



Technisch-Naturwissenschaftliche
Fakultät

Fast and Rigorous Computation of Special Functions to High Precision

DISSERTATION

zur Erlangung des akademischen Grades

Doktor

im Doktoratsstudium der

Technischen Wissenschaften

Eingereicht von:

Fredrik Johansson

Angefertigt am:

Research Institute for Symbolic Computation

Beurteilung:

Dr. Manuel Kauers (Betreuung)

Dr. Paul Zimmermann

Linz, Februar 2014

The research was funded by the Austrian Science Fund (FWF): grant no. Y464-N18.

Abstract

The problem of efficiently evaluating special functions to high precision has been considered by numerous authors. Important tools used for this purpose include algorithms for evaluation of linearly recurrent sequences, and algorithms for power series arithmetic.

In this work, we give new baby-step, giant-step algorithms for evaluation of linearly recurrent sequences involving an expensive parameter (such as a high-precision real number) and for computing compositional inverses of power series. Our algorithms do not have the best asymptotic complexity, but they are faster than previous algorithms in practice over a large input range.

Using a combination of techniques, we also obtain efficient new algorithms for numerically evaluating the gamma function $\Gamma(z)$ and the Hurwitz zeta function $\zeta(s, a)$, or Taylor series expansions of those functions, with rigorous error bounds. Our methods achieve softly optimal complexity when computing a large number of derivatives to proportionally high precision.

Finally, we show that isolated values of the integer partition function $p(n)$ can be computed rigorously with softly optimal complexity by means of the Hardy-Ramanujan-Rademacher formula and careful numerical evaluation.

We provide open source implementations which run significantly faster than previously published software. The implementations are used for record computations of the partition function, including the tabulation of several billion Ramanujan-type congruences, and of Taylor series associated with the Riemann zeta function.

Zusammenfassung

Die effiziente Auswertung spezieller Funktionen mit hoher Genauigkeit ist schon von vielen Autoren behandelt worden. Grundlegende Techniken in diesem Zusammenhang sind Algorithmen zur Berechnung linearer Rekurrenzen und Algorithmen für Potenzreihenarithmetik.

In dieser Arbeit geben wir neue Baby-step-giant-step-Algorithmen zur Berechnung linearer Rekurrenzen an, die einen teuren Parameter enthalten (zum Beispiel eine reelle Zahl mit hoher Genauigkeit), sowie Algorithmen zur Umkehrung von Potenzreihen. Unsere Algorithmen haben zwar nicht die beste asymptotische Komplexität, doch sind sie in der Praxis für einen großen Eingabebereich schneller als die bisherigen Algorithmen.

Durch eine Kombination verschiedener Techniken erhalten wir außerdem effiziente Algorithmen zur numerischen Auswertung der Gammafunktion $\Gamma(z)$ und der Hurwitz-Zeta-Funktion $\zeta(s, a)$, bzw. der Taylorentwicklungen dieser Funktionen, zusammen mit garantierten Fehlerschranken. Unsere Methoden haben quasi-optimale Komplexität, wenn man eine große Zahl von Ableitungen zu entsprechend hoher Genauigkeit berechnet.

Schließlich zeigen wir, dass isolierte Werte der Abzählfunktion $p(n)$ von Integer-Partitionen exakt und in quasi-optimaler Komplexität mit der Formel von Hardy-Ramanujan-Rademacher und sorgfältiger Numerik berechnet werden können.

Wir stellen Open-Source-Implementierungen bereit, die signifikant schneller sind als bisher publizierte Software. Diese Implementierungen werden zu Rekord-Berechnungen der Partitionsfunktion verwendet, inklusive einer Auflistung von mehreren Milliarden Ramanujan-Kongruenzen, und für Taylorentwicklungen im Zusammenhang mit der Riemannschen Zeta-Funktion.

Acknowledgements

My deepest gratitude goes to my exceptionally supportive and seemingly inexhaustible advisor, Manuel Kauers, for sharing his knowledge and grant money, assisting with practical matters, being a general source of inspiration, and giving me (for better or worse) much freedom to pursue my whims. I also emphatically thank professor Peter Paule, who in many ways has been responsible for the terrific research environment at RISC. Among my numerous other friendly colleagues at RISC, Silviu Radu deserves a special mention for the help he provided in my investigations of the partition function.

An additional 2^{10} thanks go to William Hart, Sebastian Pancratz and David Harvey, from whom I learned a great deal about implementing fast algorithms. Without their work, this thesis would necessarily have been concerned with other, less interesting topics. It has been a particular privilege to collaborate extensively with William, who (besides directly helping me improve my papers and code) has given me much encouragement and provided countless interesting discussions, be they on matters mathematical, computational, or completely off-topic.

Several other people in the computer algebra and computational number theory communities provided feedback which influenced this work. If they read this, they hopefully know who they are.

Most importantly, as always: with unbounded love to all my family, and to the memory of those who are no longer here.

Contents

1	Introduction	1
2	Arithmetic	5
2.1	Complexity of algorithms	5
2.2	Fast multiplication	6
2.3	Computing with real numbers	7
2.4	Computing with polynomials	9
3	Evaluating linearly recurrent sequences	11
3.1	Holonomic sequences and functions	11
3.2	Recurrences in matrix form	12
3.3	Binary splitting	14
3.4	Fast multipoint evaluation	16
3.5	Rectangular splitting	16
3.5.1	Variations	19
3.5.2	Several parameters	21
3.5.3	Numerical evaluation	22
3.5.4	Summation of power series	23
3.6	Computing the gamma function	24
3.6.1	Rising factorials	25
3.6.2	A one-parameter hypergeometric series	27
3.6.3	Remarks	29

4	Evaluating functions of power series	31
4.1	Bernoulli numbers	31
4.2	Defining functions of power series	32
4.3	Fast composition	33
4.4	Elementary functions	34
4.5	Fast reversion without Newton iteration	35
4.5.1	The algorithm	36
4.5.2	Complexity analysis	38
4.5.3	Benchmarks	43
4.5.4	Reversion with numerical coefficients	45
4.5.5	Remarks	46
4.6	The gamma function	47
4.7	Integer zeta values	49
4.8	The Hurwitz zeta function	53
4.8.1	Evaluation using the Euler-Maclaurin formula	54
4.8.2	Algorithmic matters	57
4.8.3	Implementation and benchmarks	61
4.8.4	Remarks	68
5	Computing the partition function	71
5.1	The pentagonal number theorem	71
5.2	The Hardy-Ramanujan-Rademacher formula	72
5.3	Evaluating the exponential sums	74
5.3.1	A simple algorithm	75
5.3.2	A fast algorithm	76
5.4	Numerical evaluation	79
5.5	Rapid computation of roots of unity	82
5.5.1	Choice of annihilating polynomial	83

5.5.2	Generating trigonometric minimal polynomials	83
5.5.3	Performance evaluation	84
5.6	Implementation	86
5.7	Multi-evaluation and congruence generation	88
5.7.1	Weaver's algorithm	89
5.7.2	Comparison of algorithms for multi-evaluation	90
5.7.3	Results	91
5.8	Remarks	92
A Implementations		93
Bibliography		97

Chapter 1

Introduction

One of the central tasks in computer algebra is the development of algorithms for computing explicit values of mathematical functions. Humans – and increasingly also computers – can look at computed data to search for patterns, leading to new conjectures and ultimately new theorems [5, 19]. Explicit values are sometimes also required as parts of proofs or in the execution of symbolic algorithms. It is important that the algorithms and implementations we use are efficient (so that we can afford to compute the data we desire) and reliable (so that we can be reasonably certain that the output is actually correct).

Functions of particular interest are called *special functions* – by tradition, these are the functions which appear frequently enough in applications to have been given names (the gamma function, the Riemann zeta function, Bessel functions, and so on). A more refined viewpoint is to consider special functions as solutions of differential, difference or functional equations of certain types. One can then hope to develop theoretical or computational tools that work for a whole class of functions, and gradually replace *ad-hoc* methods with systematic procedures that can be automated by computers.

The objective of this thesis is to study algorithms that scale well for computation of function values which are very large, or which need to be approximated to very high precision. Scalable algorithms are crucial for attacking research problems in number theory and combinatorics, but efficiency is also a concern for general use within mathematical software.

This thesis provides a modest contribution to the already extensive algorithmic theory concerning computation of special functions. We also attempt to reduce the gap between theory and practice by implementing algorithms carefully and testing their performance in the real world. We particularly hope to emphasize the idea that *numerical* computation can be done with a high level of reliability without sacrificing performance.

In chapter 2, we recall some fundamental algorithmic concepts, including fast multiplication of numbers, polynomials and matrices. This chapter also introduces notation used throughout the work.

Chapter 3 treats the problem of quickly evaluating the n -th term in a sequence defined by a linear recurrence equation. The terms can, for instance, be integers, real numbers, polynomials, or power series. A particularly important case is the class of sequences which are *holonomic* (or *P-finite*), meaning that the terms satisfy a linear recurrence equation with polynomial coefficients. A large portion of the special functions arising in applications are expressible (exactly or approximately) in terms of holonomic sequences, and algorithms that solve this problem efficiently in the general case thus lead to efficient algorithms for specific functions.

Well-known fast algorithms for evaluation of holonomic sequences include *binary splitting* and *fast multipoint evaluation*. In section 3.5, we give an algorithm which becomes efficient when the recurrence equation involves an “expensive” parameter (in a sense which is made precise), based on the baby-step giant-step technique of Paterson and Stockmeyer [87] (called *rectangular splitting* in [27]).

Our algorithm can be viewed as a generalization of the method given by Smith in [96] for computing rising factorials. Conceptually, it also generalizes an algorithm given by Smith in [95] for evaluation of hypergeometric series. Our contribution is to recognize that rectangular splitting can be applied systematically to a very general class of sequences, and in an efficient way (we provide a detailed cost analysis, noting that some care is required in the construction of the algorithm to get optimal performance).

The main intended application of rectangular splitting is high-precision numerical evaluation of special functions, where the parameter is a real or complex number (represented by a floating-point approximation). In section 3.6, we present implementation results comparing several different algorithms for numerical evaluation of the gamma function to very high precision.

Chapter 4 treats the problem of multi-evaluation of sequences, i.e. simultaneously computing all the values

$$f(0), f(1), f(2), \dots, f(n-1), f(n)$$

for a given n . This is trivial to do efficiently for linearly recurrent sequences. For more general sequences, a well-known paradigm is to identify the terms with the (formal) derivatives of a *generating function*, and evaluating an expression for the function using *power series arithmetic*. Power series manipulations also allow doing various computations with functions as analytic objects.

We first discuss standard algorithms for manipulation of power series, including formal composition and evaluation of elementary functions. In section 4.5, we then give a new

algorithm for *reversion of power series*, or equivalently, computing the coefficients in the series expansion of the compositional inverse of a given generating function. Two fast algorithms for this problem were given by Brent and Kung in 1978 [25]. Our algorithm is a baby-step giant-step version of the Lagrange inversion formula, analogous in spirit to the first Brent-Kung algorithm which is a baby-step giant-step version of Horner’s scheme. However, our algorithm is superior in two ways: it computes the reversion directly without using Newton iteration, and it allows taking advantage of structured matrix multiplication. Both differences are shown to give constant-factor speedups. Benchmarks indicate that our algorithm, although asymptotically slower than the second Brent-Kung algorithm, is the fastest algorithm in practice for reversion of power series over a large range of input sizes.

