. \\
&\quad \left. \times \sum_{r=1}^R u_{\pi_1^{(Q+1)}(k) \pi_2^{(Q+1)}(l) r} v_{\pi_2^{(Q+1)}(k') \pi_3^{(Q+1)}(l') r} w_{\pi_1^{(Q+1)}(i) \pi_3^{(Q+1)}(j) r} \right]
\end{aligned}$$

by reordering the terms,

$$= \sum_{l=1}^n \mathbf{A}_{il} \mathbf{B}_{lj} = (\mathbf{A}\mathbf{B})_{ij}$$

which follows by doing the same analysis as in the proof of Proposition 3. So $\mathbb{E}[\hat{F}^{(Q+1)}(\mathbf{A}, \mathbf{B})] = \mathbf{A}\mathbf{B}$. The claim in Proposition 6 now follows by induction. \square

2.2. In floating point arithmetic

In floating point arithmetic, we cannot achieve guarantees like those in Propositions 3 and 6. This is illustrated in the following example.

Example 7. Let $g : \mathbb{R}^{mn \times mn} \times \mathbb{R}^{mn \times mn} \rightarrow \mathbb{R}^{mn \times mn}$ be a function that computes matrix multiplication according to some EBC in floating point arithmetic. Let \hat{g} be defined analogously to \hat{f} in (8), but in terms of g instead of f . Since g is exact, we have $\kappa = 0$. We will use $\hat{G}^{(Q)}$ to denote the recursive version of \hat{g} , defined analogously to $\hat{F}^{(Q)}$. We will use $G^{(Q)}$ to denote the recursive version of g , defined analogously to $\hat{G}^{(Q)}$ but with each $s_i^{(q)}(j) = 1$ and each $\pi_i^{(q)}(j) = j$, i.e., with no randomness involved so that $G^{(Q)}$ is deterministic. We consider the same setup as in Section 2.7.10 of [18]: Let

$$\mathbf{A} = \mathbf{B} = \begin{bmatrix} 0.99 & 0.0010 \\ 0.0010 & 0.99 \end{bmatrix},$$

and suppose we are computing on a machine using 2-digit floating point arithmetic. Taking g to be Strassen's algorithm computed on this machine, we get $\|g(\mathbf{A}, \mathbf{B}) - \mathbf{A}\mathbf{B}\| \approx 0.0286$. In the definition

of \hat{g} , there are a total of 64 possible sign functions and 8 possible permutation functions. Each combination of these has the same probability of occurring, so we can readily compute $\mathbb{E}[\hat{g}(\mathbf{A}, \mathbf{B})]$. Doing this, we find that $\|\mathbb{E}[\hat{g}(\mathbf{A}, \mathbf{B})] - \mathbf{AB}\| \approx 0.0024$. So we cannot guarantee $\mathbb{E}[\hat{g}(\mathbf{A}, \mathbf{B})] = \mathbf{AB}$ in this finite precision setting. We are not even guaranteed to have $\|\mathbb{E}[\hat{g}(\mathbf{A}, \mathbf{B})] - \mathbf{AB}\| \leq \|g(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\|$: Let $\mathbf{A}^{(r)}$ and $\mathbf{B}^{(r)}$ be the same as \mathbf{A} and \mathbf{B} , respectively, but with the order of the columns reversed, i.e., $\mathbf{A}^{(r)} \stackrel{\text{def}}{=} [\mathbf{a}_{:2} \ \mathbf{a}_{:1}]$ and $\mathbf{B}^{(r)} \stackrel{\text{def}}{=} [\mathbf{b}_{:2} \ \mathbf{b}_{:1}]$. We then get $\|g(\mathbf{A}^{(r)}, \mathbf{B}^{(r)}) - \mathbf{A}^{(r)}\mathbf{B}^{(r)}\| \approx 0.0001$, but $\|\mathbb{E}[\hat{g}(\mathbf{A}^{(r)}, \mathbf{B}^{(r)})] - \mathbf{A}^{(r)}\mathbf{B}^{(r)}\| \approx 0.0024$. Despite this, randomization seems to work remarkably well in practice when an exact algorithm is computed in finite precision arithmetic, as we will see in the numerical experiments.

We now consider ABCs (which include EBCs as a special case) in finite precision arithmetic. As in [18], we use $\text{fl}(x)$ to denote the representation of $x \in \mathbb{R}$ as a floating point number, and $\text{fl}(f(x))$ to denote the result of computing $f(x)$ in floating point arithmetic. When computing the latter, the algorithm used to compute f matters. We make the standard assumption that $\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \Delta)$ where x and y are floating point numbers, op is scalar addition, subtraction or multiplication, and $|\Delta| \leq \varepsilon_{\text{machine}}$, where $\varepsilon_{\text{machine}}$ is the machine epsilon or unit roundoff. For numerical summations, e.g. $\text{fl}(\sum_{n=1}^N x_n)$, we assume that the summation is simply done sequentially, e.g.

$$\text{fl}\left(\sum_{n=1}^3 x_n\right) = \text{fl}\left(\text{fl}\left(\text{fl}(x_1) + \text{fl}(x_2)\right) + \text{fl}(x_3)\right).$$

We first present a series of results in Propositions 8–12 with the goal of understanding how algorithmic *and* numerical error combined impact the deterministic and randomized ABCs described by f and \hat{f} , respectively. Throughout these results, we make the reasonable assumption that the input matrices \mathbf{A} and \mathbf{B} , as well as the tensors \mathbf{U} , \mathbf{V} and \mathbf{W} defining the BC, are already stored in floating point format, so that e.g. $\text{fl}(\mathbf{A}) = \mathbf{A}$. Proposition 8 provides an upper bound on the numerical error for the (approximate or exact) BC f defined as in (1).

Proposition 8. *Suppose $(4n+R-1)\varepsilon_{\text{machine}} \leq 0.01$. For $\mathbf{A}, \mathbf{B} \in B_{\mu}^{(n)}$ and an (n, τ) -ABC f computed according to (5), we have*

$$\|\text{fl}(f(\mathbf{A}, \mathbf{B})) - f(\mathbf{A}, \mathbf{B})\| \leq 1.01(4n + R - 1)\sqrt{R}\varepsilon_{\text{machine}}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2},$$

where e.g. $\mathbf{U}_{::r} \in \mathbb{R}^{n \times n}$ is the r th frontal slice of \mathbf{U} , so that $\|\mathbf{U}_{::r}\|^2 = \sum_{k,l} u_{klr}^2$.

Proof. Throughout this proof, constants $\gamma, \theta, \lambda, \tilde{\gamma}, \tilde{\theta}, \tilde{\lambda}, \phi$ and ψ with subscripts are real numbers of magnitude less than or equal to $\varepsilon_{\text{machine}}$. Define

$$s_{kpr} \stackrel{\text{def}}{=} \text{fl}\left(\sum_{l=1}^p u_{klr} a_{kl}\right).$$

Then

$$\begin{aligned} s_{k1r} &= \text{fl}(u_{k1r} a_{k1}) = u_{k1r} a_{k1} (1 + \gamma_{k1r}), \\ s_{k2r} &= \text{fl}(s_{k1r} + u_{k2r} a_{k2}) = (s_{k1r} + u_{k2r} a_{k2} (1 + \gamma_{k2r})) (1 + \theta_{k2r}), \end{aligned}$$

etc. More generally,

$$s_{kpr} = \sum_{l=1}^p u_{klr} a_{kl} (1 + \gamma_{klr}) \prod_{\alpha=l}^p (1 + \theta_{k\alpha r}), \quad \theta_{k1r} \stackrel{\text{def}}{=} 0.$$

Next, define

$$\hat{s}_{pr} \stackrel{\text{def}}{=} \text{fl} \left(\sum_{k=1}^p \sum_{l=1}^n u_{klr} a_{kl} \right).$$

Then

$$\begin{aligned} \hat{s}_{1r} &= \text{fl} \left(\sum_{l=1}^n u_{1lr} a_{1l} \right) = s_{1nr}, \\ \hat{s}_{2r} &= \text{fl} \left(\hat{s}_{1r} + \text{fl} \left(\sum_{l=1}^n u_{2lr} a_{2l} \right) \right) = (s_{1nr} + s_{2nr})(1 + \lambda_{2r}), \end{aligned}$$

etc. More generally,

$$\hat{s}_{pr} = \sum_{k=1}^p s_{knr} \prod_{\alpha=k}^p (1 + \lambda_{\alpha r}), \quad \lambda_{1r} \stackrel{\text{def}}{=} 0.$$

Consequently,

$$P_r \stackrel{\text{def}}{=} \text{fl} \left(\sum_{k=1}^n \sum_{l=1}^n u_{klr} a_{kl} \right) = \hat{s}_{nr} = \sum_{k=1}^n \sum_{l=1}^n u_{klr} a_{kl} (1 + \gamma_{klr}) \left(\prod_{\alpha=l}^n (1 + \theta_{k\alpha r}) \right) \left(\prod_{\alpha=k}^n (1 + \lambda_{\alpha r}) \right).$$

Similarly,

$$Q_r \stackrel{\text{def}}{=} \text{fl} \left(\sum_{k'=1}^n \sum_{l'=1}^n v_{k'l'r} b_{k'l'} \right) = \sum_{k'=1}^n \sum_{l'=1}^n v_{k'l'r} b_{k'l'} (1 + \tilde{\gamma}_{k'l'r}) \left(\prod_{\alpha=l'}^n (1 + \tilde{\theta}_{k'\alpha r}) \right) \left(\prod_{\alpha=k'}^n (1 + \tilde{\lambda}_{\alpha r}) \right).$$

Now, define

$$z_p \stackrel{\text{def}}{=} \text{fl} \left(\sum_{r=1}^p w_{ijr} \left(\sum_{k,l=1}^n u_{klr} a_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'r} b_{k'l'} \right) \right).$$

We have

$$\begin{aligned}
z_1 &= \text{fl} \left(w_{ij1} \left(\sum_{k,l=1}^n u_{kl1} a_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'1} b_{k'l'} \right) \right) \\
&= w_{ij1} P_1 Q_1 (1 + \phi_{11}) (1 + \phi_{12}), \\
z_2 &= \text{fl} \left(z_1 + \text{fl} \left(w_{ij2} \left(\sum_{k,l=1}^n u_{kl2} a_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'2} b_{k'l'} \right) \right) \right) \\
&= (z_1 + w_{ij2} P_2 Q_2 (1 + \phi_{21}) (1 + \phi_{22})) (1 + \psi_2),
\end{aligned}$$

etc. From this, it follows that

$$\text{fl}(f(\mathbf{A}, \mathbf{B})_{ij}) = z_R = \sum_{r=1}^R w_{ijr} P_r Q_r (1 + \phi_{r1}) (1 + \phi_{r2}) \prod_{\zeta=r}^R (1 + \psi_\zeta), \quad \psi_1 \stackrel{\text{def}}{=} 0.$$

Writing this out in full and rearranging, we have

$$\begin{aligned}
\text{fl}(f(\mathbf{A}, \mathbf{B})_{ij}) &= \sum_{k,l=1}^n \sum_{k',l'=1}^n a_{kl} b_{k'l'} \sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} (1 + \gamma_{klr}) (1 + \tilde{\gamma}_{k'l'r}) \\
&\times \left(\prod_{\alpha=l}^n (1 + \theta_{k\alpha r}) \right) \left(\prod_{\alpha=l'}^n (1 + \tilde{\theta}_{k'\alpha r}) \right) \left(\prod_{\alpha=k}^n (1 + \lambda_{\alpha r}) \right) \\
&\times \left(\prod_{\alpha=k'}^n (1 + \tilde{\lambda}_{\alpha r}) \right) (1 + \phi_{r1}) (1 + \phi_{r2}) \prod_{\zeta=r}^R (1 + \psi_\zeta).
\end{aligned}$$

Since we have assumed $(4n + R - 1)\varepsilon_{\text{machine}} \leq 0.01$, it follows from Lemma 2.7.1 in [18] that there exist constants $\varepsilon_{klk'l'r}$, for $(k, l, k', l', r) \in [n]^4 \times [R]$, such that $|\varepsilon_{klk'l'r}| \leq 1.01(4n + R - 1)\varepsilon_{\text{machine}}$ and

$$\begin{aligned}
\text{fl}(f(\mathbf{A}, \mathbf{B})_{ij}) &= \sum_{k,l=1}^n \sum_{k',l'=1}^n a_{kl} b_{k'l'} \sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} (1 + \varepsilon_{klk'l'r}) \\
&= f(\mathbf{A}, \mathbf{B})_{ij} + \sum_{k,l=1}^n \sum_{k',l'=1}^n a_{kl} b_{k'l'} \sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} \varepsilon_{klk'l'r}.
\end{aligned}$$

