



Accurate and Efficient Algorithms for Floating Point Computation*

J. Demmel[†] and *P. Koev*[‡]

1 Abstract

Our goal is to find accurate and efficient algorithms, when they exist, for evaluating rational expressions containing floating point numbers, and for computing matrix factorizations (like LU, the singular value decomposition (SVD) and eigenvalue decompositions) of matrices with rational expressions as entries. More precisely, *accuracy* means the relative error in the output must be less than one (no matter how tiny the output is), and *efficiency* means that the algorithm runs in polynomial time. Our goal is challenging because our accuracy demand is much stricter than usual.

The classes of floating point expressions or matrices that we can accurately and efficiently evaluate or factor depend strongly on our model of arithmetic:

1. In the “Traditional Model” (TM), the floating point result of an operation like $a + b$ is $fl(a + b) = (a + b)(1 + \delta)$, where $|\delta|$ must be tiny.
2. In the “Long Exponent Model” (LEM) each floating point number $x = f \cdot 2^e$ is represented by the pair of integers (f, e) , and there is no bound on the sizes of the exponents e in the input data. The LEM permits accurate and efficient computation of strictly larger classes of expressions or matrices than the TM.
3. In the “Short Exponent Model” (SEM) each floating point number $x = f \cdot 2^e$ is also represented by (f, e) , but the input exponent sizes are bounded in terms

*Funding for this paper was furnished by NSF grant ACI-9813362 and DOE grant DE-FG03-94ER25219. The information presented here does not necessarily reflect the position of the Government and no official endorsement should be inferred.

[†]Mathematics Department and Computer Science Division, University of California, Berkeley, CA 94720, USA. E-mail: demmel@cs.berkeley.edu

[‡]Department of Mathematics, MIT, Cambridge, MA 02139, USA. E-mail: plamen@math.mit.edu



of the sizes of the input fractions f . We believe the SEM permits accurate and efficient computation of strictly more expressions or matrices than the LEM.

These classes will be described by factorizability properties of the rational expressions, or of the minors of the rational matrices. For each such class, we identify new algorithms that attain our goals of accuracy and efficiency. These algorithms are often exponentially faster than prior algorithms, which would simply use a conventional algorithm with sufficiently high precision.

For example, we can factorize Cauchy matrices, Vandermonde matrices, many kinds of totally positive matrices, and suitably discretized differential and integral operators in all three models much more accurately and efficiently than before. But we provably cannot add three numbers accurately in the TM, even though it is easy in the other models.

2 Introduction

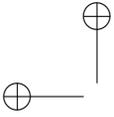
We will survey recent progress and describe open problems in the area of accurate floating point computation, in particular for matrix computations. A very short bibliography would include [18, 12, 15, 31, 24, 4, 17, 23, 5].

We consider the evaluation of multivariate rational functions $r(x)$ of floating point numbers $x = (x_1, x_2, \dots)$, and matrix computations on rational matrices $A(x)$, where each entry $A_{ij}(x)$ is such a rational function. Matrix computations will include computing determinants (and other minors), linear equation solving, performing Gaussian Elimination (GE) with various kinds of pivoting, computing the singular value decomposition (SVD), and computing eigenvalues, among others. Our goals are *accuracy* (computing each solution component with tiny relative error) and *efficiency* (the algorithm should run in time bounded by a polynomial function of the input size).

We consider three models of arithmetic, defined in the abstract, and for each one we try to classify rational expressions and matrices as to whether they can be evaluated or factored accurately and efficiently (we will say “compute(d) accurately and efficiently,” or “CAE” for short).

In the Traditional “ $1 + \delta$ ” Model (TM), we have $fl(a \otimes b) = (a \otimes b)(1 + \delta)$, $\otimes \in \{+, -, \times, \div\}$ and $|\delta| \leq \epsilon$, where $\epsilon \ll 1$ is called *machine precision*. It is the conventional model for floating point error analysis, and means that every floating point result is computed with a relative error δ bounded in magnitude by ϵ . The values of δ may be arbitrary real (or complex) numbers satisfying $|\delta| \leq \epsilon$, so that any algorithm proven to CAE in the TM must work for arbitrary real (or complex) number inputs and arbitrary real (or complex) $|\delta| \leq \epsilon$. The size of the input in the TM is the number of floating point words needed to describe it, independent of ϵ .

The Long Exponent (LEM) and Short Exponent (SEM) models, which are implementable on a Turing machine [34], make errors that may be described by the TM, but their inputs and δ 's are much more constrained. Also, we compute the size of the input in the LEM and SEM by counting the number of bits, so that higher precision and wider range take more bits.



It will turn out that problems we can provably CAE in the TM are a strict subset of those we can CAE in the LEM, which in turn we conjecture are a strict subset of those we can CAE in the SEM. In all three models we will describe the classes of rational expressions and rational matrices in terms of the factorization properties of the expressions, or of the minors of the matrices.

The reader may wonder why we insist on accurately computing tiny quantities with small relative error, since in many cases the inputs are themselves uncertain, so that one could suspect that the inherent uncertainty in the input could make even the signs of tiny outputs uncertain. First, in a number of practical problems (eg computing vibrational frequencies or energy levels) it is the smallest eigenvalues that are of physical interest. Second, it will turn out that in the TM, the class we can CAE appears to be identical to the class where all the outputs are in fact accurately determined by the inputs, in the sense that small relative changes in the inputs cause small relative changes in the outputs. In other words, it is worthwhile computing the output with some relative accuracy, because that is the accuracy to which it is determined by the input. We discuss this conjecture in section 6 below.