Next, we study the problem of efficiently computing series expansions of some higher transcendental functions to arbitrary precision. We study, in particular, the gamma function $\Gamma(z)$ in section 4.6, and the Hurwitz zeta function $\zeta(s, a)$ in section 4.8. The Hurwitz zeta function is a key function to consider for numerical evaluation, as many other transcendental functions and mathematical constants can be expressed in terms of it. We state and implement – to our knowledge for the first time – a complete algorithm to evaluate $\zeta(s, a)$ with rigorous error bounds for any values of s and a and for derivatives of arbitrary order. We make several observations regarding efficiency, and note that fast polynomial arithmetic can be exploited to achieve softly optimal complexity when computing a large number of derivatives.

Our implementation allows us to study the asymptotics of Taylor series coefficients associated with the Hurwitz zeta function. In particular, we provide explicit computations which show support for the conjectured asymptotic growth of the Keiper-Li coefficients (constituting empirical evidence for the Riemann hypothesis, although it should be stressed that this evidence is weaker than already published data based on explicit evaluation of zeros along the critical line). We also obtain asymptotic data about the Stieltjes constants. Compared to previous works, our contribution is twofold: our implementation allows us to reach coefficients of higher index (roughly 10^5 on present hardware), and our numerical values come with rigorous error bounds.

Finally, chapter 5 treats evaluation of the integer partition function $p(n)$. Fast simultaneous evaluation of the consecutive values $p(0), \dots, p(n)$ can easily be accomplished using power series methods, but computing isolated values more quickly is challenging. The values of the partition function form a non-holonomic sequence, and cannot (as far as anyone currently knows) be approached using techniques such as those discussed in chapter 3. Instead, the best known method is to compute a sufficiently accurate numerical approximation of a certain infinite series involving irrational numbers, the Hardy-Ramanujan-Rademacher (HRR) formula.

Here our contribution is to give the first comprehensive algorithmic analysis of the HRR formula. We prove that isolated values of the partition function can be computed with complexity that is nearly the best possible. This is a nontrivial result. We also discuss implementation aspects and present an implementation which runs with nearly-optimal complexity in practice, gaining more than a 100-fold speedup over previously published implementations. We are able to exactly determine values as large as $p(10^{19})$ – an integer with 3.5 billion digits. Extending previous work done by Weaver [109], we use the implementation to find several billion new Ramanujan-type congruences for the partition function.

This thesis is substantially based on a series of standalone papers penned by the author. Chapter 3 is an extended version of [54]. Chapter 4 is largely based on the papers [56] and [55]. The discussion about reversion includes an improved presentation of matrix multiplication algorithms and an added discussion about power series with numerical coefficients. The section about the power series of the gamma function is also new. Chapter 5 is based on [52], with a completely rewritten discussion regarding numerical evaluation, along with new implementation remarks. In general, the introductory content from the aforementioned papers has been reworked and expanded to unify the presentation.

Most of the algorithms covered in this thesis have been implemented using the Fast Library for Number Theory (FLINT) [47], which is joint work with William Hart and Sebastian Pancratz (as well as numerous contributors). To support asymptotically fast and provably correct computation with real and complex numbers, the author has developed the Arb library [53], based on the concept of *ball arithmetic*. We give a brief survey of the FLINT and Arb libraries in Appendix A.

Chapter 2

Arithmetic

In this chapter, we review standard techniques for computing with numbers and polynomials. We only briefly discuss concepts and notations required for the later chapters, without going into too many details. The reader can find more background material in the books by Knuth [64], von zur Gathen and Gerhard [107], and Brent and Zimmermann [27].

2.1 Complexity of algorithms

When analyzing the efficiency of an algorithm, we usually count the number of bit or word operations executed on a serial machine. This measure, which we synonymously call bit complexity or time complexity, roughly models the actual running time of an implementation on a physical computer. It can, however, be inaccurate since it discounts aspects such as parallelism and memory access costs.

In many cases, it is more convenient to count the number of ring operations ($+$, $-$, \times , equality test) in some (computable) ring R . This measure is only loosely suggestive of practical efficiency, since ring elements may take a variable amount of space to represent and require a variable amount of bit operations to operate on. In general, we can convert a bound expressed in terms of ring operations to a bit complexity bound if we are able to bound the size of the occurring elements and have a bit complexity bound for ring operations done on elements of given size.

As usual, $O(f(n))$ denotes the class of functions bounded by constant multiples of $f(n)$ for all sufficiently large n and $O^\sim(f(n))$ denotes the class of functions asymptotically bounded by $f(n)$ times polynomials in $\log(f(n))$ (see section 25.7 in [107] for precise definitions). If $g(n) = O^\sim(f(n))$, we may also write $g(n) = O(f(n)^{1+o(1)})$ or $g(n) = O(f(n)^{1+\varepsilon})$. If $f(n) = O^\sim(n)$, we call $f(n)$ softly linear or quasilinear. If $f(n)$ is a lower

bound for the complexity of solving a given problem of size n , we say that an algorithm with complexity $O^\sim(f(n))$ is softly optimal or quasioptimal. We also occasionally use the notations $\Omega(f(n))$ (for lower bounds) and $\Theta(f(n))$ (for two-sided bounds).

2.2 Fast multiplication

A useful technique when constructing algorithms is to reduce the solution of a problem to a combination of well-understood “primitive” operations. This helps theoretical analysis (we can rely on nontrivial proven complexity bounds for the primitive operations), and in practical implementations (we directly benefit from the availability of highly optimized library routines for the primitive operations). In computer algebra, the most important primitive operation is multiplication in various base rings. Bernstein [9] provides a good survey of algorithms for fast multiplication, and techniques for reducing problems to multiplication.

We frequently use the following multiplication complexity bounds. In each case, when we refer to *the* complexity of multiplication, we actually refer to the complexity with an arbitrary but fixed multiplication algorithm (which need not be the fastest possible).

- $M_{R[x]}(n)$ denotes a bound for the number of ring operations required to multiply two polynomials of degree at most n with coefficients in a ring R . The classical multiplication algorithm gives $M_{R[x]}(n) = O(n^2)$ and the Karatsuba algorithm gives $M_{R[x]}(n) = O(n^{\log_2 3}) = O(n^{1.585})$. If R contains sufficiently many roots of unity, the fast Fourier transform (FFT) allows multiplying polynomials using $M_{R[x]}(n) = O(n \log n)$ operations, which is softly optimal. When R does not have sufficiently many roots of unity, softly optimal multiplication is still possible by moving to an extension ring containing “artificial” roots of unity and then multiplying using the FFT. For an arbitrary ring R , we have $M_{R[x]}(n) = O(n \log n \log \log n)$ by the Cantor-Kaltofen theorem [30, 7].

FFT multiplication of polynomials is used in practice and exhibits the expected quasilinear complexity for many common rings R . The range of n where it becomes worthwhile is highly sensitive to implementation details.

- $M_{\mathbb{Z}}(n)$ denotes a bound for the number of bit (or word) operations required to multiply two n -bit integers. The algorithms and complexity bounds for integers are quite similar to those for polynomials (in fact, it is common in implementations to express one operation in terms of the other). In particular, we have $M_{\mathbb{Z}}(n) = O(n \log n \log \log n)$ by the Schönhage-Strassen algorithm [9]. The best current bound for integer multiplication is Fürer’s $M_{\mathbb{Z}}(n) = n \log n 2^{O(\log^* n)}$, where $\log^* n$ denotes the iterated logarithm.

Modern implementations of arbitrary-precision integer arithmetic such as GMP and MPIR [41, 79] switch between several multiplication algorithms: typically classical multiplications is used below about 10^3 bits, Karatsuba multiplication (or the higher-order Toom-Cook variants) is used for larger integers up to about 10^5 or 10^6 bits, and an FFT algorithm such as the Schönhage-Strassen algorithm is used for even larger integers. The asymptotic quasilinearity of FFT multiplication is very clearly realized in practice for numbers with millions or billions of digits.

- $\text{MM}_R(n)$ denotes a bound for the number of ring operations required to multiply two $n \times n$ matrices with coefficients in a commutative ring R . We have $\text{MM}_R(n) = O(n^3)$ using classical multiplication and $O(n^{2.807})$ using the Strassen algorithm. The best known bound is $O(n^{2.3727})$, due to Stothers [98] and Vassilevska Williams [104]. On present-day hardware, only the classical algorithm and Strassen's algorithm are considered practical. If $\text{MM}_R(n) = O(n^\omega)$ for some $2 < \omega \leq 3$, we call ω an *exponent of matrix multiplication*.

The ring is usually clear from the context, in which case we simply write $\text{M}(n)$ or $\text{MM}(n)$. We implicitly assume that the complexity bounds satisfy natural regularity conditions such as $\text{M}(n+m) \geq \text{M}(n) + \text{M}(m)$.

We sometimes consider *unbalanced* multiplications. If $\text{M}(n, m)$ denotes the complexity of multiplying two polynomials of respective degrees n and m where $n \geq m$, we can assume that $\text{M}(n, m) \leq \text{M}(n, n) = \text{M}(n)$. When $n \gg m$, we can do better by breaking the single multiplication into several balanced multiplications, giving $\text{M}(n, m) = \lceil n/m \rceil \text{M}(m) + O(n)$. The analogous statements hold for the bit complexity of unbalanced integer multiplication.

Likewise, unbalanced matrices can be multiplied by breaking them into smaller square blocks. Remarkably, Huang and Pan have shown [51] that this strategy is suboptimal for highly unbalanced matrices with the best presently known algorithms. Letting $\text{MM}(x, y, z)$ denote the complexity of multiplying a matrix of size $x \times y$ by a matrix of size $y \times z$, and taking $m = n^{1/2}$, Huang and Pan obtain $\text{MM}(m, m, n) = O(n^{1.667})$ whereas repeated multiplication of square matrices costs $m\text{MM}(m, m, m) = O(n^{1.687})$ with the Stothers-Vassilevska Williams algorithm.

2.3 Computing with real numbers

In general, an algorithm purported to work with elements of a ring R only makes sense if R is *computable* (or *effective*), meaning that any element of R can be represented using a finite number of bits and that arithmetic operations can be carried out using a finite number of steps. We typically also require that determining whether a given element is

zero can be done using a finite number of steps. Examples of computable rings include \mathbb{Z} , $\mathbb{Q}(\sqrt{-1})$ and $(\mathbb{Z}/3\mathbb{Z})[x, y]$.

The rings \mathbb{R} and \mathbb{C} are not computable in this sense, and yet we are often forced to work with them in applications. Fortunately, for many purposes, it is sufficient to replace real numbers by rational approximations. We can then carry out arithmetic operations on the approximations, and determine bounds for any propagated errors.

We can, in particular, represent a real number x by a ball $X = \hat{x} \pm r := [\hat{x} - r, \hat{x} + r]$ such that $x \in X$, where $\hat{x}, r \in \mathbb{Q}$ and $r \geq 0$. If $y = f(x)$ where $f : \mathbb{R} \rightarrow \mathbb{R}$ is some function, we may attempt to approximate y by a ball $Y = \hat{y} \pm s$ such that $f(X) \subseteq Y$. For sufficiently “nice” functions (which we need not precisely characterize here), we can always compute an output ball such that $s \rightarrow 0$ when $r \rightarrow 0$. In general, functions with discontinuities such as $f(x) = \lfloor x \rfloor$ are not “nice”. In particular, we cannot tell if a real number represented by a ball is zero (without additional information). We can, however, prove that a number resulting from the composition of “nice” functions is *not* zero, by computing a sufficiently precise approximation.