We have

$$\begin{aligned}
|\text{fl}(f(\mathbf{A}, \mathbf{B})_{ij}) - f(\mathbf{A}, \mathbf{B})_{ij}|^2 &= \left| \sum_{k,l=1}^n \sum_{k',l'=1}^n a_{kl} b_{k'l'} \sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} \varepsilon_{klk'l'r} \right|^2 \\
&\leq \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n (a_{kl} b_{k'l'})^2 \right) \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n \left(\sum_{r=1}^R u_{klr} v_{k'l'r} w_{ijr} \varepsilon_{klk'l'r} \right)^2 \right) \\
&\leq \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n (a_{kl} b_{k'l'})^2 \right) \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n \left(\sum_{r=1}^R (u_{klr} v_{k'l'r} w_{ijr})^2 \right) \left(\sum_{r=1}^R \varepsilon_{klk'l'r}^2 \right) \right) \\
&\leq \|\mathbf{A}\|^2 \|\mathbf{B}\|^2 R (1.01(4n + R - 1) \varepsilon_{\text{machine}})^2 \sum_{r=1}^R \sum_{k,l=1}^n \sum_{k',l'=1}^n (u_{klr} v_{k'l'r} w_{ijr})^2,
\end{aligned}$$

where the first and second inequality follows from the Cauchy–Schwarz inequality, and the third inequality follows from the bound on $|\varepsilon_{klk'l'r}|$. Summing over all i, j , taking a square root, and using that $\|\mathbf{A}\|, \|\mathbf{B}\| \leq \mu$, we arrive at

$$\|\text{fl}(f(\mathbf{A}, \mathbf{B})) - f(\mathbf{A}, \mathbf{B})\| \leq 1.01(4n + R - 1) \sqrt{R} \varepsilon_{\text{machine}} \mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2}.$$

□

Proposition 9 generalizes Proposition 8 to the case when the input matrices are of size $mn \times mn$.

Proposition 9. *Suppose $(4n + m + R - 2) \varepsilon_{\text{machine}} \leq 1.01$. For $\mathbf{A}, \mathbf{B} \in B_{\mu}^{(mn)}$ and an (n, τ) -ABC f computed according to (5), we have*

$$\|\text{fl}(f(\mathbf{A}, \mathbf{B})) - f(\mathbf{A}, \mathbf{B})\| \leq 1.01(4n + m + R - 2) \sqrt{R} \varepsilon_{\text{machine}} \mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2}.$$

Proof. Note that all computations when calculating $\sum_{k,l} u_{klr} \mathbf{A}_{kl}$ and $\sum_{k',l'} v_{k'l'r} \mathbf{B}_{k'l'}$ are done elementwise. The error accumulated for each entry will therefore be the same as in the case when \mathbf{A}_{kl} and $\mathbf{B}_{k'l'}$ are scalars. More precisely,

$$\text{fl} \left(\sum_{k,l=1}^n u_{klr} \mathbf{A}_{kl} \right) = \sum_{k,l=1}^n u_{klr} \mathbf{A}_{kl} \circledast \mathbf{E}_{klr}, \tag{17}$$

where \circledast is the Hadamard (elementwise) product, and $\mathbf{E}_{klr} \in \mathbb{R}^{m \times m}$ is a matrix with entries which are the product of at most $2n - 1$ terms of the form $(1 + \varepsilon)$ with $|\varepsilon| \leq \varepsilon_{\text{machine}}$. A similar statement is true for $\text{fl}(\sum_{k',l'=1}^n v_{k'l'r} \mathbf{B}_{k'l'})$. When computing each

$$\text{fl} \left(w_{ijr} \left(\sum_{k,l=1}^n u_{klr} \mathbf{A}_{kl} \right) \left(\sum_{k',l'=1}^n v_{k'l'r} \mathbf{B}_{k'l'} \right) \right),$$

we now have additional error compared to the scalar case due to the inner product computations. In order to avoid cumbersome notation, we do not write out the following computations in full. Using the result in Section 2.7.6 of [18] for the rounding error in inner products, Equation (17), Lemma 2.7.1 in [18], and letting $(f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}$ denote the element on position (α, β) in matrix block (i, j) , we can compute

$$\text{fl}((f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}) = (f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta} + \sum_{r=1}^R w_{ijr} \sum_{k,l=1}^n \sum_{k',l'=1}^n \sum_{z=1}^m u_{klr} v_{k'l'r} (\mathbf{A}_{kl})_{\alpha z} (\mathbf{B}_{k'l'})_{z\beta} \Theta_{klk'l'r\alpha\beta z},$$

where each $|\Theta_{klk'l'r\alpha\beta z}| \leq \hat{\theta} \stackrel{\text{def}}{=} 1.01(4n + m + R - 2)\varepsilon_{\text{machine}}$. Rearranging this, we have

$$\begin{aligned} & |\text{fl}((f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}) - (f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}|^2 \\ &= \left| \sum_{r=1}^R \sum_{k,l=1}^n \sum_{k',l'=1}^n \sum_{z=1}^m u_{klr} v_{k'l'r} w_{ijr} (\mathbf{A}_{kl})_{\alpha z} (\mathbf{B}_{k'l'})_{z\beta} \Theta_{klk'l'r\alpha\beta z} \right|^2 \\ &\leq \hat{\theta}^2 \left| \sum_{r=1}^R \sum_{k,l=1}^n \sum_{k',l'=1}^n |u_{klr} v_{k'l'r} w_{ijr}| \sum_{z=1}^m |(\mathbf{A}_{kl})_{\alpha z} (\mathbf{B}_{k'l'})_{z\beta}| \right|^2 \end{aligned}$$

due to the triangle inequality and definition of $\hat{\theta}$,

$$\begin{aligned} &\leq R\hat{\theta}^2 \sum_{r=1}^R \left| \sum_{k,l=1}^n \sum_{k',l'=1}^n |u_{klr} v_{k'l'r} w_{ijr}| \sum_{z=1}^m |(\mathbf{A}_{kl})_{\alpha z} (\mathbf{B}_{k'l'})_{z\beta}| \right|^2 \\ &\leq R\hat{\theta}^2 \sum_{r=1}^R \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n |u_{klr} v_{k'l'r} w_{ijr}|^2 \right) \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n \left| \sum_{z=1}^m |(\mathbf{A}_{kl})_{\alpha z} (\mathbf{B}_{k'l'})_{z\beta}| \right|^2 \right) \\ &\leq R\hat{\theta}^2 \sum_{r=1}^R \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n |u_{klr} v_{k'l'r} w_{ijr}|^2 \right) \left(\sum_{k,l=1}^n \sum_{k',l'=1}^n \left(\sum_{z=1}^m |(\mathbf{A}_{kl})_{\alpha z}|^2 \right) \left(\sum_{z=1}^m |(\mathbf{B}_{k'l'})_{z\beta}|^2 \right) \right) \\ &= R(1.01(4n + m + R - 2)\varepsilon_{\text{machine}})^2 \\ &\quad \times \left(\sum_{k,l=1}^n \sum_{z=1}^m (\mathbf{A}_{kl})_{\alpha z}^2 \right) \left(\sum_{k',l'=1}^n \sum_{z=1}^m (\mathbf{B}_{k'l'})_{z\beta}^2 \right) \left(\sum_{r=1}^R \sum_{k,l=1}^n \sum_{k',l'=1}^n u_{klr}^2 v_{k'l'r}^2 w_{ijr}^2 \right) \end{aligned}$$

where each inequality follows from an application of the Cauchy–Schwarz inequality and the final equality follows from a rearrangement of terms. Finally, since

$$\|\text{fl}(f(\mathbf{A}, \mathbf{B})) - f(\mathbf{A}, \mathbf{B})\|^2 = \sum_{\alpha,\beta=1}^m \sum_{i,j=1}^n |\text{fl}((f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}) - (f(\mathbf{A}, \mathbf{B})_{ij})_{\alpha\beta}|^2$$

and $\|\mathbf{A}\|, \|\mathbf{B}\| \leq \mu$, it follows that

$$\|\text{fl}(f(\mathbf{A}, \mathbf{B})) - f(\mathbf{A}, \mathbf{B})\| \leq 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2}.$$

□

Proposition 10 shows how the error in Proposition 9 changes when a randomized algorithm \hat{f} is used instead of its deterministic counterpart. In particular, the bound is worse by only a factor $(1 - \kappa)^{-1}$.

Proposition 10. *Suppose $(4n + m + R - 2)\varepsilon_{\text{machine}} \leq 1.01$ and $1 - \kappa > 0$. For $\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}$ and an (n, τ, κ) -RandABC \hat{f} computed according to (8), we have*

$$\begin{aligned} & \|\text{fl}(\hat{f}(\mathbf{A}, \mathbf{B})) - \hat{f}(\mathbf{A}, \mathbf{B})\| \\ & \leq 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}(1 - \kappa)^{-1}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2}. \end{aligned}$$

Proof. From (10), it follows that

$$\begin{aligned} \hat{f}(\mathbf{A}, \mathbf{B})_{ij} &= (1 - \kappa)^{-1} s_1(i) s_3(j) \sum_{r=1}^R w_{\pi_1(i)\pi_3(j)r} \\ & \times \left(\sum_{k,l=1}^n u_{\pi_1(k)\pi_2(l)r} s_1(k) s_2(l) \mathbf{A}_{kl} \right) \left(\sum_{k',l'=1}^n v_{\pi_2(k')\pi_3(l')r} s_2(k') s_3(l') \mathbf{B}_{k'l'} \right) \\ & = \sum_{r=1}^R \hat{w}_{ijr} \left(\sum_{k,l=1}^n \hat{u}_{klr} \mathbf{A}_{kl} \right) \left(\sum_{k',l'=1}^n \hat{v}_{k'l'r} \mathbf{B}_{k'l'} \right), \end{aligned} \quad (18)$$

where

$$\begin{aligned} \hat{u}_{klr} &\stackrel{\text{def}}{=} u_{\pi_1(k)\pi_2(l)r} s_1(k) s_2(l), \\ \hat{v}_{k'l'r} &\stackrel{\text{def}}{=} v_{\pi_2(k')\pi_3(l')r} s_2(k') s_3(l'), \\ \hat{w}_{ijr} &\stackrel{\text{def}}{=} (1 - \kappa)^{-1} s_1(i) s_3(j) w_{\pi_1(i)\pi_3(j)r}. \end{aligned}$$

Note that $\|\hat{\mathbf{U}}_{::r}\| = \|\mathbf{U}_{::r}\|$, $\|\hat{\mathbf{V}}_{::r}\| = \|\mathbf{V}_{::r}\|$ and $\|\hat{\mathbf{W}}_{::r}\| = (1 - \kappa)^{-1} \|\mathbf{W}_{::r}\|$. Using this fact, and applying Proposition 9 to the computation in (18), gives us the desired result. □

Propositions 11 and 12 give upper bounds for the total error, both due to algorithmic error and numerical rounding, for the deterministic and randomized algorithms, respectively.

Proposition 11. *Suppose $(4n + m + R - 2)\varepsilon_{\text{machine}} \leq 1.01$. For $\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}$ and an (n, τ) -ABC f computed according to (5), we have*

$$\begin{aligned} & \|\text{fl}(f(\mathbf{A}, \mathbf{B})) - \mathbf{AB}\| \\ & \leq 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2} + \mu^2 \tau. \end{aligned} \quad (19)$$

Proof. Follows by applying the triangle inequality and using Propositions 9 and 4 (i). □

Proposition 12. Suppose $(4n + m + R - 2)\varepsilon_{\text{machine}} \leq 1.01$. For $\mathbf{A}, \mathbf{B} \in B_\mu^{(mn)}$ and an (n, τ, κ) -RandABC \hat{f} computed according to (8), we have

$$\begin{aligned} & \|\mathbb{E}[\text{fl}(\hat{f}(\mathbf{A}, \mathbf{B}))] - \mathbf{AB}\| \\ & \leq 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}(1 - \kappa)^{-1}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2}. \end{aligned} \quad (20)$$

Proof. Follows by taking expectations of the inequality in Proposition 10 and applying Jensen's inequality [28]. \square

Note that Proposition 12 implies that

$$\mathbb{E}[\text{fl}(\hat{f}(\mathbf{A}, \mathbf{B}))] \rightarrow \mathbf{AB} \quad \text{as} \quad \varepsilon_{\text{machine}} \rightarrow 0,$$

which means that as the numerical precision increases, $\mathbb{E}[\text{fl}(\hat{f}(\mathbf{A}, \mathbf{B}))]$ approaches \mathbf{AB} . This is consistent with Proposition 3.