There are many ways to formulate the search for efficient and accurate algorithms [11, 8, 41, 7, 29, 40, 38]. Our approach differs in several ways. In contrast to either conventional floating point error analysis [29] or the model in [11], we ask that even the tiniest results have correct leading digits, and that zero be exact. In [11] the model of arithmetic allows a tiny absolute error in each operation, whereas in TM we allow a tiny relative error. Unlike [11] our LEM and SEM are conventional Turing machine models, with numbers represented as bit strings, and so we can take into account the cost of arithmetic on very large and very small numbers (i.e. those with many exponent bits). For these reasons we believe our models are closer to computational practice than the model in [11]. In contrast to [38], we (mostly) consider the input as given exactly, rather than as a sequence of ever better approximations. Finally, many of our algorithms could easily be modified to explicitly compute guaranteed interval bounds on the output [40].

3 Factorizability and Minors

We show here how to reduce the question of accurate and efficient matrix computations to accurate and efficient rational expression evaluation. The connection to LU factorization and similar operations is elementary, but the SVD requires an algorithm from [18]. (Here an accurate SVD means that the singular values are known with guaranteed relative error less than one, but the error in a singular vector u_i may be inversely proportional to the *relative gap* $\text{relgap}(i) \equiv \min_{j \neq i} |\sigma_i - \sigma_j| / \sigma_i$ between its corresponding singular value σ_i and the other singular values σ_j .)

Proposition 1. *Being able to CAE the absolute value of the determinant $|\det(A(x))|$ is necessary to be able to CAE the following matrix computations on $A(x)$: LU factorization (with or without pivoting), QR factorization, all the eigenvalues λ_i of $A(x)$, and all the singular values of $A(x)$. Conversely, being able to CAE all the minors of $(A(x))$ is sufficient to be able to CAE the following matrix computations*



on $A(x)$: A^{-1} , LU factorization (with or without pivoting), and the SVD of $A(x)$. This holds in any model of arithmetic.

Proof. First consider necessity. $|\det(A(x))|$ may be written as the product of diagonal entries of the matrices L , U and R in these factorizations, or as the product of eigenvalues or singular values. If these entries or values can be CAE, then so can their product in a straightforward way.

Now consider sufficiency. The statement about A^{-1} is just Cramer's rule, which only needs $n^2 + 1$ different minors. The statement about LU factorization depends on the fact that each nontrivial entry of L and U is a quotient of minors. The SVD is more difficult [18], and depends on the following two step algorithm: (1) Compute a *rank revealing* decomposition $A = X \cdot D \cdot Y$ where X and Y are "well-conditioned" (far from singular in the sense that $\|X\| \cdot \|X^{-1}\|$ is not too large) and (2) use a bisection-like algorithm to compute the SVD from XDY . \square

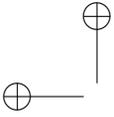
Computing $\det(A(x))$ accurately is obviously necessary for LU factorization and eigenvalues, but not necessarily for QR factorization or the SVD. The sufficiency proof can be extended to other matrix computations like the QR decomposition and pseudoinverse by considering minors of matrices like $\begin{bmatrix} I & A \\ A^T & 0 \end{bmatrix}$. Furthermore, if we can CAE the minors of $C(x) \cdot A(x) \cdot B(x)$, and $C(x)$ and $B(x)$ are well-conditioned, then we can still CAE a number of matrix factorizations, like the SVD.

The SVD can be applied to get the eigendecomposition of symmetric matrices [25], but efficiently computing accurate eigenvalues of nonsymmetric matrices seems to require more conditions. We have recently shown that computing a small set of minors accurately is sufficient for computing accurate singular values *and* eigenvalues of totally positive matrices, and that if these minors can be computed efficiently then so can their singular values and eigenvalues [33]. We discuss this further in section 5 below.

4 CAE in the Traditional Model

We begin by giving examples of expressions and matrix computations that we can CAE in the TM, and then discuss what we cannot do. The results will depend on details of the axioms we adopt, but for now we consider the minimal set of operations described in the abstract.

As long as we only do *admissible operations*, namely multiplication, division, addition of like-signed quantities, and addition/subtraction of (exact!) input data ($x \pm y$), then the worst case relative error only grows very slowly, roughly proportionally to the number of operations. It is when we subtract two like-signed approximate quantities *and* significant cancellation occurs, that the relative error can become large. So we may ask which problems we can CAE just using only admissible operations, i.e. which rational expressions factor in such a way that only admissible operations are needed to evaluate them, and which matrices have all minors with the same property.



Here are some examples [15] where we assume that the inputs are arbitrary real or complex numbers. (1) The determinant of a Cauchy matrix $C_{ij} = 1/(x_i + y_j)$ is CAE using the classical expression $\prod_{i < j} (x_j - x_i)(y_j - y_i) / \prod_{i, j} (x_i + y_j)$, as is every minor. In fact, changing one line of the classical GE routine will compute each entry of the LU decomposition accurately in about the same time as the original inaccurate version. (2) We can CAE all minors of sparse matrices, i.e. those with certain entries fixed at 0 and the rest independent indeterminates x_{ij} , if and only if the undirected bipartite graph presenting the sparsity structure of the matrix is *acyclic*; a one-line change to GE again renders it accurate. An important special case are bidiagonal matrices, which arise in the conventional SVD algorithm. (3) The eigenvalue problem for the second centered difference approximation to a Sturm-Liouville ODE or elliptic PDE on a rectangular grid (with arbitrary rectilinear boundaries) can be written as the SVD of an “unassembled” problem $G = D_1 U D_2$ where D_1 and D_2 are diagonal (depending on “masses” and “stiffnesses”) and U is *totally unimodular*, i.e. all its minors are ± 1 or 0. Again, a simple change to GE renders it accurate.