For efficiency reasons, it is preferable to choose *floating-point numbers* as rational approximations. A binary floating-point number is a rational number $a \times 2^b$ where $a \in \mathbb{Z}$ is called the *mantissa* (or *significand*) and $b \in \mathbb{Z}$ is called the *exponent*. It should be noted that some authors define the parts of a floating-point number slightly differently, but the distinction does not matter for our purposes. A nonzero floating-point number can be put in canonical form by choosing a such that $2 \nmid a$ (zero can be represented canonically by $a = b = 0$). It is convenient to adjoin the special values $-\infty$, $+\infty$ and NaN (not-a-number) to the set of floating-point numbers to allow representing limits and to allow any floating-point operation to have a well-defined result for any input.

To prevent coefficient explosion, it is crucial to *round* floating-point numbers. Rounding the mantissa of a floating-point number x to precision p bits, i.e. choosing a such that $|a| < 2^p$, generally introduces a *rounding error* of order $|x|2^{-p}$. When rounding the midpoint \hat{x} of a ball $\hat{x} \pm r$, we must take care to add an upper bound for the rounding error to the radius r (this addition can also be done as a floating-point operation, rounding upwards if necessary).

For a floating-point approximation $\hat{x} \approx x$, we define the *absolute error* as $|x - \hat{x}|$ and the *relative error* as $|x - \hat{x}|/|x|$, or when measuring in bits, the respective \log_2 values. We also define *accuracy* as the inverse of error (or the negation of the error when measuring in bits). We informally refer to the radius r of any ball $\hat{x} \pm r$ as its error, although strictly speaking this is just an upper bound for the true error of the midpoint.

In most situations, numerically evaluating a fixed expression using a working precision of p bits results in an absolute or relative accuracy of $p + O(1)$ bits. If numerical

evaluation with a working precision of p bits costs $C(p)$, we can therefore usually compute the result to an accuracy of n bits at a cost of $O(C(n))$. A common trick when the $O(1)$ term is unknown is to first make a reasonable guess, and then repeatedly double the precision until convergence. However, we need to be careful when the expression we are evaluating varies along with the precision. See in particular the remarks about polynomials below.

Addition, multiplication, division and square root of p -bit floating-point numbers have the same complexity as the corresponding integer operations, up to a constant factor (we ignore the cost of manipulating the exponent, which for most practical purposes can be assumed to lie in a fixed range).

The elementary transcendental functions (exp, log, sin, atan, etc.) can be computed to high precision in nearly linear time. The fastest known algorithm is based on the arithmetic-geometric mean (AGM), and allows approximating $\log x$ (where the number x is held fixed) to p -bit accuracy in $O(M(p) \log p)$ bit operations. One can obtain all other elementary functions with at most a constant factor slowdown using complex arithmetic and Newton iteration.

An alternative way to evaluate elementary functions is to use Taylor series. Combined with use of functional equations, one can achieve a complexity of $O(M(p) \log^2 p)$ for elementary functions, and in practice these algorithms are often faster than AGM-based methods. Moreover, such techniques also generalize to evaluation of many higher transcendental functions. We discuss techniques applicable to fast evaluation of Taylor series further in chapter 3.

The principle of attaching error bounds to floating-point numbers can be extended to numerical computations involving complex numbers, polynomials, matrices, and other objects. This technique is known as *ball arithmetic* [103].

2.4 Computing with polynomials

Arithmetic operations on polynomials can generally be reduced to polynomial multiplication and performed with complexity that is softly optimal in the size of the problem.

Polynomials in $\mathbb{Z}[x]$ of degree n and coefficients at most p bits can be multiplied using $O^\sim(np)$ bit operations (Corollary 8.27 in [107]). In practice, we often have $p = O^\sim(n)$, making the effective complexity of working with degree- n integer polynomials $O^\sim(n^2)$. Arithmetic in $\mathbb{Q}[x]$ can be reduced to arithmetic in $\mathbb{Z}[x]$ by clearing denominators.

Fast multiplication of polynomials in $\mathbb{R}[x]$ and $\mathbb{C}[x]$ represented using floating-point (or ball) coefficients can either be done directly using FFT or by truncating to integer poly-

nomials on a common power-of-two exponent and multiplying exactly in $\mathbb{Z}[x]$. The latter algorithm has the advantage that rounding error is introduced only when converting to and back from an integer polynomial, and not for each arithmetic operation in the multiplication (however, a separate computation is required to bound the propagated errors from the input polynomials). It also has the advantage that we avoid the overhead of doing a large number of arithmetic operations on individual numerical coefficients.

With numerical coefficients, the accuracy of the output polynomial can be much lower than the working precision, depending on the shape of the input polynomials and the details of error propagation in the multiplication algorithm. If all coefficients have roughly the same magnitude, then the loss of accuracy when multiplying degree- n polynomials is generically of order $\log n$ bits. In practice, it is often the case that the coefficients vary as $\log |c_k| \sim \pm k \log^{O(1)} k$. Therefore, we generically need $O^\sim(n)$ bits of precision to get an accurate result, making the effective complexity for arithmetic with degree- n polynomials $O^\sim(n^2)$. This problem can sometimes be mitigated by rescaling polynomials and splitting them into smaller blocks, as discussed by van der Hoeven [102].

Chapter 3

Evaluating linearly recurrent sequences

A sequence (usually denoted by $(c(i))_{i=0}^{\infty}$, or just $c(i)$ or c_i) is a function from the natural numbers $\mathbb{Z}_{\geq 0}$ to some set V . In this chapter, we study ways to compute a single value $c(n)$ quickly when n is large, provided that the sequence $c(i)$ satisfies a linear recurrence equation of suitable type.

The terms in a sequence are sometimes objects of inherent interest, for instance due to having combinatorial significance. Sometimes the terms in a sequence are used to approximate a limiting value. For example, the sequence of rational numbers $s_0 = 1$, $s_i = s_{i-1} + 1/i!$ converges rapidly to $e \approx 2.71828$, and one way to approximate e to high precision is to compute s_n for some large n . This can be done more efficiently than by explicitly computing $s_0, s_1, \dots, s_{n-1}, s_n$ one by one.

3.1 Holonomic sequences and functions

A sequence $(c(i))_{i=0}^{\infty}$ is called *holonomic* (or *P-finite*) of order r if it satisfies a linear recurrence equation

$$a_r(i)c(i+r) + a_{r-1}(i)c(i+r-1) + \dots + a_0(i)c(i) = 0 \quad (3.1.1)$$

where a_0, \dots, a_r are polynomials. The class of holonomic sequences enjoys many useful closure properties: for example, holonomic sequences form a ring, and if $c(i)$ is holonomic then so is the sequence of partial sums $s(n) = \sum_{i=0}^n c(i)$. A sequence is called *hypergeometric* if it is holonomic of order $r = 1$. The sequence of partial sums of a hypergeometric sequence is holonomic of order (at most) $r = 2$.

Likewise, an analytic function or formal power series $f(x)$ is called holonomic (or *D-*

finite) of order r if it satisfies a linear differential equation with polynomial coefficients

$$a_r(x)f^{(r)}(x) + a_{r-1}(x)f^{(r-1)}(x) + \dots + a_0(x)f(x) = 0 \quad (3.1.2)$$

This is equivalent to the condition that the coefficients in the Taylor series of $f(x)$ form a holonomic sequence (the orders of the respective differential and difference equations may generally be different). Holonomic functions form a ring (equivalently, the Cauchy product $d(n) = \sum_{i=0}^n c_1(i)c_2(n-i)$ of two holonomic sequences is holonomic) and yield new holonomic functions upon differentiation, integration, or right-composition by algebraic functions.

Many integer and polynomial sequences of interest in number theory and combinatorics are holonomic, and the power series expansions of many well-known special functions (such as the exponential function, sine, cosine, logarithm, arctangent, error function, and the Bessel functions) are holonomic.

Example 3.1.1. The sequence $t(k) = 1/k!$ is holonomic of order $r = 1$, satisfying the recurrence equation $(k+1)t(k+1) - t(k) = 0$ together with the initial value $t(0) = 1$. The sequence $u(k) = x^k t(k)$ is also holonomic of order $r = 1$, satisfying $(k+1)u(k+1) - xu(k) = 0$ with initial value $u(0) = 1$.

The sequence $s(k) = \sum_{i=0}^k u(i)$ is holonomic of order $r = 2$, satisfying the recurrence equation

$$(k+2)s(k+2) - (x+k+2)s(k+1) + xs(k) = 0$$

together with initial values $s(0) = 1$, $s(1) = 1 + x$. The exponential function can be computed as $\exp(x) = \lim_{n \rightarrow \infty} s(n)$.

A particularly nice property of the class of holonomic sequences and functions is that the closure properties can be carried out algorithmically. If we have algorithms for efficiently evaluating holonomic sequences of general form, it is possible to mechanically go from a specific symbolically defined holonomic function to a recurrence equation, and then to an efficient evaluation scheme for that function.

Some of the algorithms we describe below apply to more general linearly recurrent sequences, but we usually restrict the discussion to holonomic sequences for simplicity.

3.2 Recurrences in matrix form

Instead of working directly with order- r scalar recurrences such as (3.1.1), it is convenient to work with vector recurrences. Let R be a commutative ring with unity and of sufficiently large characteristic where necessary. Consider a sequence of length- r vectors

$(c(i) = (c_1(i), \dots, c_r(i))^T)_{i=0}^\infty$ satisfying a recurrence equation of the form

$$\begin{pmatrix} c_1(i+1) \\ \vdots \\ c_r(i+1) \end{pmatrix} = M(i) \begin{pmatrix} c_1(i) \\ \vdots \\ c_r(i) \end{pmatrix} \quad (3.2.1)$$

where $M \in R[k]^{r \times r}$ (or $\text{Quot}(R)(k)^{r \times r}$) and where $M(i)$ denotes entrywise evaluation. Given an initial vector $c(0)$, our goal is to evaluate the single vector $c(n)$ for some $n > 0$, where we assume that no denominator in M vanishes for $0 \leq i < n$.

A scalar recurrence of the form (3.1.1) can be rewritten as (3.2.1) by taking the vector to be $\tilde{c}(i) = (c(i), \dots, c(i+r-1))^T$ and setting M to the companion matrix

$$M = \frac{1}{a_r} \begin{pmatrix} & & & a_r \\ & & \ddots & \\ & & & a_r \\ -a_0 & -a_1 & \dots & -a_{r-1} \end{pmatrix}. \quad (3.2.2)$$

In either case, we call the sequence holonomic (of order r).

Through repeated application of the recurrence equation, $c(n)$ can be evaluated using $O(r^2n)$ arithmetic operations (or $O(rn)$ if M is companion) and temporary storage of $O(r)$ values. We call this strategy the *naive algorithm*.

The naive algorithm is not generally optimal for large n . The idea behind faster algorithms is to first compute the matrix product

$$P(0, n) = \prod_{i=0}^{n-1} M(i). \quad (3.2.3)$$

and then multiply it by the vector of initial values (matrix multiplication is of course noncommutative, and throughout this chapter the notation in (3.2.3) is understood to mean $M(n-1) \dots M(2)M(1)M(0)$).

Using a matrix product increases the cost to $O(r^\omega n)$ arithmetic operations where ω is the exponent of matrix multiplication, but we can save time for large n by exploiting the structure in the matrix product (we assume in the remainder of this chapter that r is fixed, and omit $O(r^\omega)$ factors from any complexity estimates).

The improvement is most dramatic when all matrix entries are constant, e.g. in the case of Fibonacci numbers

$$\begin{pmatrix} F_{i+1} \\ F_{i+2} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} F_i \\ F_{i+1} \end{pmatrix},$$

allowing binary exponentiation (requiring $O(\log n)$ arithmetic operations) or diagonalization to be used.