To make the bounds in (19) and (20) easier to compare, using the fact that $(1 - \kappa)^{-1} = 1 + \kappa + O(\kappa^2)$, we can rewrite (20) as

$$\begin{aligned} & \|\mathbb{E}[\text{fl}(\hat{f}(\mathbf{A}, \mathbf{B}))] - \mathbf{AB}\| \\ & \leq 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2} \\ & \quad + 1.01(4n + m + R - 2)\sqrt{R}\varepsilon_{\text{machine}}\kappa\mu^2 \sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2} + O(\varepsilon_{\text{machine}}\kappa^2). \end{aligned} \quad (21)$$

The first term on the right hand side of (19) and (21) are identical. Recall from (14) that $|\kappa| \leq n^{-5/2}\|\mathbf{y} - \mathbf{x}\| \leq n^{-5/2}\tau$. Consequently, if the quantity

$$\sqrt{\sum_{r=1}^R \|\mathbf{U}_{::r}\|^2 \|\mathbf{V}_{::r}\|^2 \|\mathbf{W}_{::r}\|^2} \quad (22)$$

is not too large, and m, n and R are of moderate size, the second term in (21) will be smaller than the second term in (19), showing that the result in Proposition 3 largely carries over to a setting with floating point arithmetic. The following example illustrates this.

Example 13. For the sake of this example, suppose we are using an approximate variant of Strassen's algorithm. Then $R = 7$ and $n = 2$. For Strassen's algorithm, the quantity in (22) is approximately 6, so we will assume that it is less than 10 for our approximate variant of the algorithm. Suppose we are multiplying two $100,000 \times 100,000$ matrices in single precision, so that $m = 50,000$ and $\varepsilon_{\text{machine}} \sim 10^{-8}$. Then, the second term in (21) is upper bounded by a quantity which is on the order of

$$10 \cdot 1.01 \cdot (4 \cdot 2 + 50000 + 7 - 2) \cdot \frac{\sqrt{7}}{2^{5/2}} \cdot 10^{-8} \mu^2 \tau \approx 0.0024 \mu^2 \tau,$$

which is much smaller than the second term in (19). The implementation of Strassen’s algorithm in [21] achieves a speed-up over an efficient implementation of the standard matrix multiplication algorithm for square matrices with as few as 1,536 rows/columns. So multiplication of matrices with 100,000 rows/columns is well beyond the problem size for which fast algorithms can outperform the standard algorithm.

3. Experiments

In this section we present some results from experiments, with additional results provided in Appendix A. We implement all experiments in Matlab with certain parts implemented in C. All our code is available online at <https://github.com/OsmanMalik/random-approximate-matrix-multiplication>.

In our experiments, we draw the matrices \mathbf{A}, \mathbf{B} from different random distributions. By *Gaussian matrix*, we mean a matrix whose elements are realizations of i.i.d. standard normal random variables. Similarly, a *uniform matrix* is one whose elements are realizations of i.i.d. Uniform(0,1) random variables. We also consider three types of random *adversarial matrices*, which were proposed by [1] and are designed to be challenging for Strassen’s algorithm.

Definition 14 (Adversarial matrices). Consider a matrix pair $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. We say that it is *type 1 adversarial* if

$$a_{ij} \sim \begin{cases} \text{Uniform}(0, 1/n^2) & \text{if } j > n/2, \\ \text{Uniform}(0, 1) & \text{otherwise,} \end{cases} \quad b_{ij} \sim \begin{cases} \text{Uniform}(0, 1/n^2) & \text{if } i < n/2, \\ \text{Uniform}(0, 1) & \text{otherwise.} \end{cases}$$

We say that the matrix pair is *type 2 adversarial* if

$$a_{ij} \sim \begin{cases} \text{Uniform}(0, n^2) & \text{if } i < n/2 \text{ and } j > n/2, \\ \text{Uniform}(0, 1) & \text{otherwise,} \end{cases} \quad b_{ij} \sim \begin{cases} \text{Uniform}(0, 1/n^2) & \text{if } j < n/2, \\ \text{Uniform}(0, 1) & \text{otherwise.} \end{cases}$$

We say that the matrix pair is *type 3 adversarial* if

$$a_{ij}, b_{ij} \sim \begin{cases} \text{Uniform}(0, 1/n^2) & \text{if } i < n/2 \text{ and } j > n/2, \text{ or if } i \geq n/2 \text{ and } j \leq n/2 \\ \text{Uniform}(0, 1) & \text{otherwise,} \end{cases}$$

Here, all the entries are assumed to be independent.

We will also consider the Hilbert matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$, which has entries $h_{ij} = 1/(i + j - 1)$, since it appears to be a particularly challenging matrix to the rescaling method which we consider in Section 3.2.

For an ABC, the corresponding randomized computation with Q recursions, $\hat{F}^{(Q)}$, was defined in (16). In this section, we will use $F^{(Q)}$ to denote the deterministic counterpart. It is defined in the same way, but with each $s_i^{(q)}(j) = 1$ and each $\pi_i^{(q)}(j) = j$, i.e., with no randomness involved. As in Example 7, for EBCs we will use $\hat{G}^{(Q)}$ and $G^{(Q)}$ to denote the corresponding randomized and deterministic computations with Q recursions.

3.1. Approximate algorithm

We create an ABC of the form (5) by taking the tensors \mathbf{U} , \mathbf{V} and \mathbf{W} corresponding to Strassen’s algorithm and perturbing them: For each of the three tensors, we add i.i.d. mean zero Gaussian noise with standard deviation 10^{-3} to each element in the tensor equal to 1 as well as to five randomly selected elements that are equal to 0. Since we do the computations in double precision, and since we add a considerable amount of noise, these experiments are designed to test Propositions 3–6, i.e., we can ignore any floating point error.

In the first experiment, we draw two Gaussian matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ and compute

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{F}_i^{(Q)}(\mathbf{A}, \mathbf{B}) - \mathbf{AB} \right\| / \|\mathbf{AB}\| \quad (23)$$

for $n \in [10^4]$ and $Q \in [3]$. Here, $\hat{F}_i^{(Q)}$ is the i th realization of $\hat{F}^{(Q)}$. Figure 1 shows the results. As expected from Propositions 3 and 6, the quantity in (23) becomes smaller as n increases.

In the second experiment, we draw two Gaussian matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ and compute

$$\|\hat{F}_i^{(Q)}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| / \|\mathbf{AB}\| \quad (24)$$

for $i \in [100]$ and $Q \in [5]$, i.e., the relative error for each of 100 trials. The box plots in Figure 2 compares the empirical distribution of (24) to the relative error for the deterministic approximate algorithm, i.e., $\|F^{(Q)}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| / \|\mathbf{AB}\|$. In this particular case, randomization does not impact the median error and there is very little variation between trials. Figure 3 repeats this experiment, but with $\mathbf{A} = \mathbf{B} = \mathbf{H}$, where \mathbf{H} is the 320×320 Hilbert matrix. In this case, the randomized scheme frequently results in a slightly larger error. In Figure A1–A5 in Appendix A we provide additional results for when \mathbf{A}, \mathbf{B} are Gaussian, uniform and type 1–3 adversarial. Those results demonstrate that randomization can both increase and decrease the error in this setting. However, as expected from Propositions 4 and 5, the difference in error between the randomized and deterministic variants is typically not substantial.

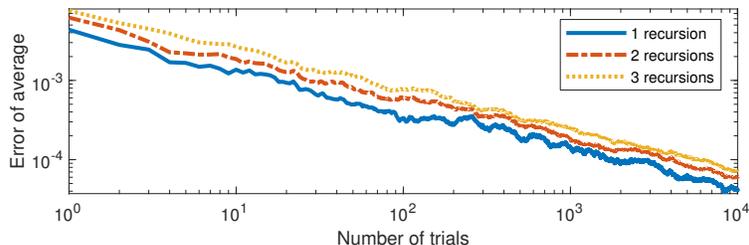


Figure 1.: Error of average for randomized ABC. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are Gaussian and remain fixed throughout the experiment. The ABC is a perturbed variant of Strassen’s algorithm.

To see how these results carry over to other ABCs, we also present results in Figure 4 from experiments which use an instance of the 12×12 APA algorithm in [4] with $\varepsilon = 1e-4$ in the representation (6). Results are shown for single recursion experiments on 12×12 matrices drawn from different distributions. As for the ABC based on the perturbed variant of Strassen’s algorithm, the randomized variant of this ABC sometimes results in a smaller and sometimes in a larger error than the deterministic counterpart.

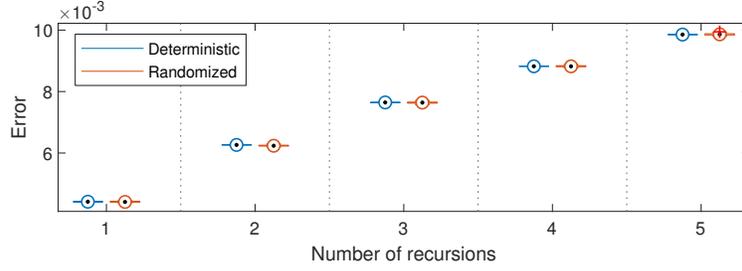


Figure 2.: Error for deterministic ABC compared to the error of the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are Gaussian and remain fixed over all realizations. Note that, unlike Figure 1, this figure shows the distribution of errors, rather than the errors of averages. The ABC is a perturbed variant of Strassen’s algorithm.

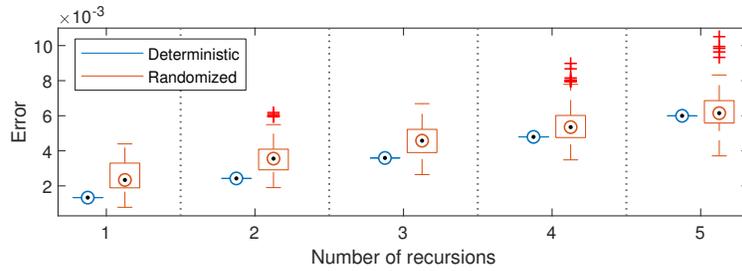


Figure 3.: Same as Figure 2, but with \mathbf{A} and \mathbf{B} equal to the 320×320 Hilbert matrix.

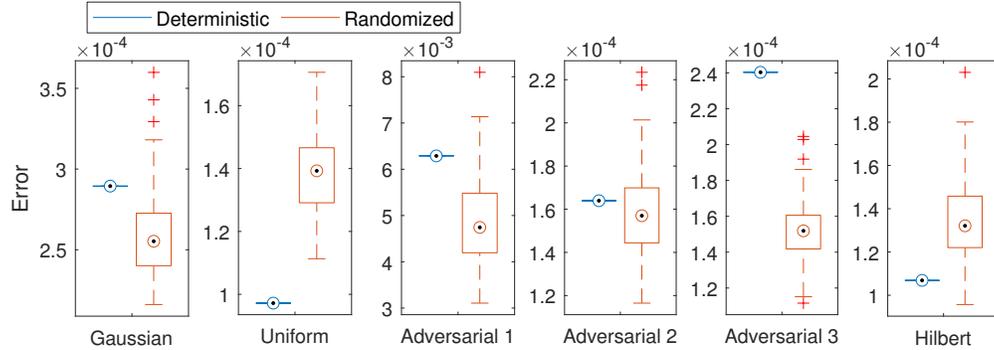


Figure 4.: Error for deterministic ABC compared to the error of the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{12 \times 12}$ are drawn from six different random distributions, one in each subplot, and remain fixed over all realizations. The ABC is an instance of the 12×12 APA algorithm of [4].

3.2. Exact algorithm in single precision floating point arithmetic

We first consider Strassen’s algorithm, without any perturbations so that it is exact, when the computations are done in single precision floating point arithmetic. In error computations, we use the double precision product for \mathbf{AB} computed using the standard algorithm as the true value of the product. Recall that by “standard algorithm” we mean the $O(n^3)$ algorithm.