In contrast, one can show that it is impossible in the TM to add three numbers $x + y + z$ accurately in constant time; the proof involves showing that for *any* algorithm the rounding errors δ and inputs x, y, z can be chosen to have an arbitrarily large relative error [18]. This depends on the δ 's being permitted to be arbitrary real numbers in our model. This impossibility is in stark contrast to the bit model, although computing accurate floating sums as efficiently as possible is more subtle [16].

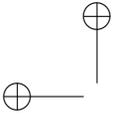
Vandermonde matrices $V_{ij} = x_i^{j-1}$ are also more subtle. Since the product of a Vandermonde matrix and the Discrete Fourier Transform (DFT) is Cauchy, and we can compute the SVD of a Cauchy, we can compute the SVD of a Vandermonde [13]. This fits in our TM model because the roots of unity in the DFT need only be known approximately, and so may be computed in the TM model. In contrast, one can use the result in the last paragraph to show that the inverse of a Vandermonde cannot be computed accurately. Similarly, polynomial Vandermonde matrices with $V_{ij} = P_i(x_j)$, P_i a (normalized) orthogonal polynomial, also permit accurate SVDs [19] but probably not inverses.

These results, with others to be discussed later, are summarized in Table 1.

5 Adding Nonnegativity to the Traditional Model

If we further restrict the domain of (some) inputs to be nonnegative, then many more accurate and efficient computations are possible, $x + y + z$ as a trivial example. A more interesting example is the set of weakly diagonally dominant M-matrices, which arise as discretizations of PDEs; they must be represented as offdiagonal entries and the row sums [31, 20].

More interesting is the class of *totally positive (TP) matrices*, all of whose minors are positive. Numerous structure theorems show how to represent such matrices as products of much simpler TP matrices [15, Sec. 9] [1, 30]. Accurate formulas for the (nonnegative) minors of these simpler matrices combined with the



Cauchy-Binet theorem yield accurate formulas for the minors of the original TP matrix, but typically at an exponential cost.

The most important structure theorem states that all TP matrices have a *bidiagonal or Neville factorization*, an LDU factorization where L (U) is the product of unit lower (upper) bidiagonal nonnegative matrices, and D is diagonal and positive. There are n^2 free positive parameters defining these matrices that parameterize the set of all TP matrices. If these positive parameters, which are products and quotients of minors of the matrix, can be computed accurately and efficiently, then many subsequent linear algebra operations (like the SVD, solving $Ax = b$, inversion, and the eigenvalue problem) can also be solved accurately and efficiently [33]. This all applies to nonsingular totally nonnegative, oscillatory, and sign-regular matrices (inverses of TP matrices) [26].

In fact, the set of totally positive matrices for which it is possible to compute the SVD and eigenvalues (as well as inverses, etc.) accurately and efficiently is closed under a large set of operations: Given any two TP matrices A and B with Neville factorizations, the Neville factorization of $A \cdot B$, the scaled inverse $SA^{-1}S$ (where $S = \text{diag}(\text{sign}(\pm 1))$) and Schur complements of A can be CAE, so that the SVDs and eigenvalues of these new matrices can be CAE [32].

A great many TP matrices have efficient and accurate algorithms for computing their Neville factorizations. For example, the Björck-Pereyra method [6] was invented for Vandermonde systems, and has been extended to Cauchy [9], Cauchy-Vandermonde [36], and generalized Vandermonde matrices [22].

But for TP generalized Vandermonde matrices $G_{ij} = x_i^{\mu_j}$, where the μ_j form an increasing nonnegative sequence of integers, the problem is much more subtle. $\det(G)$ is known to be the product of $\prod_{i < j} (x_j - x_i)$ and a *Schur function* [35] $s_\lambda(x_i)$, where the sequence $\lambda = (\lambda_j) = (\mu_{n+1-j} - (n-j))$ is called a *partition*. Schur functions are polynomials with nonnegative integer coefficients, so since their arguments x_i are nonnegative, they can certainly be computed accurately. However, straightforward evaluation would have an exponential cost $O(n^{|\lambda|})$, $|\lambda| = \sum_j \lambda_j$. But by exploiting combinatorial identities satisfied by Schur functions along with techniques of divide-and-conquer and memoization, the cost of evaluating the determinant (or complete Neville factorization) can be reduced to polynomial time $n^2 \Lambda$ where $\Lambda \leq 2(\lambda + 1) \cdot (\#\lambda_i > 0) \cdot \prod_{j > 1} (\lambda_j + 1)^2$. The λ_i are counted as part of the size of the input in this case [31, 21, 22].

Combined with the results on computing the SVD and eigenvalues in [33, 32], this positively settles an open question in [31, 14]: we can now accurately and efficiently compute the eigenvalues and singular values of TP generalized Vandermonde matrices. Questions remain, however, about the accuracy of the associated vectors.