The case of constant coefficients is rather special, however. Techniques that work with nonconstant coefficients include *binary splitting*, which gives a near-optimal asymptotic speedup for certain types of coefficients, and *fast multipoint evaluation* which is the asymptotically fastest known algorithm for holonomic sequences with general coefficients. We discuss these algorithms in the following sections, and finally discuss new algorithms based on the idea of *rectangular splitting*.

From this point, we may view the problem of evaluating a term in a holonomic sequence as that of evaluating a matrix product (3.2.3) for some $M \in R[k]^{r \times r}$. It is not a restriction to demand that the entries of M are polynomials: if $M = \tilde{M}/q$, we can write $P(0, n) = \tilde{P}(0, n)/Q(0, n)$ where

$$\tilde{P}(0, n) = \prod_{i=0}^{n-1} q(i)\tilde{M}(i)$$

and

$$Q(0, n) = \prod_{i=0}^{n-1} q(i).$$

This reduces the problem to evaluating two denominator-free products, where the second product has order 1.

It is important to stress the usefulness of the matrix approach to working with sequences, as it often allows replacing clever formula manipulations with trivial matrix arithmetic. This viewpoint is enunciated by Bernstein [9] who provides an extensive bibliography and lists numerous applications of the idea. We can either implement computations directly in terms of matrices, or view a matrix algorithm as a “meta-algorithm” from which we can construct and analyze algorithms for specific sequences.

3.3 Binary splitting

In the *binary splitting* algorithm, we recursively compute a product of square matrices $P(a, b) = \prod_{i=a}^{b-1} M(i)$ as $P(m, b)P(a, m)$ where $m = \lfloor (a + b)/2 \rfloor$. Here the entries of $M(i)$ need not necessarily be polynomials of i as in the holonomic case, but can be arbitrary functions of i . Binary splitting becomes useful when the entries of the individual factors $M(i)$ are small, and the entries of partial products grow with the number of factors. Then this scheme balances the sizes of the subproducts in a way that allows us to take advantage of fast multiplication.

For example, take $M(i) \in R[x]^{r \times r}$ where the polynomials in $M(i)$ have bounded degree (note that the variable x used here has a different role than the variable k used in the previous section). Then $P(a, b)$ has entries in $R[x]$ of degree $O(b-a)$, and binary splitting can be shown to compute $P(0, n)$ using $O(M(n) \log n)$ operations in R , using $O(n)$ extra

storage. If the matrix entries are not constant polynomials, the output has degree $\Theta(n)$, making the binary splitting softly optimal. This is a significant improvement over the naive algorithm, which in general uses $O(n^2)$ coefficient operations to generate the n -th entry in a holonomic sequence of polynomials.

Analogously, binary splitting reduces the bit complexity from $O^\sim(n^2)$ to $O^\sim(n)$ for computing matrix products over \mathbb{Z} or \mathbb{Q} (or any algebraic number field) whenever the bit size of $M(i)$ for $0 \leq i < n$ grows at most as $\log^{O(1)} n$. In particular, this result applies to the evaluation of holonomic sequences, in which $M(i)$ has entries of size $O(\log n)$.

For example, we can use binary splitting to evaluate the sequence of partial sums $s(n)$ of the Taylor series of $\exp(x)$ (per Example 3.1.1). If x is a rational number, binary splitting allows us to compute $\exp(x)$ to a precision of p bits using $O^\sim(p)$ bit operations. This observation was made in 1976 by Brent [21].

If x is a real number with a precision of p bits, binary splitting does not provide a speedup over the $O^\sim(p^2)$ naive algorithm if applied directly to the same sum, but as Brent noted in [21], we can compute $\exp(x)$ in softly optimal time $O^\sim(p)$ as

$$\exp(x) = \exp(x_1) \exp(x_2) \cdots \exp(x_m)$$

where

$$x = x_1 + x_2 + \dots + x_m$$

and x_i extracts 2^i bits from the binary expansion of x (this balances the number of terms required in the Taylor series against the bit size of the terms). A variation of this idea works for computing other elementary functions with softly optimal complexity.

In fact, for any function $f(x)$ satisfying a holonomic differential equation with rational (or algebraic) coefficients, we can compute $f(x)$ in softly optimal time $O^\sim(p)$ for a fixed real or complex number x , by repeatedly translating the differential equation to points x_1, \dots, x_m as in Brent's algorithm for the exponential. This method was developed by Chudnovsky and Chudnovsky [32], who named it the *bit-burst algorithm*, and independently with improvements by van der Hoeven [100].

The bit-burst algorithm is used for evaluation of elementary functions in several libraries, and the general version for holonomic functions has been analyzed and implemented by Mezzarobba [75, 76, 77].

The binary splitting algorithm has been rediscovered or reapplied many times in different forms. Further discussion of the method and extensive bibliographic information can be found in [13], [44], [9] and [27].

3.4 Fast multipoint evaluation

The *fast multipoint evaluation* method is useful for fast evaluation of holonomic sequences when all arithmetic operations are assumed to have uniform cost. It is useful, for example, when all operations are done on floating-point numbers with the same precision, or when evaluating sequences over finite fields.

Fast multipoint evaluation allows evaluating a polynomial $p \in R[k]$ of degree d simultaneously at d points using $O(M(d) \log d)$ operations in R , with temporary storage of $O(d \log d)$ coefficients. The idea is that evaluating $p(c)$ is equivalent to computing the remainder of p upon polynomial division by $k - c$, so several remainders can be computed quickly as remainders by $(k - c_1) \cdots (k - c_n)$ in a tree-like fashion.

Applying fast multipoint evaluation to a polynomial matrix product, we obtain Algorithm 3.4.1, which is due to Chudnovsky and Chudnovsky [31].

Algorithm 3.4.1 Polynomial matrix product using fast multipoint evaluation

Input: $M \in R[k]^{r \times r}$, $n = m \times w$

Output: $\prod_{i=0}^{n-1} M(i)$

- 1: $[T_0, T_1, \dots, T_{m-1}] \leftarrow [M(k), M(k+1), \dots, M(k+m-1)]$
▷ Compute entrywise Taylor shifts of the matrix
 - 2: $U \leftarrow \prod_{i=0}^{m-1} T_i$ ▷ Binary splitting in $R[k]^{r \times r}$
 - 3: $[V_0, V_1, \dots, V_{w-1}] \leftarrow [U(0), U(m), \dots, U((w-1)m)]$ ▷ Fast multipoint evaluation
 - 4: **return** $\prod_{i=0}^{w-1} V_i$ ▷ Repeated multiplication in $R^{r \times r}$
-

We assume for simplicity of presentation that n is a multiple of the parameter m (in general, we can take $w = \lfloor n/m \rfloor$ and multiply by the remaining factors naively).

Taking $m \sim n^{1/2}$, Algorithm 3.4.1 requires $O(M(n^{1/2}) \log n)$ arithmetic operations in the ring R , using $O(n^{1/2} \log n)$ temporary storage during the fast multipoint evaluation step.

Bostan, Gaudry and Schost [18] improve the algorithm to obtain an $O(M(n^{1/2}))$ operation bound, which is the best available result for evaluating the n -th term of a holonomic sequence over a general ring. Algorithm 3.4.1 and some of its applications are studied further by Ziegler [112].

3.5 Rectangular splitting

We now consider holonomic sequences whose recurrence equation involves coefficients from a commutative ring C with unity as well as an additional, distinguished parameter x . The setting is as in the previous sections, but with $R = C[x]$. In other words, we

are considering holonomic sequences of polynomials of the parameter. This is equivalent to considering sequences of rational functions of the parameter, since we can clear denominators. We make the following definition.

Definition 3.5.1. A holonomic sequence $(c(n) \equiv c(x, n))_{n=0}^{\infty}$ is parametric over C (with parameter x) if it satisfies a linear recurrence equation of the form (3.2.1) with $M \in R[k]$ where $R = C[x]$.

Let H be a commutative C -algebra, and let $c(x, k)$ be a parametric holonomic sequence defined by a recurrence matrix $M \in C[x][k]^{r \times r}$ together with an initial vector $c(z, 0) \in H^r$.

Given some $z \in H$ and $n \in \mathbb{N}$, we wish to compute the single vector $c(z, n) \in H^r$ efficiently subject to the assumption that operations in H are expensive compared to operations in C . Accordingly, we distinguish between:

- *Coefficient operations* in C
- *Scalar operations* in H (additions in H and multiplications $C \times H \rightarrow H$)
- *Nonscalar multiplications* $H \times H \rightarrow H$

As a general principle, we wish to avoid nonscalar multiplications.

Example 3.5.2. The sequence of rising factorials

$$c(x, n) = x^{\overline{n}} = x(x+1) \cdots (x+n-1)$$

is first-order holonomic (hypergeometric) with the defining recurrence equation

$$c(x, n+1) = (n+x)c(x, n),$$

and parametric over $C = \mathbb{Z}$. In some applications, we wish to evaluate $c(z, n)$ for $z \in H$ where $H = \mathbb{R}$ or $H = \mathbb{C}$.

The Paterson-Stockmeyer algorithm [87] solves the problem of evaluating a polynomial $P(x) = \sum_{i=0}^{n-1} p_i x^i$ with $p_i \in C$ at $x = z \in H$ using a reduced number of nonscalar multiplications. The idea is to write the polynomial as a rectangular array

$$\begin{aligned} P(x) &= (p_0 + \dots + p_{m-1}x^{m-1}) \\ &\quad + (p_m + \dots + p_{2m-1}x^{m-1})x^m \\ &\quad + (p_{2m} + \dots + p_{3m-1}x^{m-1})x^{2m} \\ &\quad + \dots \end{aligned} \tag{3.5.1}$$

After computing a table containing x^2, x^3, \dots, x^{m-1} , the inner (rowwise) evaluations can be done using only scalar multiplications, and the outer (columnwise) evaluation with respect to x^m can be done using about n/m nonscalar multiplications. With $m \sim n^{1/2}$, this algorithm requires $O(n^{1/2})$ nonscalar multiplications and $O(n)$ scalar operations.

A straightforward application of the Paterson-Stockmeyer algorithm to evaluate each entry of $\prod_{i=0}^{n-1} M(x, i) \in C[x]^{r \times r}$ yields Algorithm 3.5.1 and the corresponding complexity estimate of Theorem 3.5.3.

Algorithm 3.5.1 Polynomial matrix product and evaluation using rectangular splitting

Input: $M \in C[x][k]^{r \times r}$, $z \in H$, $n = m \times w$

Output: $\prod_{i=0}^{n-1} M(z, i) \in H^{r \times r}$

- 1: $[T_0, \dots, T_{n-1}] \leftarrow [M(x, 0), \dots, M(x, n-1)]$
▷ Evaluate matrix w.r.t. k , giving $T_i \in C[x]^{r \times r}$
 - 2: $U \leftarrow \prod_{i=0}^{n-1} T_i$
▷ Binary splitting in $C[x]^{r \times r}$
 - 3: $V \leftarrow U(z)$
▷ Evaluate U entrywise using Paterson-Stockmeyer with step length m
 - 4: **return** V
-

Theorem 3.5.3. *The n -th entry in a parametric holonomic sequence can be evaluated using $O(n^{1/2})$ nonscalar multiplications, $O(n)$ scalar operations, and $O(M(n) \log n)$ coefficient operations.*

Proof. We call Algorithm 3.5.1 with $m \sim n^{1/2}$. Letting $d = \max \deg_k M$ and $e = \max \deg_x M$, computing T_0, \dots, T_{n-1} takes $O(nde) = O(n)$ coefficient operations. Since $\deg_x T_i \leq e$, generating U using binary splitting costs $O(M(n) \log n)$ coefficient operations. Each entry in U has degree at most $ne = O(n)$, and can thus be evaluated using $O(n^{1/2})$ nonscalar multiplications and $O(n)$ scalar operations with the Paterson-Stockmeyer algorithm. □

If we only count nonscalar multiplications, Theorem 3.5.3 is an asymptotic improvement over fast multipoint evaluation which uses $O(n^{1/2} \log^{2+o(1)} n)$ nonscalar multiplications ($O(n^{1/2} \log^{1+o(1)} n)$ with the improvement of Bostan, Gaudry and Schost).