In the first experiment, we draw two Gaussian matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ and compute the quantity in (23), but with each $\hat{F}_i^{(Q)}$ replaced by $\hat{G}_i^{(Q)}$, where $\hat{G}_i^{(Q)}$ is the i th realization of $\hat{G}^{(Q)}$, for $n \in [10^4]$ and $Q \in [3]$. Figure 5 shows the results, where we also have included the error for the standard algorithm computed in single precision as a reference. Although it is clear that the randomized algorithms do not converge to the exact correct answer, it seems like their expectations perform better than the standard algorithm. Although the figure only shows the result for a specific random pair (\mathbf{A}, \mathbf{B}) , we get qualitatively similar results every time we draw a new Gaussian matrix pair. The fact that the average of a single, or a few, outcomes of the randomized algorithms performs worse than the standard algorithm is to be expected, since Strassen’s algorithm is more susceptible to numerical error than the standard algorithm. Figure 6 repeats this experiment, but with $\mathbf{A} = \mathbf{B} = \mathbf{H}$, where \mathbf{H} is the 80×80 Hilbert matrix. The results are similar to those in Figure 5. Figures A6–A10 in the appendix provide additional results for when \mathbf{A}, \mathbf{B} are Gaussian, uniform and type 1–3 adversarial.

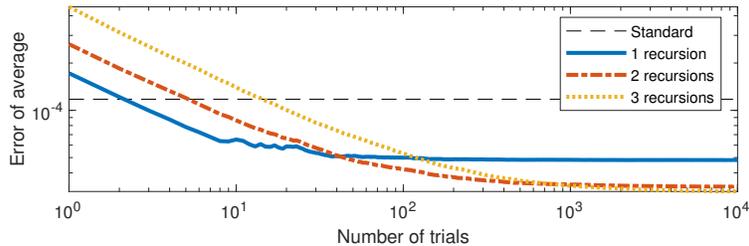


Figure 5.: Error of average for randomized EBC compared to the standard $O(n^3)$ algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are Gaussian and remain fixed throughout the experiment.

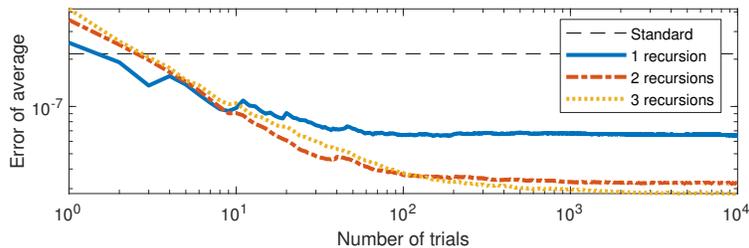


Figure 6.: Same as Figure 5, but with \mathbf{A} and \mathbf{B} equal to the 80×80 Hilbert matrix.

In the second experiment, we first compare the deterministic and three randomized versions of Strassen’s algorithm, as well as the rescaling scheme proposed in [1]. The first randomized version uses both random permutations and random signs. The two other randomized versions use only random signs and only random permutations, respectively. We include these two variants to better understand how random signs and random permutations each impact the performance. These three randomized methods will be referred to as “fully randomized,” “random sign” and “random permutation,” respectively. The algorithm that only uses random permutations corresponds to the method suggested in [11].

The purpose of the rescaling scheme in [1] is to improve numerical stability of fast EBCs, and has

two steps: Outside scaling and inside scaling. With outside scaling, $\mathbf{C} = \mathbf{A}\mathbf{B}$ is computed via

$$\mathbf{C}_{\text{outside}} \stackrel{\text{def}}{=} \mathbf{D}_{\mathbf{A}} G^{(Q)}(\mathbf{D}_{\mathbf{A}}^{-1}\mathbf{A}, \mathbf{B}\mathbf{D}_{\mathbf{B}}^{-1})\mathbf{D}_{\mathbf{B}}, \quad (25)$$

where $\mathbf{D}_{\mathbf{A}} \stackrel{\text{def}}{=} \text{diag}(\max_j |a_{ij}|)$ and $\mathbf{D}_{\mathbf{B}} \stackrel{\text{def}}{=} \text{diag}(\max_i |b_{ij}|)$. With inside scaling, \mathbf{C} is instead computed via

$$\mathbf{C}_{\text{inside}} \stackrel{\text{def}}{=} G^{(Q)}(\mathbf{A}\mathbf{D}, \mathbf{D}^{-1}\mathbf{B}), \quad (26)$$

where $\mathbf{D} \stackrel{\text{def}}{=} \text{diag}(\sqrt{\max_j |b_{kj}| / \max_i |a_{ik}|})$. In exact arithmetic, the scaling matrices in (25) and (26) will cancel out, in which case $\mathbf{C}_{\text{outside}} = \mathbf{C}_{\text{inside}} = \mathbf{C}$. Both outside and inside scaling can be applied at the same time, as well as multiple times in an alternating fashion. The rescaling scheme that works best in numerical experiments in [1] does outside-inside rescaling twice. We use the same rescaling scheme in our experiments, which we refer to as ‘‘Rescaled 2x O-I.’’ See Section 6 in [1], in particular Algorithm 3, for further details on the rescaling method.

For the experiment, we draw two random matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ and compute the quantity in (24), but with each $\hat{F}_i^{(Q)}$ replaced by $\hat{G}_i^{(Q)}$, for $i \in [100]$ and $Q \in [5]$. Figures 7–11 show the results for Gaussian, uniform and type 1–3 adversarial matrices. Figure 12 shows the result when \mathbf{A} and \mathbf{B} are both Hilbert matrices. These figures include the error of the standard algorithm computed in single precision as a reference.

When the matrices are Gaussian (Figure 7), all algorithms perform roughly the same with little variation between trials. For uniform matrices (Figure 8), the rescaling method performs about the same as the deterministic method. The fully randomized method has a lower error than the deterministic method, and it seems like this improvement comes from the random signs. For type 1 adversarial matrices (Figure 9), the rescaling method does remarkably well, achieving the same accuracy as the standard algorithm. The fully randomized algorithm outperforms the deterministic algorithm, and it seems like this improvement is coming from the random signs. For type 2 adversarial matrices (Figure 10), the fully randomized method will sometimes perform worse than the deterministic algorithm, but has a lower median error for 2 or more recursions. The rescaling method also improves on the deterministic method, although the median error for the fully randomized method is lower for 4 recursions or more. Type 3 adversarial matrices were specifically proposed in [1] to show a situation when rescaling does not work. This is clear in Figure 11, where the rescaling method has the same error as the deterministic method. Our fully randomized method, however, has a lower error, and it seems like both the random signs and random permutations contribute to this performance improvement. When the matrices are Hilbert matrices (Figure 12), the rescaling method does very poorly, with a much larger error than the deterministic method. Once again, our randomized method reduces the error compared to the deterministic method, with both the random signs and random permutations contributing to the improved performance. Figures A11–A15 in the appendix provide additional results for when \mathbf{A}, \mathbf{B} are Gaussian, uniform, and type 1–3 adversarial.

In Figure 13, we give the results for some experiments which use an EBC for 12×12 matrix multiplication derived as in [3] from the 12×12 APA algorithm in [4] via the approach discussed in Section 1.1.² The APA algorithm in question has an error tensor whose entries are polynomials with maximum degree $d = 6$. Consequently, 7 distinct values of ε are required in (7) to derive an exact scheme. For this purpose, we choose $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_7$ to be $0.1, 0.2, \dots, 0.7$. Our results are for

²There appears to be a few typos in the definition of $w_r^{(s)}$ in Equation (5.2) in [3], which defines the APA scheme. We encourage the reader to consult our code for a corrected definition.

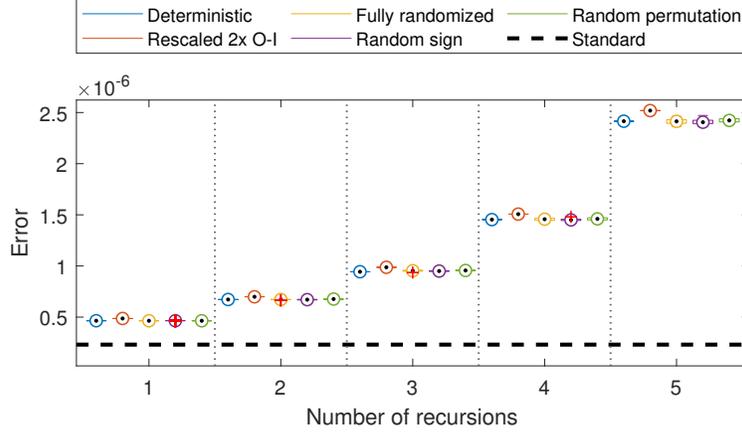


Figure 7.: Error for different variants of the Strassen EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are Gaussian and remain fixed over all realizations.

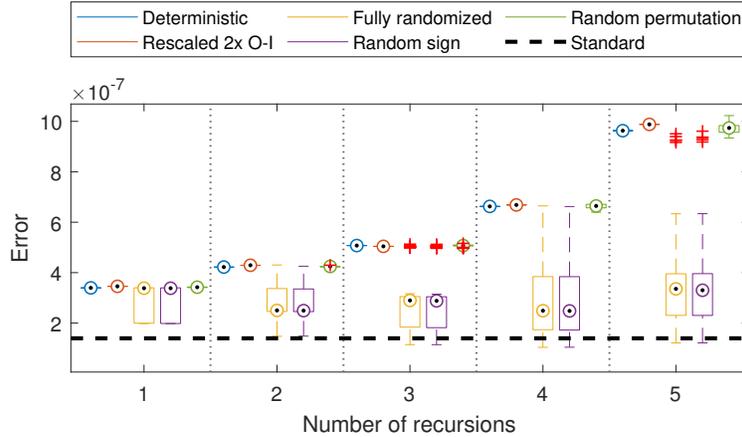


Figure 8.: Same as Figure 7, but with uniform matrices.

single recursion experiments on 12×12 matrices drawn from different distributions. The results are similar to those for the Strassen EBC, except for the Gaussian case where the variability of the randomized algorithms is much higher. The variability of the randomized methods also appears to be somewhat higher on the other matrix types as well. Overall, our randomized methods perform favorably compared to the deterministic method, especially for the adversarial matrices and the Hilbert matrix. The rescaling method once again performs very well on type 1 adversarial matrices, but poorly on the Hilbert matrix.

Although our method has some variability in performance due to being randomized, it seems to reduce the error compared to the deterministic method more reliably on a wider range of matrices than the rescaling method does, especially when more recursions are used. Our method also works well with the Hilbert matrix, which leads to large errors for the rescaling method. One benefit of our randomized method is that it can be done exactly even in low precision arithmetic, since it

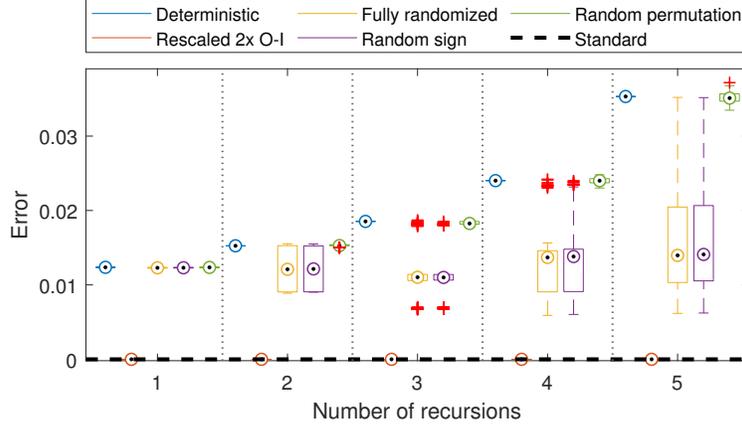


Figure 9.: Same as Figure 7, but with type 1 adversarial matrices.

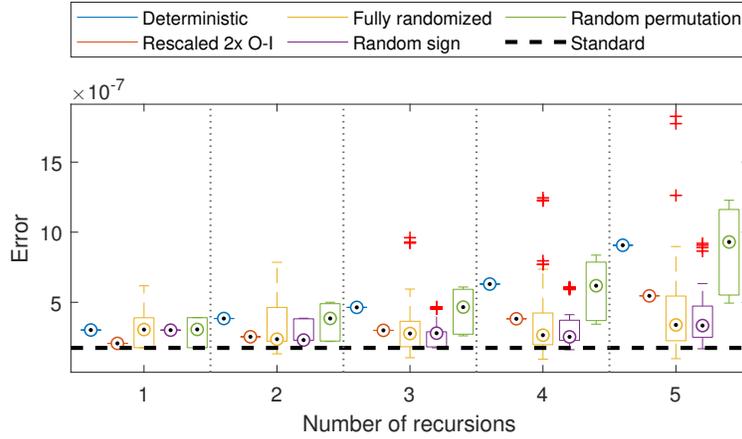


Figure 10.: Same as Figure 7, but with type 2 adversarial matrices.

only involves permutation of rows/columns and the flipping of signs. The rescaling method, on the other hand, involves diagonal matrices with floating point numbers along the diagonals, which adds another potential source for numerical error in the algorithm. These experiments also indicate that both random signs and random permutations may separately help to reduce the error, and that combining the two seems to lower the error further.