In [14] we made a conjecture (Conjecture 1) in an attempt to characterize the set of homogeneous polynomials that can be evaluated accurately in the TM. Briefly, we conjectured that $f(x_1, \dots, x_n)$ could be evaluated accurately (on a reasonable subset of the unit sphere) if and only if each of its factors were of the form x_i , $x_i - x_j$, $x_i + x_j$ or bounded away from 0 on the domain of evaluation. These properties are indeed sufficient, but not necessary, as the simple example $x_1^2 + (x_2 + x_3)^2$ shows. Indeed, any formulation of a condition for efficient evaluation must depend on more than the variety determined by f , since $f(x_1, x_2, x_3, x_4) = x_1^4 + x_1^2 \cdot (x_2 + x_3)^2$



and $g(x_1, x_2, x_3, x_4) = x_1^4 + x_1^2 \cdot (x_2 + x_3 + x_4)^2$ determine the same real variety $\{x_1 = 0\}$, yet it seems that only f can be evaluated accurately for all inputs, since any accurate algorithm for g would require the accurate evaluation of $(x_2 + x_3 + x_4)^2$ when x_1 is tiny enough, apparently contradicting the result in [18] cited in Section 4. Necessary conditions remain elusive.

6 Summary of Costs in the Traditional Model

Table 1 below summarizes our state of knowledge of the cost of accurate matrix computations in the TM. There is one row for the most important classes of matrices discussed so far, as well as others. Furthermore, some totally positive (TP) matrix classes are considered as well, since this can completely change the cost of accurate computation. Details of matrix classes are explained in footnotes.

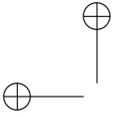
There is one column in Table 1 for each matrix computation we consider. “Any minor” means any single arbitrary minor of the matrix A . GENP, GEPP and GECP stand for Gaussian elimination with no pivoting, partial pivoting, and complete pivoting, resp. SVD and EVD mean the singular value and eigenvalue decompositions, respectively; we only consider the problem of computing all singular value or eigenvalues accurately, the more difficult question of vectors is omitted (see the open questions in section 9). NF refers to Neville (bidiagonal) factorization, which is primarily of interest for TP matrices. Other details are explained in footnotes.

The table entry n^2 means $O(n^2)$, for example. The table entry “exp” means exponential time, based on formulas for evaluating arbitrary Schur polynomials [31, 35] in the case of TP (generalized) Vandermonde matrices, and based on the Cauchy-Binet formula for minors of products of matrices in the case of “any TP”. The table entry “No” means impossible, because it would require the computation of an expression including $x+y+z$, which we know to be impossible. It is interesting that sometimes the SVD is possible in $O(n^3)$ time, but that inversion is impossible (Vandermondes).

Our results have an asterisk (*) in front of the citation. (Sometimes we cite one of our own papers if it is a convenient place to find a concise description of a classical result; there is no asterisk in this case).

7 Extending the Traditional Model

So far we have considered the simplest version of the TM, where (1) we have only the input data, and no additional constants available, (not even integers, let alone arbitrary rationals or reals), (2) the input data is given exactly (as opposed to within a factor of $1 + \delta$), and (3) there is no way to “round” a real number to an integer, and so convert the problem to the LEM or SEM models. We note that in [11], (1) integers are available, (2) the input is rounded, and (3) there is no way to “round” to an integer. Changes to these model assumptions will affect the classes of problems we can solve. For example, if we (quite reasonably) were to permit exact integers as input, then we could CAE expressions like $x - 1$, and otherwise



presumably not. If we went further and permitted exact rational numbers, then we could also CAE $9x^2 - 1 = 9(x - \frac{1}{3})(x + \frac{1}{3})$. Allowing algebraic numbers would make $x^2 - 2 = (x - \sqrt{2})(x + \sqrt{2})$ CAE.

If inputs were not given exactly, but rather first multiplied by a factor $1 + \delta$, then we could no longer always accurately compute $x \pm y$ where x and y are inputs, eliminating Cauchy matrices and most others. But the problems we can solve with exact inputs in the TM still have an attractive property with inexact inputs: Small relative changes in the inputs *often* cause only a small relative change in the outputs, independent of their magnitudes. The output relative errors may be larger than the input relative error by a factor called a *relative condition number* κ_{rel} , which is at most a polynomial function of $\max(1/\text{rel_gap}(x_i, \pm x_j))$. Here $\text{rel_gap}(x_i, \pm x_j) = |x_i \mp x_j| / (|x_i| + |x_j|)$ is the *relative gap* between inputs x_i and $\pm x_j$, and the maximum is taken over all factors $x_i \mp x_j$ dividing the rational expression being computed. So if all the inputs differ in several of their leading digits, all the leading digits of the outputs are determined accurately. We note that κ_{rel} can be large, depending on the expression being evaluated, but it can only be unbounded when a relative gap goes to zero.

If a problem has this attractive property, we say that it possesses a relative perturbation theory. In practical situations, where only a few leading digits of the inputs x_i are known, this property justifies the use of algorithms that try to compute the output as accurately as we do.

In [14] we stated a conjecture (Conjecture 2) similar to Conjecture 1, that a relative perturbation theory for the evaluation of f existed if and only if f factored in a certain simple way. As with Conjecture 1, these conditions are in fact sufficient but not necessary for a relative perturbation theory to exist.