Algorithm 3.5.1 is not ideal in practice since the polynomials in U grow to degree $O(n)$. Their coefficients also grow to $O(n \log n)$ bits when $C = \mathbb{Z}$ (for example, in the case of rising factorials, the coefficients are the Stirling numbers of the first kind $S(n, k)$ which grow to a magnitude between $(n-1)!$ and $n!$).

This problem can be mitigated by repeatedly applying Algorithm 3.5.1 to successive subproducts $\prod_{i=a}^{a+\tilde{n}} M(z, i)$ where $\tilde{n} \ll n$, but the nonscalar complexity is then no longer the best possible.

A better strategy is to apply rectangular splitting to the matrix product itself, leading to Algorithm 3.5.2. We can then reach the same operation complexity while only working with polynomials of degree $O(n^{1/2})$, and over $C = \mathbb{Z}$, having coefficients of bit size $O(n^{1/2} \log n)$.

Theorem 3.5.4. *For any choice of m , Algorithm 3.5.2 requires $O(m + n/m)$ non-scalar multiplications, $O(n)$ scalar operations, and $O((n/m)\mathbf{M}(m)\log m)$ coefficient operations. In particular, the complexity bounds stated in Theorem 3.5.3 also hold for Algorithm 3.5.2 with $m \sim n^{1/2}$. Moreover, Algorithm 3.5.2 only requires storage of $O(m)$ elements of C and H , and if $C = \mathbb{Z}$, the coefficients have bit size $O(m \log m)$.*

Proof. This follows by applying a similar argument as used in the proof of Theorem 3.5.3 to the operations in the inner loop of Algorithm 3.5.2, noting that U has entries of degree $m \deg_x M = O(m)$ and that the matrix multiplication $S \times V$ requires $O(1)$ nonscalar multiplications and scalar operations (recalling that we consider r fixed). \square

Algorithm 3.5.2 Improved polynomial matrix product and evaluation using rectangular splitting

Input: $M \in C[x][k]^{r \times r}$, $z \in H$, $n = m \times w$

Output: $\prod_{i=0}^{n-1} M(z, i) \in H^{r \times r}$

- 1: Compute power table $[z^j, 0 \leq j \leq m \deg_x M]$
 - 2: $V \leftarrow 1_{H^{r \times r}}$ ▷ Start with the identity matrix
 - 3: **for** $i \leftarrow 0$ to $w - 1$ **do**
 - 4: $[T_0, \dots, T_{m-1}] \leftarrow [M(x, im + j)]_{j=0}^{m-1}$ ▷ Evaluate matrix w.r.t. k , giving $T_j \in C[x]^{r \times r}$
 - 5: $U \leftarrow \prod_{j=0}^{m-1} T_j$ ▷ Binary splitting in $C[x]^{r \times r}$
 - 6: $S \leftarrow U(z)$ ▷ Evaluate w.r.t. x using power table
 - 7: $V \leftarrow S \times V$ ▷ Multiplication in $H^{r \times r}$
 - 8: **return** V
-

3.5.1 Variations

Many variations of Algorithm 3.5.2 are possible. Instead of using binary splitting directly to compute U , we can generate the bivariate matrix

$$W_m = \prod_{i=0}^{m-1} M(x, k + i) \in C[x][k]^{r \times r} \quad (3.5.2)$$

at the start of the algorithm, and then obtain U by evaluating W_m at $k = 0, m, 2m, \dots$. We may also work with differences of two successive U (for small m , this can introduce cancellation resulting in slightly smaller polynomials or coefficients). Combining both

variations, we end up with Algorithm 3.5.3 in which we expand and evaluate the bivariate polynomial matrices

$$\Delta_m = \prod_{i=0}^{m-1} M(x, k + m + i) - \prod_{i=0}^{m-1} M(x, k + i) \in C[x][k]^{r \times r}.$$

This version of the rectangular splitting algorithm can be viewed as a generalization of an algorithm used by Smith [96] for computing rising factorials (we consider the case of rising factorials further in Section 3.6.1).

In fact, the author of this thesis first found Algorithm 3.5.3 by generalizing Smith's algorithm, and only later discovered Algorithm 3.5.2 by "interpolation" between Algorithm 3.5.1 and Algorithm 3.5.3.

Algorithm 3.5.3 Polynomial matrix product and evaluation using rectangular splitting (variation)

Input: $M \in C[x][k]^{r \times r}$, $z \in H$, $n = m \times w$

Output: $\prod_{i=0}^{n-1} M(z, i) \in H^{r \times r}$

- 1: Compute power table $[z^j, 0 \leq j \leq m \deg_x M]$
 - 2: $\Delta \leftarrow \prod_{i=0}^{m-1} M(x, k + m + i) - \prod_{i=0}^{m-1} M(x, k + i)$ ▷ Binary splitting in $C[x][k]^{r \times r}$
 - 3: $V \leftarrow S \leftarrow \prod_{i=0}^{m-1} M(z, i)$ ▷ Evaluate w.r.t. k , and w.r.t. x using power table
 - 4: **for** $i \leftarrow 0$ to $w - 2$ **do**
 - 5: $S \leftarrow S + \Delta(z, mi)$ ▷ Evaluate w.r.t. k , and w.r.t. x using power table
 - 6: $V \leftarrow S \times V$
 - 7: **return** V
-

The efficiency of Algorithm 3.5.3 is theoretically somewhat worse than that of Algorithm 3.5.2. Since $\deg_x W_m = O(m)$ and $\deg_k W_m = O(m)$, W_m has $O(m^2)$ terms (likewise for Δ_m), making the space complexity higher and increasing the number of coefficient operations to $O((n/m)m^2)$ for the evaluations with respect to k . However, this added cost may be negligible in practice. Crucially, when $C = \mathbb{Z}$, the coefficients have similar bit sizes to those in Algorithm 3.5.2.

Generating W_m or Δ_m at the start of the algorithm also adds some cost, but this is cheap compared to the evaluations when n is large enough: binary splitting over $C[x][k]$ costs $O(M(m^2) \log m)$ coefficient operations by Lemma 8.2 and Corollary 8.28 in [107]. This is essentially the same as the total cost of binary splitting in Algorithm 3.5.2 when $m \sim n^{1/2}$.

We also note that a small improvement to Algorithm 3.5.2 is possible if $M(x, k + m) = M(x + m, k)$: instead of computing U from scratch using binary splitting in each loop iteration, we can update it using a Taylor shift. At least in sufficiently large characteristic, the Taylor shift can be computed using $O(M(m))$ coefficient operations with the

convolution algorithm of Aho, Steiglitz and Ullman [1], saving a factor $O(\log n)$ in the total number of coefficient operations. In practice, basecase Taylor shift algorithms may also be beneficial (see [106]).

In lucky cases, the polynomial coefficients (in either Algorithm 3.5.2 or 3.5.3) might satisfy a recurrence relation, allowing them to be generated using $O(n)$ coefficient operations (and avoiding the dependency on polynomial arithmetic).

3.5.2 Several parameters

The rectangular splitting technique can be generalized to sequences $c(x_1, \dots, x_v, k)$ depending on several parameters. In Algorithm 3.5.2, we simply replace the power table by a v -dimensional array of the possible monomial combinations. Then we have the following result (ignoring coefficient operations).

Theorem 3.5.5. *The n -th entry in a holonomic sequence depending on v parameters can be evaluated with rectangular splitting using $O(m^v + n/m)$ nonscalar multiplications and $O((n/m)m^v)$ scalar multiplications. In particular, taking $m = n^{1/(v+1)}$, $O(n^{1-1/v})$ nonscalar multiplications and $O(n^{2v/(1+v)})$ scalar multiplications suffice.*

Proof. If $d_i = \deg_{x_i} M \leq d$, the entries of a product of m successive shifts of M are C -linear combinations of $x_1^{e_{1,j}} \dots x_n^{e_{n,j}}$, $0 \leq e_{i,j} \leq md_i \leq md$, so there is a total of $O(m^v)$ powers. □

Unfortunately, this gives rapidly diminishing returns for large v . When $v > 1$, the number of nonscalar multiplications according to Theorem 3.5.5 is asymptotically worse than with fast multipoint evaluation, and reducing the number of nonscalar multiplications requires us to perform more than $O(n)$ scalar multiplications, as shown in Table 3.1. Nevertheless, rectangular splitting could perhaps still be useful in some settings where the cost of nonscalar multiplications is sufficiently large.

v	m	Nonscalar	Scalar
1	$n^{1/2}$	$O(n^{0.5})$	$O(n)$
2	$n^{1/3}$	$O(n^{0.666\dots})$	$O(n^{1.333\dots})$
3	$n^{1/4}$	$O(n^{0.75})$	$O(n^{1.5})$
4	$n^{1/5}$	$O(n^{0.8})$	$O(n^{1.6})$

Table 3.1: Step size m minimizing the number of nonscalar multiplications for rectangular splitting involving v parameters.

3.5.3 Numerical evaluation

Assume that we want to evaluate $c(x, n)$ where the underlying coefficient ring is $C = \mathbb{Z}$ (or $\overline{\mathbb{Q}}$) and the parameter x is a real or complex number represented by a floating-point approximation with a precision of p bits.

The naive algorithm clearly uses $O(nM_{\mathbb{Z}}(p))$ bit operations to evaluate $c(x, n)$, or $O^{\sim}(np)$ with FFT multiplication.

In Algorithm 3.5.2, the nonscalar multiplications cost $O((m + n/m)M_{\mathbb{Z}}(p))$ bit operations. The coefficient operations cost $O(mn)$ bit operations (assuming the use of fast polynomial arithmetic), which becomes negligible if p grows faster than m .

Finally, the scalar multiplications (which are unbalanced) cost

$$O\left(np \frac{M_{\mathbb{Z}}(m \log m)}{m \log m}\right)$$

bit operations. Taking $m \sim n^{\alpha}$ for $0 < \alpha < 1$, we get an asymptotic speedup with classical or Karatsuba multiplication (see Table 3.2) provided that p grows sufficiently rapidly along with n .

With FFT multiplication, the scalar multiplications become as expensive as the nonscalar multiplications, and rectangular splitting therefore does not give an asymptotic improvement.

Mult. algorithm	Scalar multiplications	Naive
Classical	$O^{\sim}(n^{1+\alpha}p)$	$O^{\sim}(np^2)$
Karatsuba	$O^{\sim}(n^{1+0.585\alpha}p)$	$O^{\sim}(np^{1.585})$
FFT	$O^{\sim}(np)$	$O^{\sim}(np)$

Table 3.2: Bit complexity of scalar multiplications in Algorithm 3.5.2 and total bit complexity of the naive algorithm

However, due to the overhead of FFT multiplication, rectangular splitting is still likely to save a constant factor over the naive algorithm. In practice, one does not necessarily get the best performance by choosing $m \approx n^{0.5}$ to minimize the number of nonscalar multiplications alone; the best m has to be determined empirically.