3.3. Randomization of ABC derived from APA algorithm

In Section 1.1 we discussed APA algorithms and how to derive EBCs from them by taking a linear combination of a few instances of the APA algorithm following an idea of Bini [3]. Although the EBCs derived in this fashion are exact mathematically, they may suffer from greater numerical error than the standard $O(n^3)$ algorithm due to cancellation. In this subsection, we compare an instance of the EBC in [3], which is derived from the APA algorithm in [4], to a randomized version of the ABC we get by fixing the error parameter in the same APA algorithm. The goal with these experiments

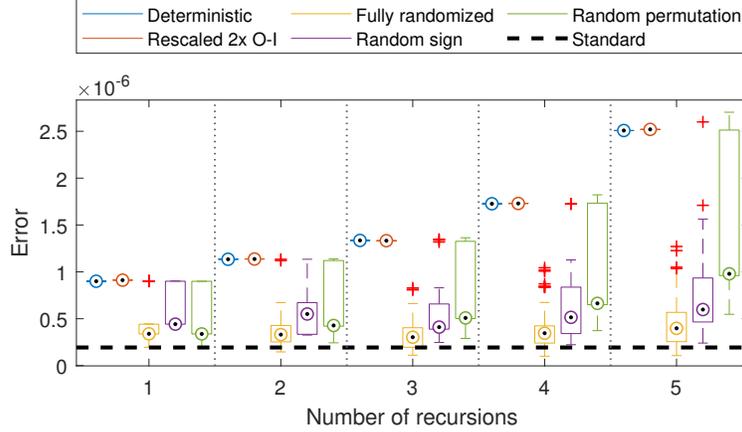


Figure 11.: Same as Figure 7, but with type 3 adversarial matrices.

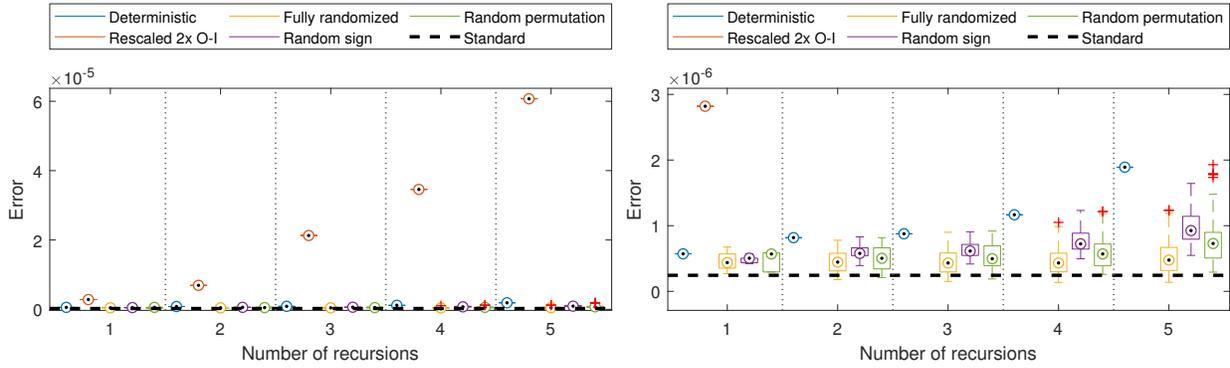


Figure 12.: Same as Figure 7, but with Hilbert matrices. The left and right plots show the same results, with the right plot zoomed in closer to the interesting portion. The left plot is included to give a sense of the size of the errors for the rescaling method compared to the other methods.

is to see if the expectation of our randomized ABC can perform better than the EBC. For the EBC, which is given by (7), we choose $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_7$ to be $0.1, 0.2, \dots, 0.7$. To get an ABC, we choose $\varepsilon = 1.5e-2$ in the APA algorithm. This ABC is then randomized as in Definition 2. The experiments are similar to that in Figure 1, but with an added baseline which shows the performance of the EBC. The experiments are done for a single recursion on 12×12 matrices. All computations are done in single precision arithmetic. Figure 14 shows the results, which each subplot showing the outcome for a different kind of matrix. Based on these result, the expectation of the randomized ABC seems to perform better than the EBC.

4. Conclusion

In this paper we have suggested an approach for randomizing formulas for bilinear computation of matrix products which does not increase the asymptotic computational complexity. We have considered the implications of this approach when there are two sources of error: The first due to

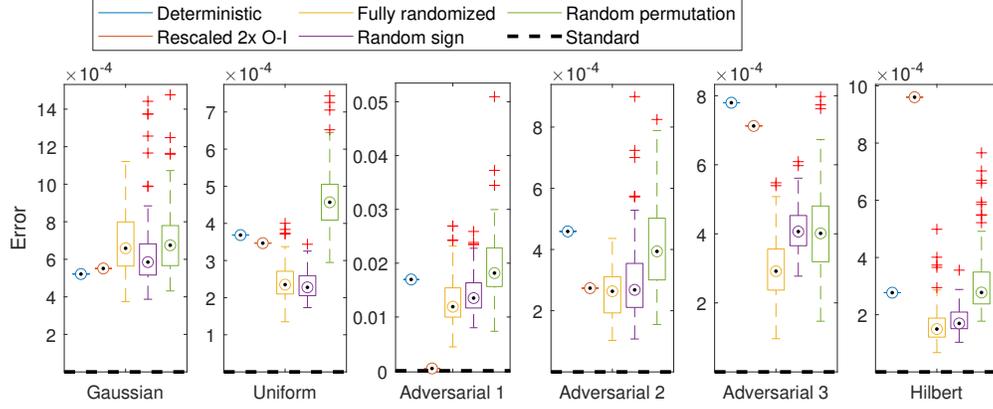


Figure 13.: Error for different variants of the Bini [3] EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{12 \times 12}$ are drawn from six different random distributions, one in each subplot, and remain fixed over all realizations.

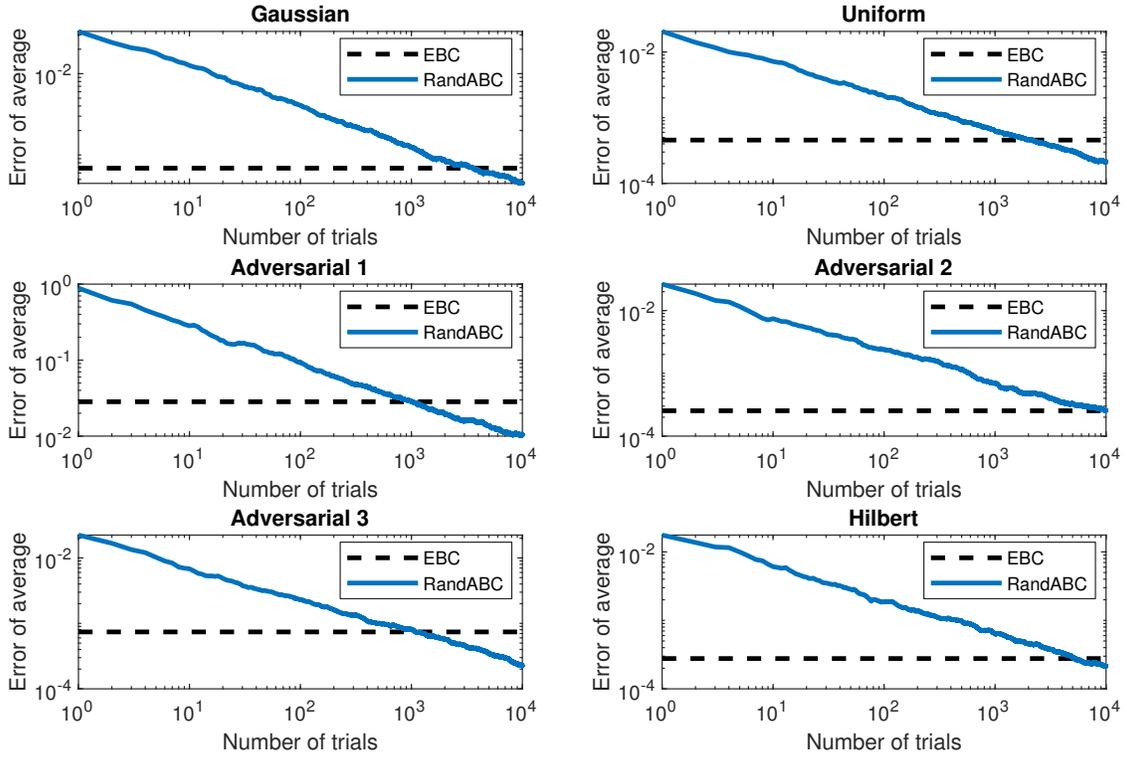


Figure 14.: Comparison of EBC in [3] to randomized ABC derived from APA algorithm in [4].

the algorithm itself being only approximately correct, and the second due to numerical error from using floating point arithmetic. We believe that our results are encouraging, and provide ideas for improving the properties of matrix multiplication when these error sources are present separately or together.

An interesting area for future research is to investigate other methods of randomization (e.g.

combining the random sign changes in this paper with the fast Hadamard transform) and see if such a method can further improve the results. It would also be interesting to investigate how fast randomized approximate methods for matrix multiplication can be used in applications. One potential application areas is for computations in neural networks, where some amount of error in the computation is acceptable, and where randomization may help improve robustness of the trained model. It would also be interesting to investigate how our randomization scheme can be used as a component in fast linear algebra algorithms, such as recursive dense matrix inversion; see [13] for details.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Grant No. ECCS-1810314.

References

- [1] G. Ballard, A. Benson, A. Druinsky, B. Lipshitz, and O. Schwartz, *Improving the Numerical Stability of Fast Matrix Multiplication*, SIAM Journal on Matrix Analysis and Applications 37 (2016), pp. 1382–1418.
- [2] A.R. Benson and G. Ballard, *A Framework for Practical Parallel Fast Matrix Multiplication*, in *ACM SIGPLAN Notices*, Vol. 50. ACM, 2015, pp. 42–53.
- [3] D. Bini, *Relations between exact and approximate bilinear algorithms. Applications*, Calcolo 17 (1980), pp. 87–97.
- [4] D. Bini, M. Capovani, F. Romani, and G. Lotti, *$O(n^{2.7799})$ complexity for $n \times n$ approximate matrix multiplication*, Information Processing Letters 8 (1979), pp. 234–235.
- [5] D. Bini and G. Lotti, *Stability of Fast Algorithms for Matrix Multiplication*, Numerische Mathematik 36 (1980), pp. 63–72.
- [6] D. Bini, G. Lotti, and F. Romani, *Approximate solutions for the bilinear form computational problem*, SIAM Journal on Computing 9 (1980), pp. 692–697.
- [7] M. Bläser, *On the complexity of the multiplication of matrices of small formats*, Journal of Complexity 19 (2003), pp. 43–60.
- [8] R.P. Brent, *Algorithms for matrix multiplication*, Tech. Rep. STAN-CS-70-157, Stanford University, 1970.
- [9] R.P. Brent, *Error Analysis of Algorithms for Matrix Multiplication and Triangular Decomposition using Winograd’s Identity*, Numerische Mathematik 16 (1970), pp. 145–156.
- [10] P. Bürgisser, M. Clausen, and M.A. Shokrollahi, *Algebraic Complexity Theory*, Vol. 315, Springer Science & Business Media, 2013.
- [11] R.R. Castrapel and J.L. Gustafson, *Precision improvement method for the Strassen/Winograd matrix multiplication method*, U.S. Patent No. 7209939B2 (2007), pp. 1–11.
- [12] H. De Silva, J.L. Gustafson, and W.F. Wong, *Making Strassen Matrix Multiplication Safe*, in *2018 IEEE 25th International Conference on High Performance Computing (HiPC)*. IEEE, 2018, pp. 173–182.
- [13] J. Demmel, I. Dumitriu, and O. Holtz, *Fast linear algebra is stable*, Numerische Mathematik 108 (2007), pp. 59–91.
- [14] J. Demmel, I. Dumitriu, O. Holtz, and R. Kleinberg, *Fast matrix multiplication is stable*, Numerische Mathematik 106 (2007), pp. 199–224.