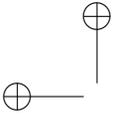
8 CAE in the Long and Short Exponent Models

Now we consider standard Turing machines, where input floating point numbers $x = f \cdot 2^e$ are stored as the pair of integers (f, e) , so the size of x is $\text{size}(x) = \#\text{bits}(f) + \#\text{bits}(e)$. We distinguish two cases, the Long Exponent Model (LEM) where f and e may each be arbitrary integers, and the Short Exponent Model (SEM), where the length of e is bounded depending on the length of f . In the simplest case, when $e = 0$ (or lies in a fixed range) then the SEM is equivalent to taking integer inputs, where the complexity of problems is well understood. This is more generally the case if $\#\text{bits}(e)$ grows no faster than a polynomial function of $\#\text{bits}(f)$.

In particular it is possible to CAE the determinant of an integer (or SEM) matrix each of whose entries is an independent floating point number [10]. This is not possible as far as we know in the LEM (even for tridiagonal matrices with independent floating point entries) which accounts for a large complexity gap between the two models.

We start by illustrating some differences between the LEM and SEM, and then describe the class of problems that we can CAE in the LEM.

First, the number of bits in an expression with LEM inputs can be exponen-



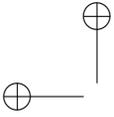
tially larger than the number of bits in the same expression when evaluated with SEM inputs. For example, the number of nonzero bits in $\prod_{i=1}^n (1 + 2^{e_i})$ grows like $O(n)$ if the e_i are bounded, but like 2^n if the e_i are arbitrary. In fact, LEM arithmetic can encode symbolic algebra, because if e_1 and e_2 have no overlapping bits, then we can recover e_1 and e_2 from their product $2^{e_1} \cdot 2^{e_2} = 2^{e_1+e_2}$, letting us effectively represent the indeterminates x_1 and x_2 by 2^{e_1} and 2^{e_2} , resp.

Second, the error of many conventional matrix algorithms is typically proportional to the condition number $\kappa(A) = \|A\| \cdot \|A^{-1}\|$. This means that a conventional algorithm run with $O(\log \kappa(A))$ extra bits of precision will compute an accurate answer. It turns out that if $A(x)$ has rational entries in the SEM model, then $\log \kappa(A)$ is at most a polynomial function of the input size, so conventional algorithms run in high precision will CAE the answer. However $\log \kappa(A)$ for LEM matrices can be exponentially larger, so this approach does not work. The simplest example is $\log \kappa(\text{diag}(1, 2^e)) = e = 2^{\#\text{bits}(e)}$. On the other hand $\log \log \kappa(A(x))$ is a lower bound on the complexity of any algorithm, because this is a lower bound on the number of exponent bits in the answer. One can show that $\log \log \kappa(A(x))$ grows at most polynomially large in the size of the input [18].

Finally, we consider the problem of computing an arbitrary bit in the simple expression $p = \prod_{i=1}^n (1 + x_i)$. When the x_i are in the SEM, then p can be computed exactly in polynomial time. However when the x_i are in the LEM, then one can prove that computing an arbitrary bit of p is as hard as computing the permanent, a well-known combinatorially difficult problem. Here is another apparently simple problem not known to even be in NP: testing singularity of a floating point matrix. In the SEM, we can CAE the determinant. But in the LEM, the obvious choice of a “witness” for singularity, a null vector, can have exponentially many bits in it, even if the matrix is just tridiagonal. We conjecture that deciding singularity of an LEM matrix is NP-hard.

So how do we compute efficiently in the LEM? The idea is to use *sparse arithmetic*, or to represent only the nonzero bits in the number. (A long string of 1s can be represented as the difference of two powers of 2 and similarly compressed). In contrast, in the SEM one uses *dense arithmetic*, storing all fraction bits of a number. For example, in sparse arithmetic $2^e + 1$ takes $O(\log e)$ bits to store in sparse arithmetic, but e bits in dense arithmetic. This idea is exploited in practical floating point computation, where extra precise numbers are stored as arrays of conventional floating point numbers, with possibly widely different exponents [39].

Now we describe the class of rational functions that we can CAE in the LEM. We say the rational function $r(x)$ is in *factored form* if $r(x) = \sum_{i=1}^n p_i(x_1, \dots, x_k)^{e_i}$, where each e_i is an integer, and $p_i(x_1, \dots, x_k)$ is written as an explicit sum of nonzero monomials. We say $\text{size}(r)$ is the number of bits needed to represent it in factored form. Then by (1) computing each monomial in each p_i exactly, (2) computing the leading bits of their sum p_i using sparse arithmetic (the cost is basically sorting the bits [16]), and (3) computing the leading bits of the product of the $p_i^{e_i}$ by conventional rounded multiplication or division, one can evaluate $r(x)$ accurately in time a polynomial in $\text{size}(r)$ and $\text{size}(x)$. In other words, the class of rational expression that we can CAE are those that we can express in factored form in polynomial space.



Now we consider matrix computations. It follows from the last paragraph that if each minor $r(x)$ of $A(x)$ can be written in factored form of a size polynomial in the size of $A(x)$, then we can CAE all the matrix computations that depend on minors. So the question is which matrix classes $A(x)$ have all their minors (or just the ones needed for a particular matrix factorization) expressible in a factored form no more than polynomially larger than the size of $A(x)$. The obvious way to write a minor like $\det(A(x))$, with the Laplace expansion of $n!$ terms, is clearly exponentially larger than $A(x)$, so it is only specially structured $A(x)$ that will work.