Algorithm 3.4.1 is asymptotically faster than the naive algorithm as well as rectangular splitting, with a bit complexity of $O^{\sim}(n^{1/2}p)$. It should be noted that this estimate does not reflect the complexity required to obtain a given *accuracy*. As observed by Köhler and Ziegler [66], fast multipoint evaluation can exhibit poor numerical stability, suggesting that p might have to grow at least as fast as n to get accuracy proportional to p .

When x and all coefficients in M are positive, rectangular splitting introduces no subtractions that can cause catastrophic cancellation, and the reduction of nonscalar multiplications even improves stability compared to the naive algorithm, making $O(\log n)$ guard bits sufficient to reach p -bit accuracy. When sign changes are present, evaluating degree- m polynomials in expanded form can reduce accuracy, typically requiring use of $O^\sim(m)$ guard bits. In this case Algorithm 3.5.2 is a marked improvement over Algorithm 3.5.1.

3.5.4 Summation of power series

A common situation is that we wish to evaluate a truncated power series

$$f(x) \approx s(x, n) = \sum_{k=0}^n c(k)x^k, \quad n = O(p) \quad (3.5.3)$$

where $c(k)$ is a holonomic sequence taking rational (or algebraic) values and x is a real or complex number. In this case the Paterson-Stockmeyer algorithm is applicable, but might not give a speedup when applied directly as in Algorithm 3.5.1 due to the growth of the coefficients. Since $d(k) = c(k)x^k$ and $s(x, n)$ are holonomic sequences with x as parameter, Algorithm 3.5.2 is applicable.

Smith noted in [95] that when $c(k)$ is hypergeometric (Smith considered the Taylor expansions of elementary functions in particular), the Paterson-Stockmeyer technique can be combined with scalar divisions to remove accumulated factors from the coefficients. This keeps all scalars at a size of $O(\log n)$ bits, giving a speedup over naive evaluation when non-FFT multiplication is used (and when scalar divisions are assumed to be roughly as cheap as scalar multiplications). This algorithm is studied in more detail by Brent and Zimmermann [27].

At least conceptually, Algorithm 3.5.2 can be viewed as a generalization of Smith's hypergeometric summation algorithm to arbitrary holonomic sequences depending on a parameter (both algorithms can be viewed as means to eliminate repeated content from the associated matrix product). The speedup is not quite as good since we only reduce the coefficients to $O(n^{1/2} \log n)$ bits versus Smith's $O(\log n)$. However, even for hypergeometric series, Algorithm 3.5.2 can be slightly faster than Smith's algorithm for small n (e.g. $n \lesssim 100$) since divisions tend to be more expensive than scalar multiplications in implementations.

Algorithm 3.5.2 is also more general: for example, we can use it to evaluate the generalized hypergeometric function

$${}_pF_q \left(\begin{matrix} a_1, \dots, a_p \\ b_1, \dots, b_q \end{matrix} \middle| z \right) = \sum_{k=0}^{\infty} \frac{a_1^{\overline{k}} \cdots a_p^{\overline{k}}}{b_1^{\overline{k}} \cdots b_q^{\overline{k}}} \frac{w^k}{k!} \quad (3.5.4)$$

where a_i, b_i, w (as opposed to w alone) are rational functions of the real or complex parameter x .

An interesting question, which we do not attempt to answer here, is whether there is a larger class of parametric sequences other than hypergeometric sequences and their sums for which we can reduce the number of nonscalar multiplications to $O(n^{1/2})$ while working with coefficients that are strictly smaller than $O(n^{1/2} \log n)$ bits.

If all coefficients in (3.5.3) including the parameter x are rational or algebraic numbers and the series converges, $f(x)$ can be evaluated to p -bit precision using $O^\sim(p)$ bit operations using binary splitting. Combined with the bit-burst technique, an $O^\sim(p)$ bit complexity can also be achieved for arbitrary real or complex x .

For high-precision evaluation of elementary functions, the bit-burst algorithm typically only becomes worthwhile at a precision of several thousand digits, while implementations typically use Smith's algorithm for summation of hypergeometric series at lower precision. We expect that Algorithm 3.5.2 can be used in a similar fashion for a larger class of special functions.

When $c(k)$ in (3.5.3) involves real or complex numbers, binary splitting no longer gives a speedup. In this case, we can use fast multipoint evaluation to compute (3.5.3) using $O^\sim(p^{1.5})$ bit operations (Borwein [16] discusses the application to numerical evaluation of hypergeometric functions). This method does not appear to be widely used in practice, presumably owing to its high overhead and relative implementation difficulty. Although rectangular splitting is not as fast asymptotically, its ease of implementation and low overhead makes it an attractive alternative.

3.6 Computing the gamma function

In this section, we consider two holonomic sequences depending on a numerical parameter: rising factorials, and the partial sums of a certain hypergeometric series defining the incomplete gamma function. In both cases, our goal is to accelerate numerical evaluation of the gamma function at very high precision.

We have implemented the algorithms in the present section using floating-point ball arithmetic (with rigorous error bounding) as part of the Arb library. All benchmark results were obtained on a 2.0 GHz Intel Xeon E5-2650 CPU.

3.6.1 Rising factorials

Rising factorials of a real or complex argument appear when evaluating the gamma function via the asymptotic Stirling series

$$\log \Gamma(x) = \left(x - \frac{1}{2}\right) \log x - x + \frac{\log 2\pi}{2} + \sum_{k=1}^{N-1} \frac{B_{2k}}{2k(2k-1)x^{2k-1}} + R_N(x).$$

To compute $\Gamma(x)$ with p -bit accuracy, we choose a positive integer n such that there is an N for which $|R_N(x+n)| < 2^{-p}$, and then evaluate $\Gamma(x) = \Gamma(x+n)/x^{\bar{n}}$. It is sufficient to choose n such that the real part of $x+n$ is of order βp where $\beta = (2\pi)^{-1} \log 2 \approx 0.11$.

The efficiency of the Stirling series can be improved by choosing n slightly larger than the absolute minimum in order to reduce N . For example, $\operatorname{Re}(x+n) \approx 2\beta p$ is a good choice. A faster rising factorial is doubly advantageous: it speeds up the argument reduction, and making larger n cheap allows us to get away with fewer Bernoulli numbers.

Smith [96] uses the difference of four consecutive terms

$$\begin{aligned} (x+k+4)^{\bar{4}} - (x+k)^{\bar{4}} &= (840 + 632k + 168k^2 + 16k^3) \\ &\quad + (632 + 336k + 48k^2)x \\ &\quad + (168 + 48k)x^2 \\ &\quad + 16x^3 \end{aligned}$$

to reduce the number of nonscalar multiplications to compute $x^{\bar{n}}$ from $n-1$ to about $n/4$. This is precisely Algorithm 3.5.3 specialized to the sequence of rising factorials and with a fixed step length $m=4$.

Consider Smith's algorithm with a variable step length m . Using the binomial theorem and some rearrangements, the polynomials can be written down explicitly as

$$\Delta_m = (x+k+m)^{\bar{m}} - (x+k)^{\bar{m}} = \sum_{v=0}^{m-1} x^v \sum_{i=0}^{m-v-1} k^i C_m(v, i) \quad (3.6.1)$$

where

$$C_m(v, i) = \sum_{j=i+1}^{m-v} m^{j-i} S(m, v+j) \binom{v+j}{v} \binom{j}{i} \quad (3.6.2)$$

and where $S(m, v+j)$ denotes an unsigned Stirling number of the first kind. This formula can be used to generate Δ_m efficiently in practice without requiring bivariate polynomial arithmetic. In fact, the coefficients can be generated even cheaper by taking advantage of the recurrence (found by M. Kauers)

$$(v+1)C_m(v+1, i) = (i+1)C_m(v, i+1). \quad (3.6.3)$$

We have implemented several algorithms for evaluating the rising factorial of a real or complex number. For tuning parameters, we empirically determined simple formulas that give nearly optimal performance for different combinations of $n, p < 10^5$ (typically within 20% of the speed with the best tuning value found by a brute force search):

- In Algorithm 3.4.1, $m = n^{0.5}$.
- Algorithm 3.5.1 is applied to subproducts of length $\tilde{n} = \min(2n^{0.5}, 10p^{0.25})$, with $m = \tilde{n}^{0.5}$.
- In Algorithms 3.5.2 and 3.5.3, $m = \min(0.2p^{0.4}, n^{0.5})$.

Our implementation of Algorithm 3.5.3 uses (3.6.2) instead of binary splitting, and Algorithm 3.5.2 exploits the symmetry of x and k to update the matrix U using Taylor shifts instead of repeated binary splitting.

Figure 3.6.1 compares the running times where x is a real number with a precision of $p = 4n$ bits. This input corresponds to that used in our Stirling series implementation of the gamma function.

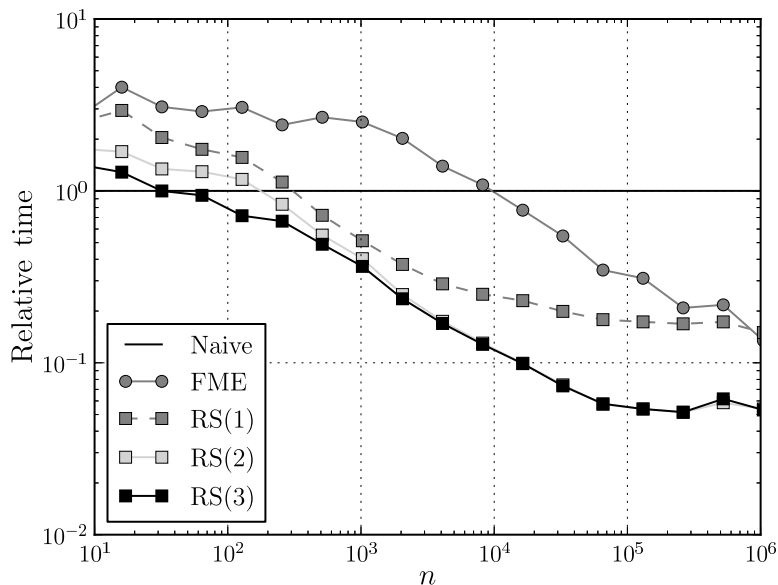


Figure 3.1: Timings of rising factorial algorithms, normalized against the naive algorithm. FME: fast multipoint evaluation (Algorithm 3.4.1), RS(1): Algorithm 3.5.1, RS(2): Algorithm 3.5.2, RS(3): Algorithm 3.5.3.

On this benchmark, Algorithms 3.5.2 and 3.5.3 are the best by far, gaining a 20-fold speedup over the naive algorithm for large n (the speedup levels off around $n = 10^5$, which is expected since this is the approximate point where FFT integer multiplication

kicks in). Algorithm 3.5.3 is slightly faster than Algorithm 3.5.2 for $n < 10^3$, even narrowly beating the naive algorithm for n as small as $\approx 10^2$.

Algorithm 3.4.1 (fast multipoint evaluation) has the most overhead of all algorithms and only overtakes the naive algorithm around $n = 10^4$ (at a precision of 40,000 bits). Despite its theoretical advantage, it is slower than rectangular splitting up to n exceeding 10^6 .