- [15] P. Drineas, R. Kannan, and M.W. Mahoney, *Fast Monte Carlo algorithms for matrices I: Approximating matrix multiplication*, SIAM Journal on Computing 36 (2006), pp. 132–157.
- [16] B. Dumitrescu, *Improving and estimating the accuracy of Strassen’s algorithm*, Numerische Mathematik 79 (1998), pp. 485–499.
- [17] V. Elser, *A Network That Learns Strassen Multiplication*, The Journal of Machine Learning Research 17 (2016), pp. 3964–3976.
- [18] G.H. Golub and C.F. Van Loan, *Matrix Computations*, 4th ed., Johns Hopkins University Press, Baltimore, 2013.
- [19] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, *Deep Learning with Limited Numerical Precision*, in *International Conference on Machine Learning*. 2015, pp. 1737–1746.
- [20] M. Hopkins, M. Mikaitis, D.R. Lester, and S. Furber, *Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ODEs*, arXiv preprint arXiv:1904.11263 (2019).
- [21] J. Huang, C.D. Yu, and R.A. van de Geijn, *Implementing Strassen’s Algorithm with CUTLASS on NVIDIA Volta GPUs*, arXiv:1808.07984 [cs] (2018).
- [22] S. Huss-Lederman, E.M. Jacobson, J.R. Johnson, A. Tsao, and T. Turnbull, *Implementation of Strassen’s Algorithm for Matrix Multiplication*, in *Supercomputing’96: Proceedings of the 1996 ACM/IEEE Conference on Supercomputing*. IEEE, 1996, pp. 32–32.
- [23] R.W. Johnson and A.M. McLoughlin, *Noncommutative Bilinear Algorithms for 3×3 Matrix Multiplication*, SIAM Journal on Computing 15 (1986), pp. 595–603.
- [24] I. Kaporin, *The aggregation and cancellation techniques as a practical tool for faster matrix multiplication*, Theoretical Computer Science 315 (2004), pp. 469–510.
- [25] J.D. Laderman, *A noncommutative algorithm for multiplying 3×3 matrices using 23 multiplications*, Bulletin of the American Mathematical Society 82 (1976), pp. 126–128.
- [26] R. Pagh, *Compressed Matrix Multiplication*, ACM Transactions on Computation Theory 5 (2013), pp. 9:1–9:17.
- [27] S.I. Resnick, *A Probability Path*, Modern Birkhäuser Classics, Birkhäuser Basel, 2014.
- [28] M. Schaefer, *Note on the k -dimensional Jensen inequality*, The Annals of Probability (1976), pp. 502–504.
- [29] A.V. Smirnov, *The bilinear complexity and practical algorithms for matrix multiplication*, Computational Mathematics and Mathematical Physics 53 (2013), pp. 1781–1795.
- [30] V. Strassen, *Gaussian Elimination is not Optimal*, Numerische Mathematik 13 (1969), pp. 354–356.
- [31] N. Wang, J. Choi, D. Brand, C.Y. Chen, and K. Gopalakrishnan, *Training Deep Neural Networks with 8-Bit Floating Point Numbers*, in *Advances in Neural Information Processing Systems*. 2018, pp. 7675–7684.

Appendix A. Additional experiments

First, we repeat the second experiment we did in Section 3.1 for the Strassen ABC with multiple Gaussian, uniform, and type 1–3 adversarial matrices. We use the same setup as in the main manuscript: For each experiment, we create an approximate algorithm by perturbing Strassen’s algorithm, we draw random matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$, and then we compute the quantity in (24) for $i \in [100]$ and $Q \in [5]$ and compare it to the relative error for the deterministic approximate algorithm. Figures A1–A5 show the results. In the case of Gaussian matrices (Figure A1), the results look very similar to those in the main manuscript, with almost no difference in error between the deterministic and the randomized approximate algorithms. Figures A2–A5 (uniform and type 1–3

adversarial) show that randomization can both increase and decrease the error. Note that both the deterministic and randomized approximate algorithms do particularly poorly on type 1 adversarial matrices.

Next, we repeat the first experiment in Section 3.2 with multiple Gaussian, uniform and type 1–3 adversarial matrices. We use the same setup as in the main manuscript: For each experiment, we draw random matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ and compute the quantity in (23), but with each $\hat{F}_i^{(Q)}$ replaced by $\hat{G}_i^{(Q)}$, for $n \in [10^4]$ and $Q \in [3]$. Figures A6–A10 show the results. Although using more recursions seems to increase the error for a single realization of the algorithm, it also seems to reduce the error of the expectation of the computed matrix product.

Finally, we repeat the second experiment in Section 3.2 for the Strassen EBC with multiple Gaussian, uniform and type 1–3 adversarial matrices. We use the same setup as in the main manuscript: For each experiment, we draw random matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ and compute the quantity in (24), but with each $\hat{F}_i^{(Q)}$ replaced by $\hat{G}_i^{(Q)}$, for $i \in [100]$ and $Q \in [5]$. We do the same for variants of the algorithm that only use random permutations or random sign functions. We compare these to the deterministic version $G^{(Q)}$, and also include the error of the standard algorithm computed in single precision as a reference. Figures A11–A15 show the results, which look very similar to those presented in the main manuscript.

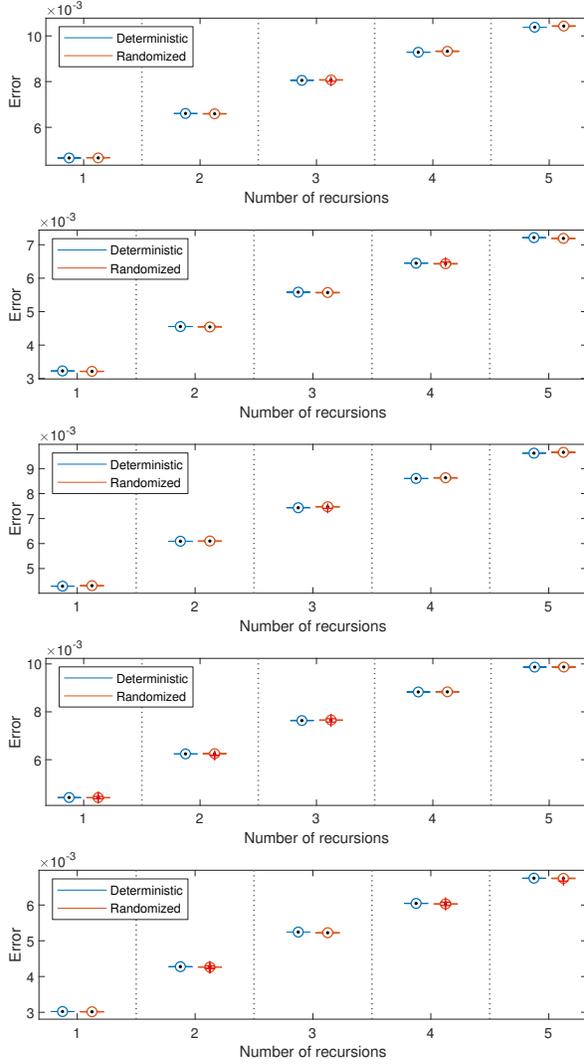


Figure A1.: (Five subplots above) Error for deterministic ABC compared to the error for the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **Gaussian**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

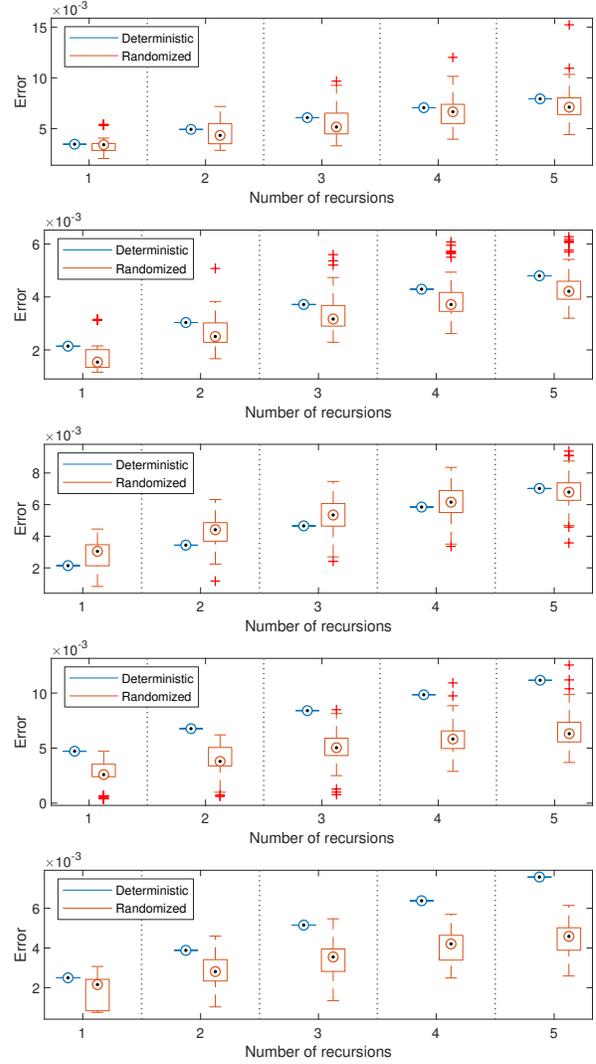


Figure A2.: (Five subplots above) Error for deterministic ABC compared to the error for the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **uniform**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

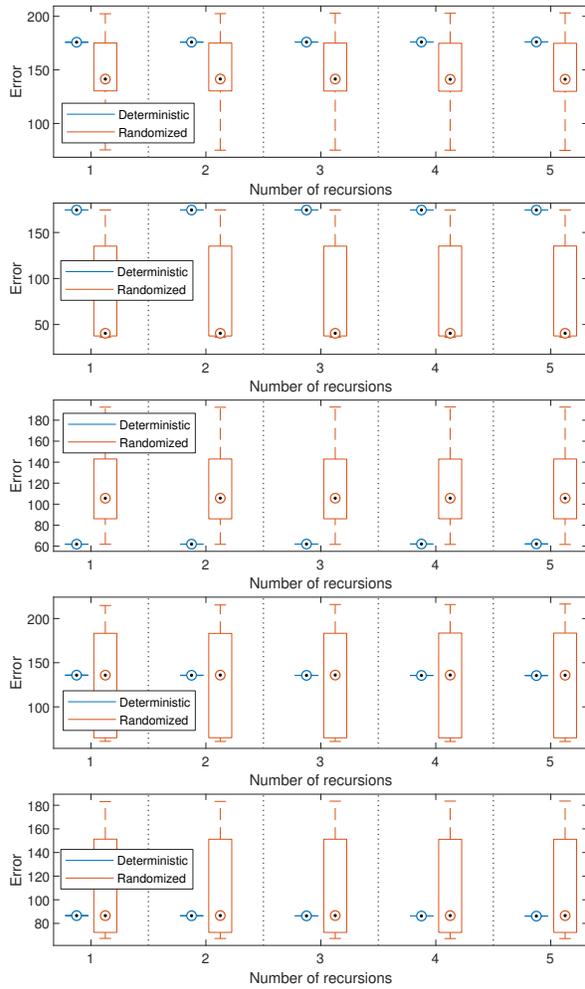


Figure A3.: (Five subplots above) Error for deterministic ABC compared to the error for the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 1 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

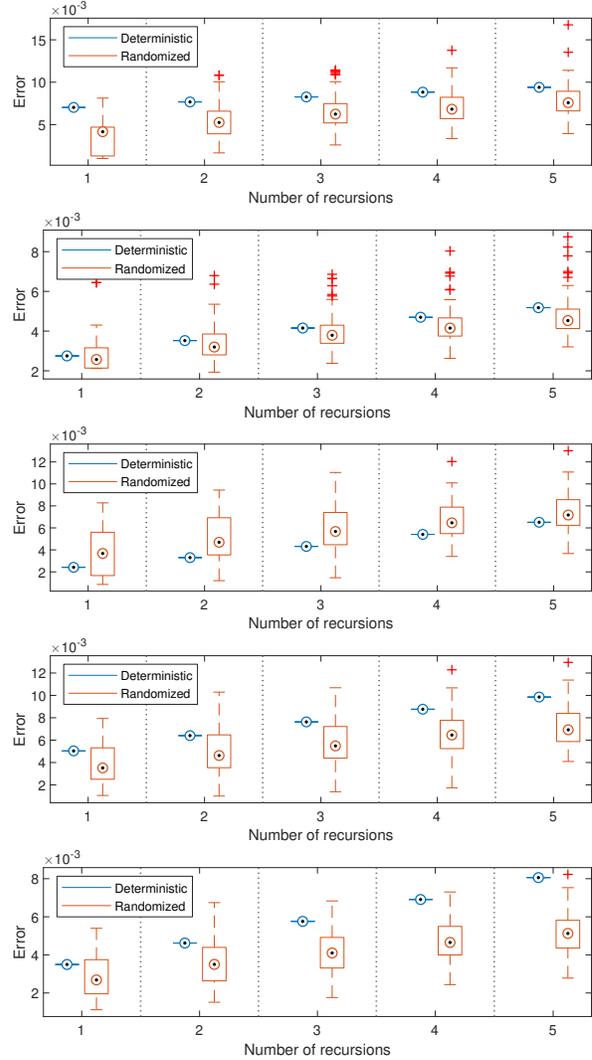


Figure A4.: (Five subplots above) Error for deterministic ABC compared to the error for the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 2 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