All the matrices that we could CAE in the TM are also possible in the LEM. The most obvious classes of $A(x)$ that we can CAE in the LEM that were impossible in the TM are gotten by replacing all the indeterminates in the TM examples by arbitrary rational expressions of polynomial size. For example, the entries of an M-matrix can be polynomial-sized rational expressions in other quantities. Another class are Green's matrices (inverses of tridiagonals), which can be thought of as discretized integral operators, with entries written as $A_{ij} = x_i \cdot y_j$.

The obvious question is whether A each of whose entries is an independent number in the LEM falls in this class. We conjecture that it does not, as mentioned before.

9 Conclusions and Open Problems

Our goal has been to identify rational expressions (or matrices) that we can evaluate accurately (or on which we can perform accurate matrix computations), in polynomial time. Accurately means that we want to get a relative error less than 1, and polynomial time means in a time bounded by a polynomial function of the input size.

We have defined three reasonable models of arithmetic, the Traditional Model (TM), the Long Exponent Model (LEM) and the Short Exponent Model (SEM), and tried to identify the classes of problems that can or cannot be computed accurately and efficiently for each model. The TM can be used as a model to do proofs that also hold in the implementable LEM and SEM, but since it ignores the structure of floating point numbers as stored in the computer, it is strictly weaker than either the LEM or SEM. In other words, there are problems (like adding $x + y + z$) that are provably impossible in the TM but straightforward in the other two models.

We also believe that the LEM is strictly weaker than the SEM, in the sense that there appear to be computations (like computing the determinant of a general, or even tridiagonal, matrix) that are possible to do accurately in polynomial time in the SEM but not in the LEM. In the SEM, essentially all problems that can be written down in polynomial space can be solved accurately in polynomial time. For the LEM, only expressions that can be written in *factored form* in polynomial space can be computed efficiently in polynomial time.

A number of open problems and conjectures were mentioned in the paper, including the question marks in Table 1. We mention three additional ones here.

1. It is known that a small relative change $\eta \ll 1$ in any entry of a bidiagonal



matrix can only perturb a singular vector by approximately $\eta/\text{relgap}(i)$ (see the beginning of section 3). If this same result were true for the product of two nonnegative bidiagonals, namely that a small relative change η in any entry could only perturb a singular vector of the product by $\eta/\text{relgap}(i)$, then we could conclude that the singular vectors of totally positive matrices in Table 1 could be computed as accurately as expected, namely with accuracy proportional to machine precision divided by $\text{relgap}(i)$.

2. What can be said about the nonsymmetric eigenvalue problem, other than for totally positive matrices, all of which have positive and simple eigenvalues? In other words, what matrix properties, perhaps related to minors, guarantee that all eigenvalues of a nonsymmetric matrix can be computed accurately?
3. How do the complexity classes change when we consider randomness, either in the problems considered, or the algorithms, and tolerate a small probability of large relative error?

Acknowledgements

The authors acknowledge Benjamin Diamant, Zlatko Drmač, Stan Eisenstat, Ming Gu, Yozo Hida, William Kahan, Ivan Slapničar and Kresimir Veselić for their collaboration over many years in developing this material.

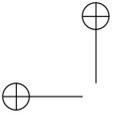


Table 1, Part 1 - Cost of Accurate Computation in TM

Type of Matrix	$\det A$	A^{-1}	Any minor	GENP	GEPP	GECP
Cauchy	n^2 [27]	n^2 [27]	n^2 [27]	n^2 [12]	n^3 [12]	n^3 *[12]
TP Cauchy	n^2 [27]	n^2 [27]	n^2 [27]	n^3 [12]	n^3 [12]	n^3 *[12]
Vandermonde	n^2 [6]	No *[18]	No *[18]	No *[18]	No *[18]	No *[18]
TP Vandermonde ⁷⁾	n^2 [6]	n^3 *[28, 31]	exp *[31]	n^2 *[31]	n^2 *[31]	exp *[31]
Confluent Vandermonde	n^2 [28]	No *[18]	No *[18]	No *[18]	No *[18]	No *[18]
TP Confluent Vandermonde	n^2 [28]	n^3 [28]	? [28]	n^3 *[31]	? [28]	? [28]
Three-term Vandermonde Orth. Poly.	n^2 [28]	? [28]	? [28]	? [28]	? [28]	? [28]
3-term Vand. Orth. Poly. + other cond. ⁵⁾	n^2 [28]	n^3 [28]	? [28]	? [28]	? [28]	? [28]
Generalized Vandermonde	No *[18]	No *[18]	No *[18]	No *[18]	No *[18]	No *[18]
TP generalized Vandermonde ⁴⁾	Λn^2 *[21, 22]	Λn^3 *[21, 22]	exp *[21, 22]	Λn^2 *[21, 22]	Λn^2 *[21, 22]	exp *[21, 22]
Any TP ⁶⁾	n [33]	n^3 *[33]	exp *[33]	n^3 *[33]	exp *[33]	exp *[33]
Wkly diag.dom. M-matrices	n^3 [37]	n^3 [3, 2]	No *[18]	n^3 *[20]	n^3 *[20]	n^3 *[20]

- 1) Each entry is meant in a big-O sense, i.e. n^2 means $O(n^2)$.
- 2) Below each entry is a citation to where the result (old or new) appears.
- 3) Our results have an asterisk (*) in front of the citation.
- 4) If $G_{ij} = x_i^{j-1+\lambda_{n-j}}$ then $\Lambda \leq 2(\lambda_1 + 1)(\lambda_2 + 1)^2 \dots (\lambda_p + 1)^2 p$.
- 5) Also need $0 < x_1 < x_2 < \dots < x_n$, $d_i > 0$, $c_i \geq 0$ for all i , and $b_i - x_k \leq 0$ for all $i + k \leq n - 1$.
- 6) Any TP matrix, specified by the entries of its bidiagonal decomposition, or by combining such matrices by taking products, Schur complements and scaled inverses as described in Section 5.
- 7) The contribution to computing the inverse of a TP Vandermonde in [31] improves the complexity of the result in [28] from $2n^3$ to $\frac{5}{12}n^3$.