Table 3.6.1 shows absolute timings for evaluating $\Gamma(x)$ where x is a small real number in Pari/GP 2.5.4, and our implementation in Arb (we omit results for MPFR 3.1.1 and Mathematica 9.0, which were both slower than Pari). Both implementations use the Stirling series, caching the Bernoulli numbers to speed up multiple evaluations. The better speed of Arb for a repeated evaluation (where the Bernoulli numbers are already cached) is mainly due to the use of rectangular splitting to evaluate the rising factorial. The total speedup is smaller than it would be for computing the rising factorial alone since we still have to evaluate the Bernoulli number sum in the Stirling series. The gamma function implementations over \mathbb{C} have similar characteristics.

Decimals	Pari/GP	(first)	Arb	(first)
100	0.000088		0.00010	
300	0.00048		0.00036	
1000	0.0057		0.0025	
3000	0.072	(9.2)	0.021	(0.090)
10000	1.2	(324)	0.25	(1.4)
30000	15	(8697)	2.7	(22)
100000			39	(433)
300000			431	(7131)

Table 3.3: Timings in seconds for evaluating $\Gamma(x)$ where x is a small real number (timings for the first evaluation, including Bernoulli number generation, is shown in parentheses).

3.6.2 A one-parameter hypergeometric series

The gamma function can be approximated via the (lower) incomplete gamma function as

$$\Gamma(z) \approx \gamma(z, N) = z^{-1} N^z e^{-N} {}_1F_1(1, 1 + z, N). \quad (3.6.4)$$

Borwein [16] noted that applying fast multipoint evaluation to a suitable truncation of the hypergeometric series in (3.6.4) allows evaluating the gamma function of a fixed real or complex argument to p -bit precision using $O^\sim(p^{1.5})$ bit operations, which is the best known result for general z (if z is algebraic, binary splitting evaluation of the same series achieves a complexity of $O^\sim(p)$, as noted by Brent [22]).

Let $t_k = N^k / (z(z+1)\cdots(z+k))$ and $s_n = \sum_{k=0}^n t_k$, giving

$$\lim_{n \rightarrow \infty} s_n = {}_1F_1(1, 1+z, N)/z.$$

For $z \in [1, 2]$, choosing $N \approx p \log 2$ and $n \approx (e \log 2)p$ gives an error of order 2^{-p} (it is easy to compute strict bounds). The partial sums satisfy the order-2 recurrence

$$\begin{pmatrix} s_k \\ t_{k+1} \end{pmatrix} = \frac{M(k)}{q(k)} \frac{M(k-1)}{q(k-1)} \cdots \frac{M(0)}{q(0)} \begin{pmatrix} 0 \\ 1/z \end{pmatrix} \quad (3.6.5)$$

where

$$M(k) = \begin{pmatrix} 1+k+z & 1+k+z \\ 0 & N \end{pmatrix}, \quad q(k) = 1+k+z. \quad (3.6.6)$$

The matrix product (3.6.5) may be computed using fast multipoint evaluation or rectangular splitting. We note that the denominators are identical to the top left entries of the numerator matrices, and therefore do not need to be computed separately.

Figure 3.6.2 compares the performance of the Stirling series (with fast argument reduction using rectangular splitting) and three different implementations of the ${}_1F_1$ series (naive summation, fast multipoint evaluation, and rectangular splitting using Algorithm 3.5.2 with $m = 0.2n^{0.4}$) for evaluating $\Gamma(x)$ where x is a real argument close to unity.

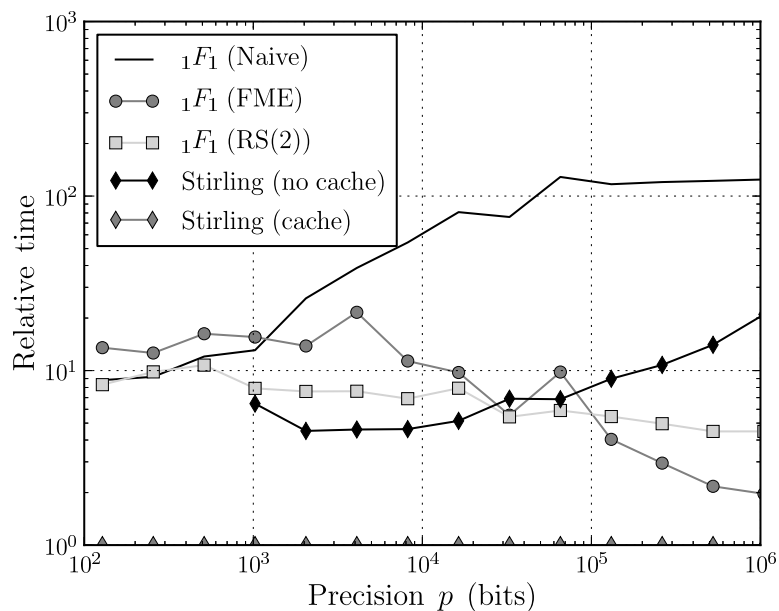


Figure 3.2: Timings of gamma function algorithms, normalized against the Stirling series with Bernoulli numbers cached. FME: fast multipoint evaluation (Algorithm 3.4.1), RS(2): Algorithm 3.5.2.

Both fast multipoint evaluation and rectangular splitting speed up the hypergeometric series compared to naive summation. Using either algorithm, the hypergeometric series is competitive with the Stirling series for a single evaluation at precisions above roughly 10,000 decimal digits.

Algorithm 3.4.1 performs better than on the rising factorial benchmark, and is faster than Algorithm 3.5.2 above 10^5 bits. A possible explanation for this difference is that roughly $n \approx 2p$ terms are added in the hypergeometric series (where p is the precision in bits), compared to $n \approx p/4$ for the rising factorial, and rectangular splitting favors higher precision and fewer terms.

The speed of Algorithm 3.5.2 is remarkably close to that of Algorithm 3.4.1 even for p as large as 10^6 . Despite being asymptotically slower, the simplicity of rectangular splitting combined with its lower memory consumption and better numerical stability (in our implementation, Algorithm 3.5.2 only loses a few significant digits, while Algorithm 3.4.1 loses a few percent of the number of significant digits) makes it an attractive option for extremely high-precision evaluation of the gamma function.

Once the Bernoulli numbers have been cached after the first evaluation, the Stirling series still has a clear advantage up to precisions exceeding 10^6 bits. We may remark that our implementation of the Stirling series has been optimized for multiple evaluations: by choosing larger rising factorials and generating the Bernoulli numbers dynamically without storing them, both the speed and memory consumption for a single evaluation could be improved.

3.6.3 Remarks

We have shown that rectangular splitting can be profitably applied to evaluation of a general class of linearly recurrent sequences. When used for numerical evaluation of special functions, our benchmark results indicate that rectangular splitting can be faster than either naive evaluation or fast multipoint evaluation over a wide precision range (between approximately 10^3 and 10^6 bits).

Two natural questions are whether our approach can be extended to more general classes of sequences, and whether it can be optimized further, perhaps for more specific classes of sequences. A partial answer to the first question is that rectangular splitting can be applied to any product of polynomial matrices $\prod_i M(i)$, $M(i) \in R[x]$, not just those where the entries of $M(i)$ are polynomials in i . Much of the analysis we have presented is valid assuming only that the entries grow at most as polynomials in i , in particular assuming that $\deg_x M(i)$ is bounded. Rectangular splitting is analogous to binary splitting in this sense, whereas the fast multipoint evaluation algorithm does not seem to generalize in such a way.

Chapter 4

Evaluating functions of power series

Interesting number sequences can often be interpreted as coefficients in the series expansions of analytic functions. Power series are a powerful mathematical and computational paradigm: rather than viewing $f(z), f'(z), \dots$ as separate values, we treat the power series $f(z) + f'(z)x + \dots$ as a single object and phrase transformations of the coefficients as operations done on power series. This provides at least two benefits: firstly, we often make life easier by avoiding messy explicit formulas for the individual coefficients, and secondly, we can take advantage of fast algorithms for power series arithmetic.

4.1 Bernoulli numbers

The Bernoulli numbers are defined by the generating function

$$\frac{x}{e^x - 1} = \sum_{k=0}^{\infty} B_k \frac{x^k}{k!}. \quad (4.1.1)$$

Bernoulli numbers appear in various expansions of special functions. We have already used them in the calculation of the gamma function, and will later use them when calculating the Hurwitz zeta function.

An efficient way to generate Bernoulli numbers is to construct a truncation of the power series

$$\frac{e^x - 1}{x} = \sum_{k=0}^{\infty} \frac{x^k}{(k+1)!} \quad (4.1.2)$$

and compute its multiplicative inverse using Newton iteration (see section 4.4). This allows us to compute B_0, \dots, B_n using $O(M(n))$ ring operations. It allows us, in particular, to generate the Bernoulli numbers as exact fractions in time $O^\sim(n^2)$, which is

softly optimal since the numerators and denominators take $\Omega(n^2)$ bits of space to write down.

Moreover, if we want to compute the Bernoulli numbers modulo a prime $p > n$, doing all arithmetic modulo p reduces the complexity to $O^\sim(n \log p)$ which again is softly optimal. This algorithm extends to other numbers with generating functions given by compositions of arithmetic operations and elementary functions.

Recursive algorithms with a bit complexity of $O^\sim(n^3)$ are often adequate in practice. Some alternatives are discussed in [23]. An interesting algorithm, used in unpublished work of Bloemen [10], is to compute B_n via $\zeta(n)$ by direct approximation of the sum $\sum_{k=0}^{\infty} k^{-n}$, recycling the powers to process several n simultaneously. This algorithm has suboptimal complexity $O^\sim(n^3)$, but the implied constant is extremely small. For multi-evaluation of Bernoulli numbers in FLINT and Arb, Bloemen's method turned out to be faster than power series inversion for n as large as 10^5 .

4.2 Defining functions of power series

As usual, $R[[x]]$ denotes the ring of formal power series over a commutative ring R , and

$$R[[x]]/\langle x^n \rangle \cong R[x]/\langle x^n \rangle$$

denotes the ring of formal power series truncated to degree less than n . Given a formal power series $F = \sum_{k=0}^{\infty} f_k x^k$, we sometimes use the notation $[x^k]F$ for the coefficient f_k . The arithmetic operations (addition, multiplication, reciprocal and division) are defined in the usual way.

If $F, G \in R[[x]]$ and $[x^0]G = 0$, then the *composition* of $F = \sum_{k=0}^{\infty} f_k$ and $G = \sum_{k=0}^{\infty} g_k$ is defined to be the power series $H = F(G) = f_0 + f_1 G + f_2 G^2 + \dots$. Since $[x^n]H$ only depends on f_k, g_k for $k \leq n$, formal composition of power series commutes with truncation.

Let f be an analytic function defined on some domain $U \subseteq \mathbb{C}$ and let

$$G = g_0 + g_1 x + g_2 x^2 + \dots \in \mathbb{C}[[x]]$$

where $G_0 \in U$. Then we define $f(G)$ as the object

$$f(G) = f(g_0) + f'(g_0)(G - g_0) + \frac{f''(g_0)}{2!}(G - g_0)^2 \dots \in \mathbb{C}[[x]].$$

If G is the Taylor series expansion of an analytic function g at some point z and $g(z) \in U$, then $f(G)$ is the Taylor series expansion of $f(g)$ at z (note, however, that $f(G)$ is well-defined even if G has zero radius of convergence). By standard uniqueness theorems from complex analysis, identities involving functions also extend to identities of power

series: if it holds for all $z \in U$ that $f(z) = h(z)$, then it holds for all $G \in \mathbb{C}[[x]]$ with $[x^0]G \in U$ that $f(G) = h(G)$.