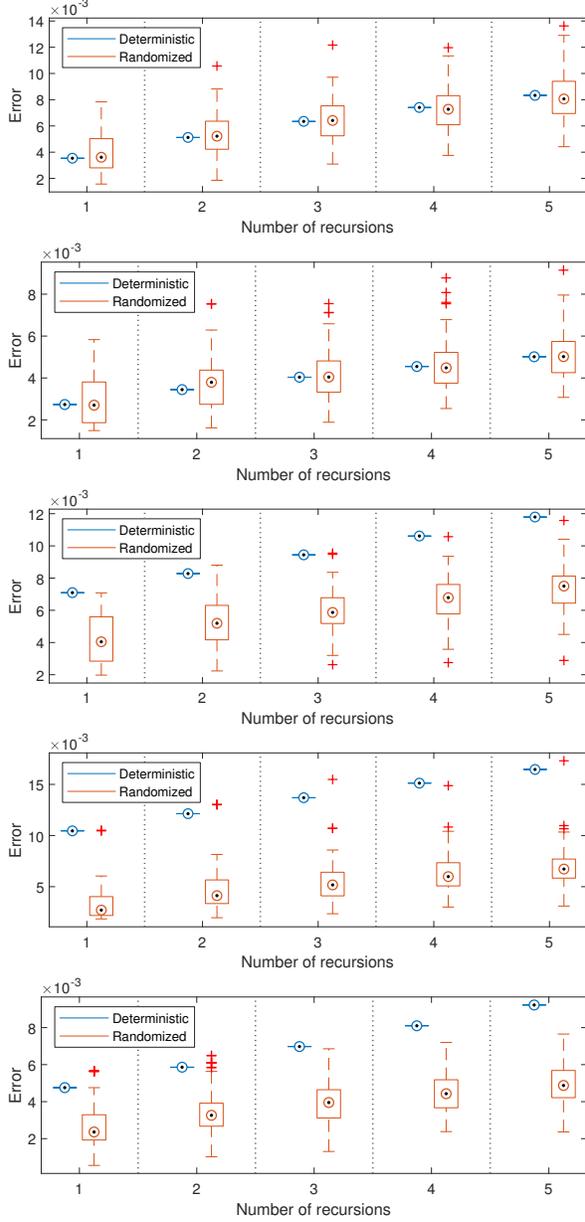


Figure A5.: (Five subplots above) Error for deterministic ABC compared to the error for the randomized counterpart, over 100 realizations of the randomized algorithm. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 3 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

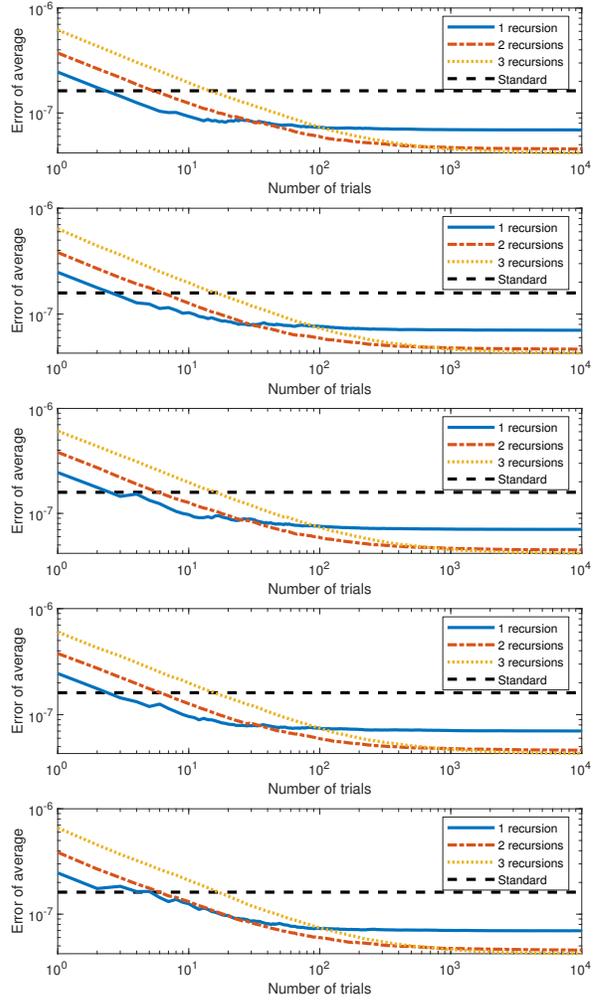


Figure A6.: (Five subplots above) Error for average of randomized EBC compared to the standard algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are **Gaussian**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

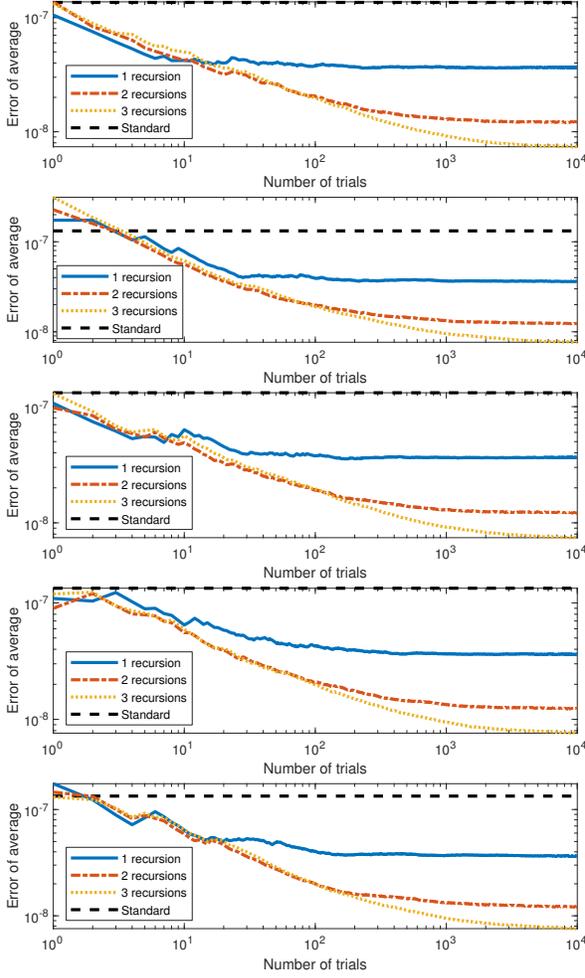


Figure A7.: (Five subplots above) Error for average of randomized EBC compared to the standard algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are **uniform**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

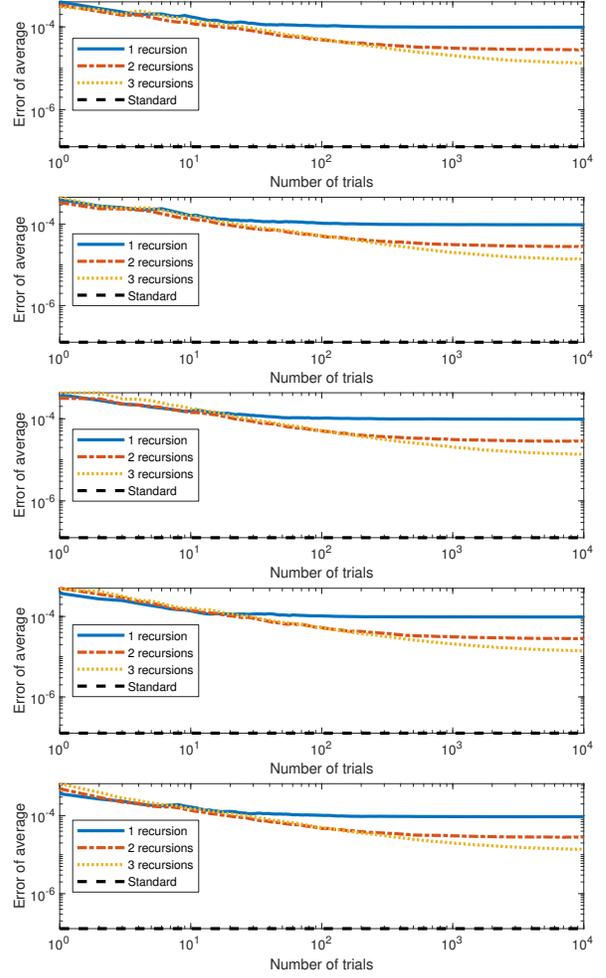


Figure A8.: (Five subplots above) Error for average of randomized EBC compared to the standard algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are **type 1 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

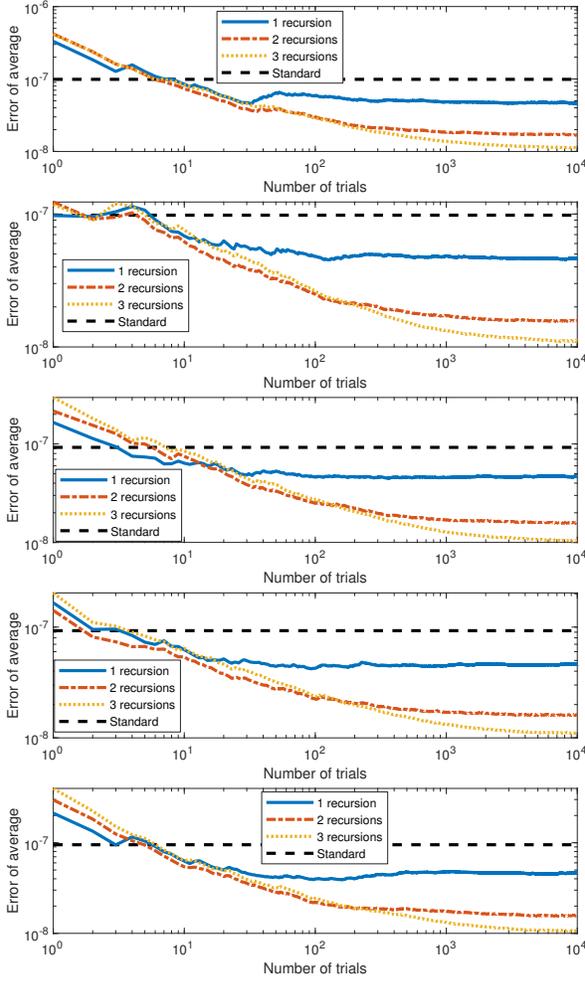


Figure A9.: (Five subplots above) Error for average of randomized EBC compared to the standard algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are **type 2 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

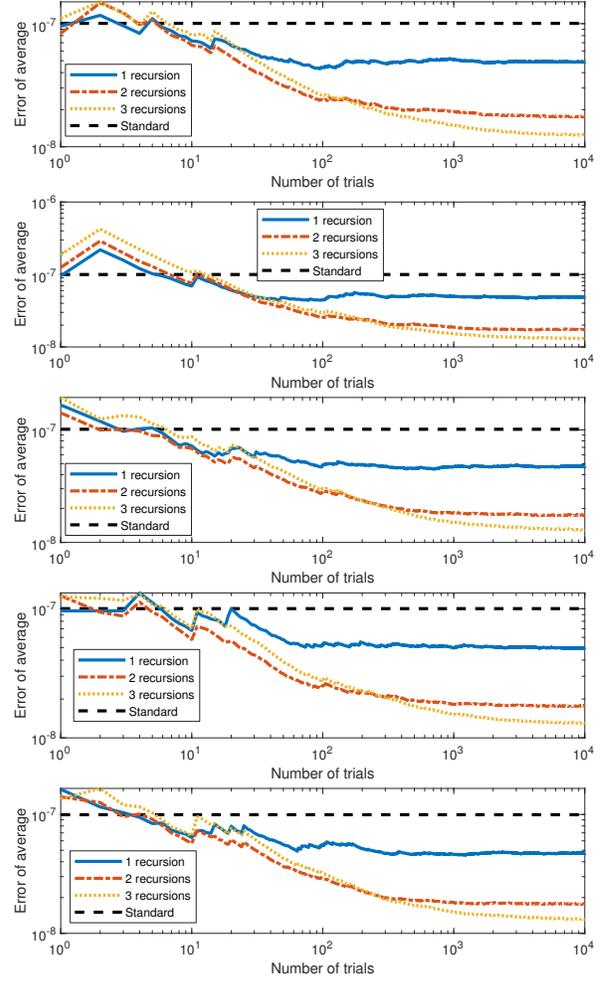


Figure A10.: (Five subplots above) Error for average of randomized EBC compared to the standard algorithm in single precision floating point arithmetic. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{80 \times 80}$ are **type 3 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