Table 1, Part 2 - Cost of Accurate Computation in TM

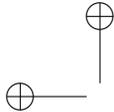
Type of Matrix	SVD	EVD	NF	$Ax=b$ F ⁷⁾	$Ax=b$ B ⁷⁾
Cauchy	n^3 *[12]	?	n^2 [9]	n^2 [9]	?
TP Cauchy	n^3 *[12]	n^3 *[33]	n^2 [9]	n^2 [9]	n^2 [9]
Vandermonde	n^3 *[12, 19]	?	n^2 [6]	No *[18]	?
TP Vandermonde	n^3 *[12, 19]	n^3 *[33]	n^2 [6]	n^2 [28]	n^2 [28]
Confluent Vandermonde	?	?	n^2 [28]	No *[18]	?
TP Confluent Vandermonde	n^3 *[33]	n^3 *[33]	n^2 [28]	n^2 [28]	n^2 [28]
Three-term Vandermonde Orth. Poly.	n^3 *[19]	?	?	?	?
3-term Vand. Orth. Poly. + other cond. ⁵⁾	n^3 *[19]	?	?	n^2 [28]	?
Generalized Vandermonde	No *[18]	?	No *[18]	No *[18]	?
TP generalized Vandermonde ⁴⁾	Λn^3 *[33]	Λn^3 *[33]	Λn^2 *[21, 22]	Λn^2 *[21, 22]	Λn^2 *[21, 22]
Any TP ⁶⁾	n^3 *[33]	n^3 *[33]	0	n^2 [28]	n^2 [28]
Wkly diag.dom. M-matrices	n^3 *[20]	?	?	?	?

- 1) Each entry is meant in a big-O sense, i.e. n^2 means $O(n^2)$.
- 2) Below each entry is a citation to where the result (old or new) appears.
- 3) Our results have an asterisk (*) in front of the citation.
- 4) If $G_{ij} = x_i^{j-1+\lambda_{n-j}}$ then $\Lambda \leq 2(\lambda_1 + 1)(\lambda_2 + 1)^2 \dots (\lambda_p + 1)^2 p$.
- 5) Also need $0 < x_1 < x_2 < \dots < x_n$, $d_i > 0$, $c_i \geq 0$ for all i , and $b_i - x_k \leq 0$ for all $i + k \leq n - 1$.
- 6) Any TP matrix, specified by the entries of its bidiagonal decomposition, or by combining such matrices by taking products, Schur complements and scaled inverses as described in Section 5.
- 7) “ $Ax = b$: F” means x satisfies a *Forward* error bound $|x - \hat{x}| \leq O(\epsilon)|A^{-1}||b|$ (so $|x - \hat{x}| \leq O(\epsilon)|x|$, when $\text{sign}(x_i) = (\pm 1)^i$ has alternating sign pattern and A is TP). “ $Ax = b$: B” means x has *Backward* error $|A - \hat{A}| \leq O(\epsilon)|A|$.

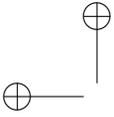


Bibliography

- [1] T. Ando. Totally positive matrices. *Lin. Alg. Appl.*, 90:165–219, 1987.
- [2] S. A. Attahiru, X. Junggong, and Q. Ye. Accurate computation of the smallest eigenvalue of a diagonally dominant M-matrix. preprint, 2002.
- [3] S. A. Attahiru, X. Junggong, and Q. Ye. Entrywise perturbation theory for diagonally dominant M-matrices with applications. *Numer. Math.*, 90:401–414, 2002.
- [4] J. Barlow and J. Demmel. Computing accurate eigensystems of scaled diagonally dominant matrices. *SIAM J. Num. Anal.*, 27(3):762–791, June 1990.
- [5] J. Barlow, B. Parlett, and K. Veselic, editors. *Proceedings of the International Workshop on Accurate Solution of Eigenvalue Problems*, volume 309 of *Linear Algebra and its Applications*. Elsevier, 2001.
- [6] Å. Björck and V. Pereyra. Solution of Vandermonde systems of equations. *Math. Comp.*, 24(112):893–903, 1970.
- [7] L. Blum, F. Cucker, M. Shub, and S. Smale. *Complexity and Real Computation*. Springer, 1997.
- [8] L. Blum, M. Shub, and S. Smale. On a theory of computation and complexity over the real numbers: NP-completeness, recursive functions and universal machines. *AMS Bulletin (New Series)*, 21(1):1–46, July 1989.
- [9] T. Boros, T. Kailath, and V. Olshevsky. A fast Björck-Pereyra-type algorithm for parallel solution of Cauchy linear equations. *Lin. Alg. Appl.*, 302/303:265–293, 1999.
- [10] K. Clarkson. Safe and effective determinant evaluation. In *33rd Annual Symp. on Foundations of Comp. Sci.*, pages 387–395, 1992.
- [11] F. Cucker and S. Smale. Complexity estimates depending on condition and roundoff error. *J. ACM*, 46(1):113–184, Jan 1999.
- [12] J. Demmel. Accurate SVDs of structured matrices. *SIAM J. Mat. Anal. Appl.*, 21(2):562–580, 1999.