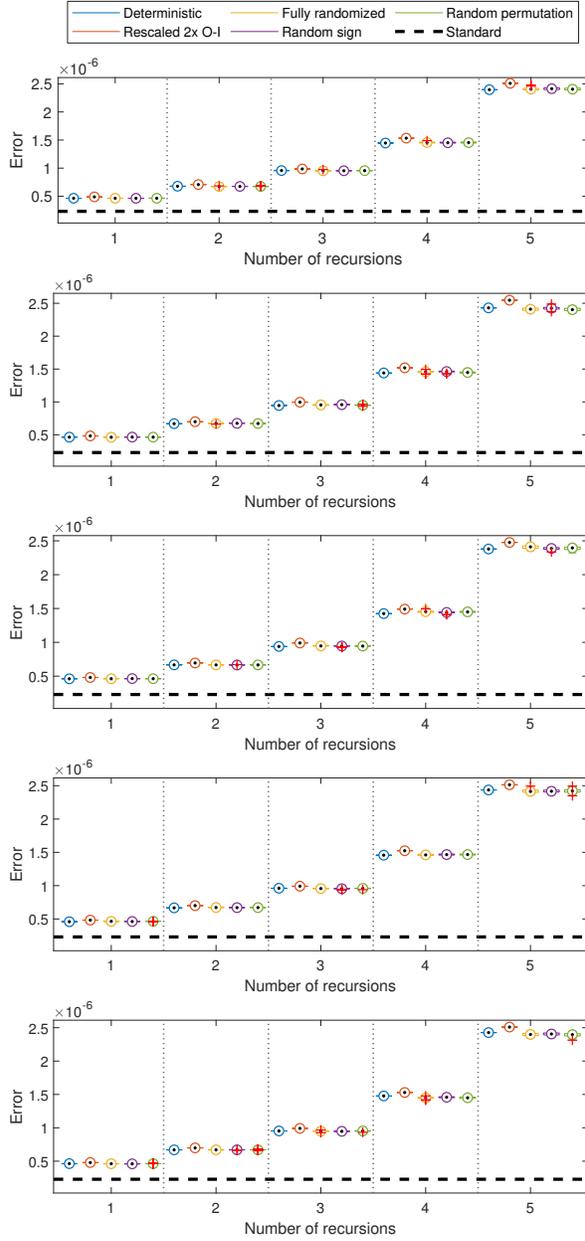


Figure A11.: (Five subplots above) Error for different variants of the EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **Gaussian**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

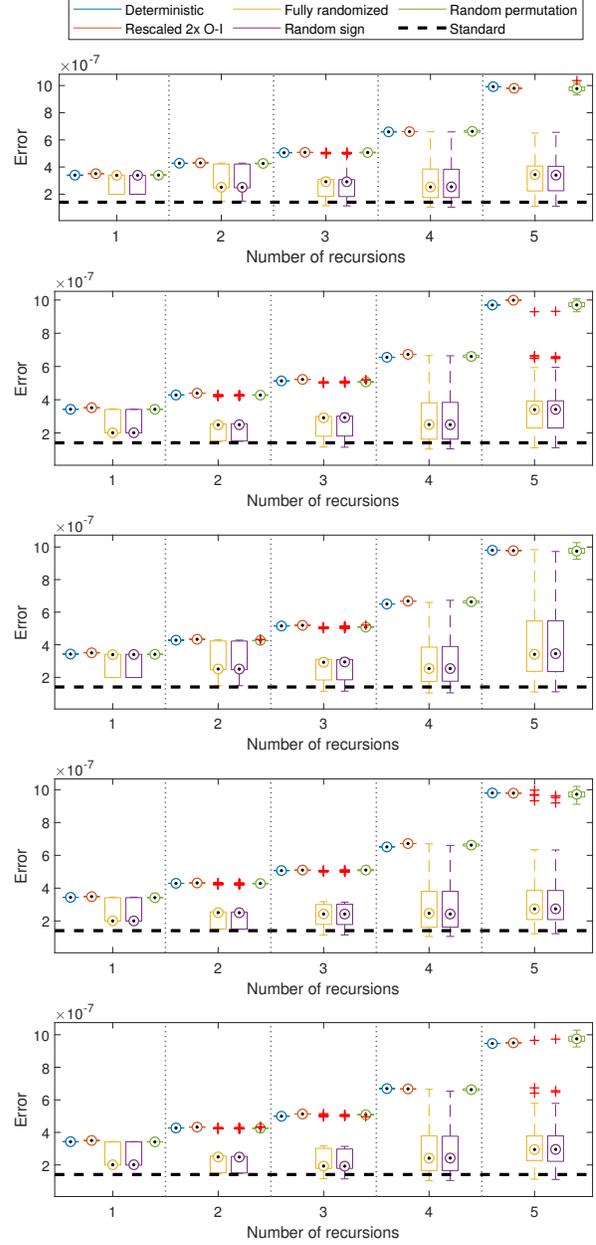


Figure A12.: (Five subplots above) Error for different variants of the EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **uniform**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

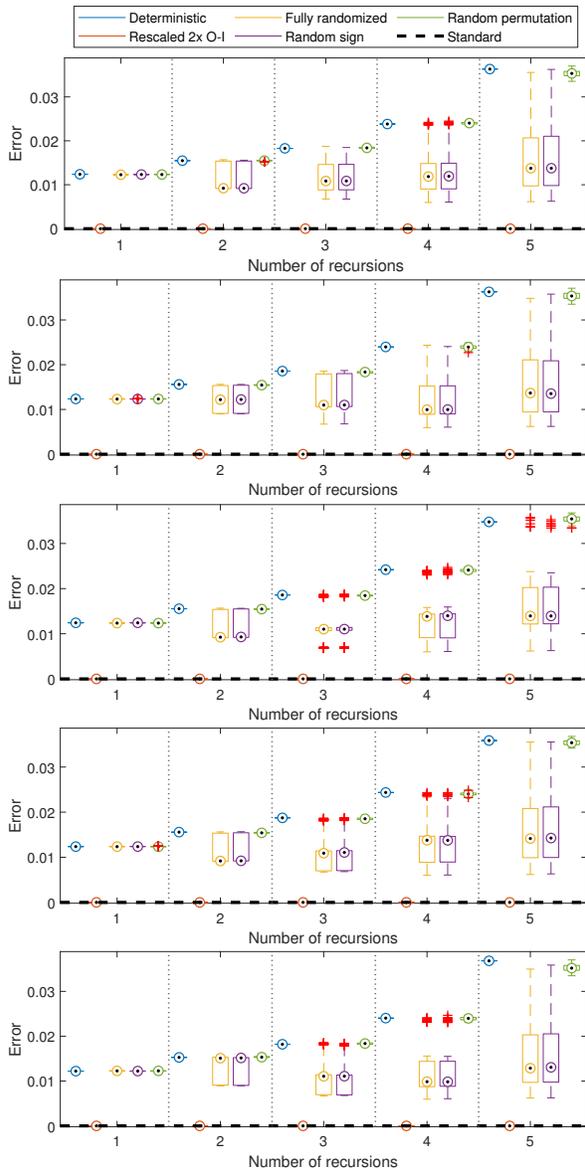


Figure A13.: (Five subplots above) Error for different variants of the EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 1 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

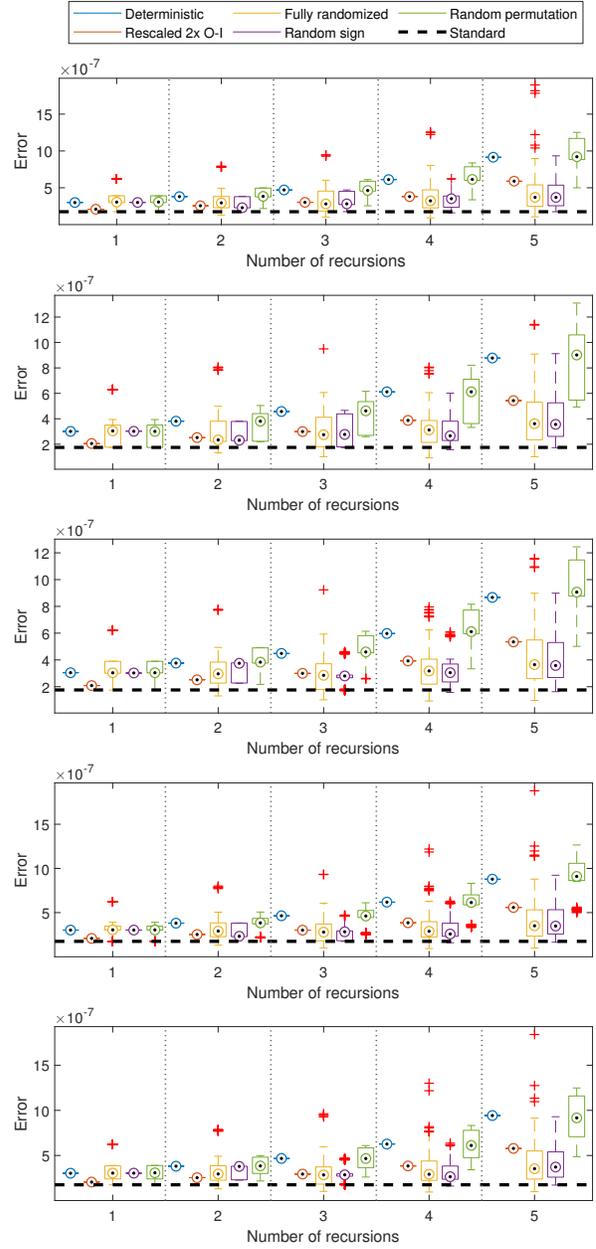


Figure A14.: (Five subplots above) Error for different variants of the EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 2 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .

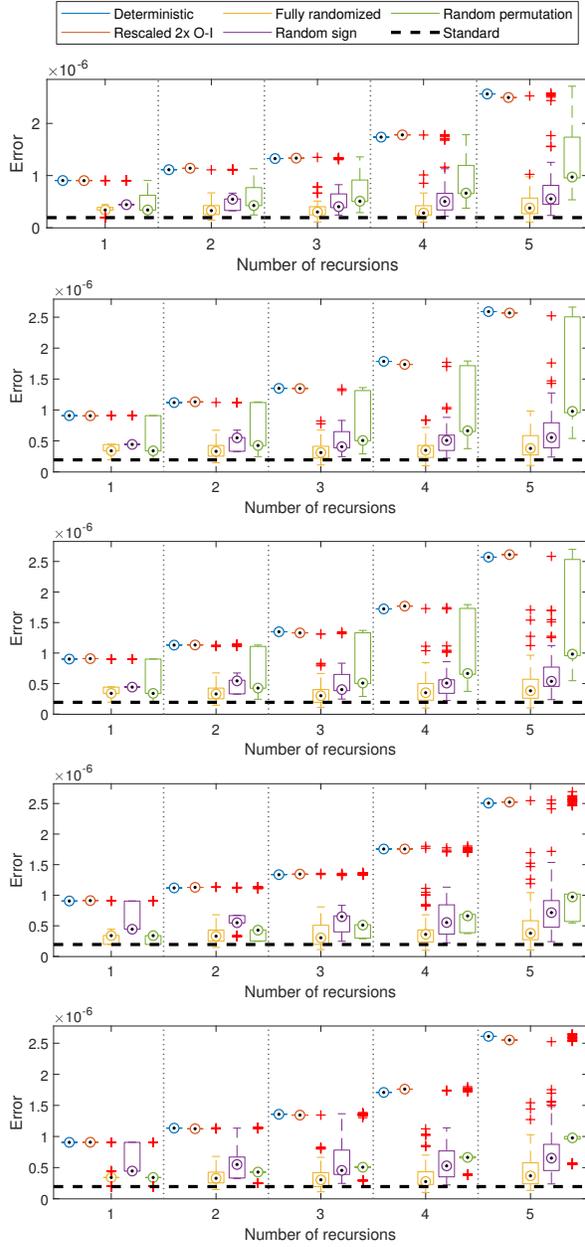


Figure A15.: (Five subplots above) Error for different variants of the EBC in single precision floating point arithmetic, over 100 realizations of the randomized algorithms. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{320 \times 320}$ are **type 3 adversarial**, and each subplot corresponds to one realization of the pair (\mathbf{A}, \mathbf{B}) .