- [13] J. Demmel. Accurate SVDs of structured matrices. *SIAM J. Mat. Anal. Appl.*, 21(2):562–580, 1999.
- [14] J. Demmel. The complexity of accurate floating point computation. In *Proceedings of the 2002 International Congress of Mathematicians*, Beijing, 2002.
- [15] J. Demmel, M. Gu, S. Eisenstat, I. Slapničar, K. Veselić, and Z. Drmač. Computing the singular value decomposition with high relative accuracy. *Lin. Alg. Appl.*, 299(1–3):21–80, 1999.
- [16] J. Demmel and Y. Hida. Accurate and efficient floating point summation. to appear in *SIAM J. Sci. Comp.*; www.cs.berkeley.edu/~demmel/AccurateSummation.pdf, 2002.
- [17] J. Demmel and W. Kahan. Accurate singular values of bidiagonal matrices. *SIAM J. Sci. Stat. Comput.*, 11(5):873–912, September 1990.
- [18] J. Demmel and P. Koev. Necessary and sufficient conditions for accurate and efficient rational function evaluation and factorizations of rational matrices. In V. Olshevsky, editor, *Special Issue on Structured Matrices in Mathematics, Computer Science and Engineering*, volume 281 of *Contemporary Mathematics*, pages 117–145. AMS, 2001.
- [19] J. Demmel and P. Koev. Accurate SVDs of polynomial Vandermonde matrices involving orthonormal polynomials. to appear in *Lin. Alg. Appl.*; www-math.mit.edu/~plamen/files/chebvand.ps, 2002.
- [20] J. Demmel and P. Koev. Accurate SVDs of weakly diagonally dominant M-matrices. submitted to *Num. Math.*; www-math.mit.edu/~plamen/files/mmat.ps, 2002.
- [21] J. Demmel and P. Koev. Efficient and accurate evaluation of Schur and Jack polynomials. in preparation, www-math.mit.edu/~plamen/files/p-schur.pdf, 2002.
- [22] J. Demmel and P. Koev. The accurate and efficient solution of a totally positive generalized Vandermonde linear system. submitted to *SIAM J. Mat. Anal. Appl.*, www-math.mit.edu/~plamen, 2003.
- [23] J. Demmel and K. Veselić. Jacobi’s method is more accurate than QR. *SIAM J. Mat. Anal. Appl.*, 13(4):1204–1246, 1992.
- [24] I. S. Dhillon. *A New $O(n^2)$ Algorithm for the Symmetric Tridiagonal Eigenvalue/Eigenvector Problem*. PhD thesis, University of California, Berkeley, California, May 1997.
- [25] F. M. Dopico, J. M. Molera, and J. Moro. An orthogonal high relative accuracy algorithm for the symmetric eigenproblem. to appear in *SIAM J. Mat. Anal. Appl.*, www.uc3m.es/uc3m/dpto/MATEM/molera/indice.html, 2003.



- [26] F. Gantmacher and M. Krein. *Oszillationsmatrizen, Oszillationskerne, und kleine Schwingungen mechanischer Systeme*. Akademie-Verlag, Berlin, 1960.
- [27] I. Gohberg, T. Kailath, and V. Olshevsky. Fast Gaussian elimination with partial pivoting for matrices with displacement structure. *Math. Comp.*, 64(212):1557–1576, 1995.
- [28] N. J. Higham. Stability analysis of algorithms for solving confluent Vandermonde-like systems. *SIAM J. Mat. Anal. Appl.*, 11(1):23–41, 1990.
- [29] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, PA, 1996.
- [30] S. Karlin. *Total Positivity*. Stanford University Press, 1968.
- [31] P. Koev. *Accurate and efficient computations with structured matrices*. PhD thesis, University of California, Berkeley, California, May 2002.
- [32] P. Koev. Accurate computations with totally nonnegative matrices. in preparation, math.mit.edu/~plamen, 2003.
- [33] P. Koev. Accurate eigenvalues and SVDs of totally positive matrices. in preparation, www-math.mit.edu/~plamen/files/tpeig.pdf, 2003.
- [34] H. Lewis and C. Papadimitriou. *Elements of the theory of computation*. Prentice Hall, 2nd edition edition, 1997.
- [35] I. G. MacDonald. *Symmetric functions and Hall polynomials*. Oxford University Press, 2nd edition, 1995.
- [36] J. J. Martínez and J. M. Peña. Fast algorithms for Björck-Pereyra type for solving Cauchy-Vandermonde linear systems. *Appl. Num. Math.*, 26(3):343–352, 1998.
- [37] C. O’Cinneide. Relative-error for the LU decomposition via the GTH algorithm. *Numer. Math.*, 73:507–519, 1996.
- [38] M. Pour-El and J. Richards. *Computability in Analysis and Physics*. Springer-Verlag, 1989.
- [39] D. Priest. Algorithms for arbitrary precision floating point arithmetic. In P. Kornerup and D. Matula, editors, *Proceedings of the 10th Symposium on Computer Arithmetic*, pages 132–145, Grenoble, France, June 26-28 1991. IEEE Computer Society Press.
- [40] Reliable Computing (a journal). Kluwer. www.cs.utep.edu/interval-comp/rcjournal.html.
- [41] S. Smale. Some remarks on the foundations of numerical analysis. *SIAM Review*, 32(2):211–220, June 1990.