Since $G - g_0 = O(x)$, $[x^n]f(G)$ only depends on the finitely many terms $g_0 \dots g_n$. This means that the operation of evaluating a function of a power series commutes with power series truncation, and therefore, identities of functions over $\mathbb{C}[[x]]$ are also valid over $\mathbb{C}[[x]]/\langle x^n \rangle$.

In many situations when working with numerical power series, we do need to work with non-formal (i.e. analytic) limits, which are defined coefficientwise. We thus say that $F = \lim_{n \rightarrow \infty} F_n$ if $[x^k]F_n \rightarrow [x^k]F$ for all k . It is useful to have a convenient notation for coefficientwise bounds of power series: if $F = \sum_k f_k x^k \in \mathbb{C}[[x]]$, we define $|F| = \sum_k |f_k| x^k$ (this notation does not overload with the previously-given definition of applying a function to a power series, since $|\cdot|$ is not analytic). If $\forall k : |f_k| \leq |g_k|$, we write $|F| \leq |G|$ (in this case, G is said to be a majorant series of F). It is easy to check that $|F + G| \leq |F| + |G|$ and $|FG| \leq |F||G|$ for all F, G .

If we parse a formula for a function such that all compositions are formal, then “computing a function of a power series” is effectively the same thing as “computing a power series of a function”. For example, the expression $\log(3 + \sin(x))$, interpreted as a power series, should be read as $F(G)$ where $F = \log(3 + x)$ and $G = \sin(x)$, and not, for example, as $F(G)$ where $F = \log(1 + x)$ and $G = 2 + \sin(x)$. The latter composition neither converges formally nor analytically.

4.3 Fast composition

Classical algorithms for formal composition of power series truncated to length n require $O(n^3)$ operations on elements in the coefficient ring [64]. This bound can be improved to $O(nM(n))$ by use of fast polynomial multiplication. In [25], Brent and Kung gave two asymptotically faster algorithms for composition.

The first algorithm (BK 2.1) is a baby-step giant-step version of Horner’s rule, with the additional improvement that the scalar multiplications are grouped into a single matrix multiplication.

Suppose we want to compute $f(g) \in R[[x]]/\langle x^n \rangle$. Choosing a step length $m \geq n^{1/2}$, let A be the $m \times m$ matrix populated by the coefficients in f in row-major order, let B be the $m \times n$ matrix whose rows are the first m powers of g as vectors of coefficients, and let $c_{i,j}$ denote the coefficient at row i and column j in the $m \times n$ matrix $C = AB$

(indexed from zero). Then

$$f(g) = \sum_{i=0}^{m-1} \left(\sum_{j=0}^{n-1} c_{i,j} x^j \right) (g^m)^i.$$

For example, with $n = 9, m = 3$, C is the 3×9 matrix obtained by multiplying

$$A = \begin{pmatrix} [x^0]f & [x^1]f & [x^2]f \\ [x^3]f & [x^4]f & [x^5]f \\ [x^6]f & [x^7]f & [x^8]f \end{pmatrix}$$

with

$$B = \begin{pmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ [x^0]g & [x^1]g & [x^2]g & [x^3]g & \dots & [x^8]g \\ [x^0]g^2 & [x^1]g^2 & [x^2]g^2 & [x^3]g^2 & \dots & [x^8]g^2 \end{pmatrix}.$$

BK 2.1 allows computing a length- n composition in $O(n^{1/2}(\mathbf{M}(n) + \mathbf{MM}(n^{1/2})))$ operations. This bound reduces to $O(n^{1/2}\mathbf{M}(n) + n^2)$ with classical matrix multiplication, and to $O(n^{1/2}\mathbf{M}(n) + n^{1.91})$ with the Strassen algorithm. The last term can be improved to $O(n^{1.688})$ with the Coppersmith-Winograd algorithm [107], or to $O(n^{1.68632})$ with the recent bound for $\mathbf{MM}(n)$ by Stothers and Vassilevska Williams [98, 104]. The best available bound for the last term is $O(n^{1.667})$, using the improved technique for multiplication of nonsquare matrices of Huang and Pan [51].

The second algorithm (BK 2.2) uses formal Taylor expansion to break the composition into smaller pieces of carefully chosen size, and requires $O((n \log n)^{1/2}\mathbf{M}(n))$ operations. This is asymptotically slower than BK 2.1 when classical multiplication ($\mathbf{M}(n) = O(n^2)$) or Karatsuba multiplication ($\mathbf{M}(n) = O(n^{\log_2 3}) = O(n^{1.585})$) is used, but faster when FFT polynomial multiplication ($\mathbf{M}(n) = O(n \log^{1+o(1)} n)$) is available [9, 107].

4.4 Elementary functions

As noted by Brent and Kung, many special left-compositions, including the evaluation of reciprocals, square roots, and elementary transcendental functions of power series, can be done in merely $O(\mathbf{M}(n))$ operations.

The main tool is Newton's method for root-finding. Suppose $F \in R[[x]]$ with $[x^0]F = 0$ and $[x^1]F \neq 0$, and suppose $h_k \in R[[x]]$ satisfies $F(h_k) \equiv 0 \pmod{x^n}$. Then

$$h_{k+1} = h_k - \frac{F(h_k)}{F'(h_k)} \equiv 0 \pmod{x^{2n}}.$$

By truncating each iterate h_k at the number of correct terms, Newton's iteration allows computing a root of F to order $O(x^n)$ using $O(n) + C(n) + C(n/2) + C(n/4) + \dots$ ring operations, where $C(n)$ is the cost of the composition $F(h_k)/F'(h_k)$ to order $O(x^n)$.

Computing an inverse function $h = f^{-1}(g)$ is equivalent to finding a root of $F(h) = f(h) - g$. If left-composition of f has complexity at least $M(n)$, then the cost of left-composition by f^{-1} is the same up to a constant factor. In particular, composition by algebraic functions has complexity $O(M(n))$. (In the special case of computing a reciprocal $1/f$, needed in general for the division in $F(h_k)/F'(h_i)$, the Newton iteration reduces to $h_{k+1} = 2h_k - fh_k^2$ which only depends on multiplication.)

The logarithm and inverse trigonometric functions of power series can be computed in $O(M(n))$ operations as formal integrals of algebraic functions, i.e.

$$\begin{aligned}\log(f(x)) &= \int \frac{f'(x)}{f(x)} dx \\ \operatorname{atan}(f(x)) &= \int \frac{f'(x)}{1+f(x)^2} dx \\ \operatorname{asin}(f(x)) &= \int \frac{f'(x)}{(1-f(x)^2)^{1/2}} dx \\ \operatorname{acos}(f(x)) &= - \int \frac{f'(x)}{(1-f(x)^2)^{1/2}} dx,\end{aligned}$$

In each case, the corresponding scalar function of $[x^0]f$ should be added as the constant of integration. Finally, the exponential function and trigonometric functions can be computed in $O(M(n))$ operations from their inverse functions using Newton iteration.

Recent research has focused on speeding special compositions by constant factors by eliminating redundancy from Newton iteration [8, 45, 48]. Algorithms with quasilinear complexity are also known for certain right-compositions, including right-composition by algebraic functions [101] and some special functions [20]. Improved composition algorithms over special rings have also been investigated [6, 59, 93]. However, the algorithms of Brent and Kung have remained the best known in the general case.

4.5 Fast reversion without Newton iteration

Reversion of power series is the problem of finding the compositional inverse series of a given power series. Just as composition of formal power series (locally) represents composition of functions, *reversion of power series* represents inversion of functions. For example, if

$$f(x) = \exp(x) - 1 = x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

and

$$g(x) = \log(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots,$$

then f and g are reversions of each other as formal power series, and analytic inverse functions of each other in a neighborhood of $x = 0$.

We now consider computing the reversion of a general power series f , i.e. not having special form. Any algorithm for composition can be used for reversion (and vice versa) via Newton iteration, with at most a constant factor slowdown, as observed by Brent and Kung. In particular, combining Newton iteration with algorithms BK 2.1 or BK 2.2 from [25] gives two algorithms for reversion with respective complexity $O(n^{1/2}(M(n) + MM(n^{1/2})))$ and $O((n \log n)^{1/2}M(n))$.

Here we give an algorithm for reversion analogous to BK 2.1 and likewise requiring $O(n^{1/2}(M(n) + MM(n^{1/2})))$ operations, but achieving a constant factor speedup. The speedup ratio depends on the asymptotics of $M(n)$ and $MM(n)$ and is in the range between 1.2 and 2.6 for polynomial and matrix multiplication algorithms used in practice. Our algorithm also allows incorporating the complexity refinement of Huang and Pan.

Whereas BK 2.1 can be viewed as a baby-step giant-step version of Horner's rule, our algorithm can be viewed as a baby-step giant-step version of the Lagrange inversion formula, avoiding Newton iteration entirely (apart from a single $O(M(n))$ reciprocal computation). It is somewhat surprising that such an algorithm has been overlooked until now, with all publications following Brent and Kung apparently having taken Newton iteration as the final word on the subject matter.

4.5.1 The algorithm

Our setting is the ring of truncated power series $R[[x]]/\langle x^n \rangle$ over a commutative coefficient ring R in which the integers $1, \dots, n-1$ are cancellable (i.e. nonzero and not zero divisors). For example, we may take $R = \mathbb{Z}$ or $R = \mathbb{Z}/p\mathbb{Z}$ with prime $p \geq n$. We recall the Lagrange inversion formula ([64], p. 527). If $f(x) = x/h(x)$ where $h(0)$ is a unit in R , then the compositional inverse or reversion $f^{-1}(x)$ satisfying $f(f^{-1}(x)) = f^{-1}(f(x)) = x$ exists and its coefficients are given by

$$[x^k]f^{-1}(x) = \frac{1}{k}[x^{k-1}]h(x)^k.$$

The straightforward way to evaluate n terms of $f^{-1}(x)$ with the Lagrange inversion formula is to compute $h(x)$ (this requires $O(M(n))$ operations with Newton iteration) and then compute the powers h^2, h^3, \dots successively, for a total of $(n + O(1))M(n)$ operations.

Our observation is that it is redundant to compute all the powers of h given that we only are interested in a single coefficient from each. To do better, we choose some $1 \leq m < n$ and precompute h, h^2, h^3, \dots, h^m . For $0 \leq k < n$, we can then write h^k as h^{i+j} where $0 \leq j < m$ and $i = lm$ for some $0 \leq l \leq \lceil n/m \rceil$, only requiring $h^m, h^{2m}, h^{3m}, \dots$ to be computed subsequently. Determining a single coefficient in $h^k = h^i h^j$ can then be done in $O(n)$ operations using the definition of the Cauchy product. Picking $m \approx n^{1/2}$

minimizes the number of polynomial multiplications required.

We give a detailed account of this procedure as Algorithm 4.5.1. We note that most of the polynomial arithmetic is done to length $n - 1$ rather than length n , as the initial coefficient always is zero.

Algorithm 4.5.1 Fast Lagrange inversion

Input: $f = a_1x + a_2x^2 + \dots + a_{n-1}x^{n-1}$ where $n > 1$ and a_1 is a unit in R

Output: $g = b_1x + \dots + b_{n-1}x^{n-1}$ such that $f(g(x)) = g(f(x)) = x \bmod x^n$

- 1: $m \leftarrow \lceil \sqrt{n-1} \rceil$
 - 2: $h \leftarrow x/f \bmod x^{n-1}$
 - 3: **for** $1 \leq i < m$ **do**
 - 4: $h^{i+1} \leftarrow h^i \times h \bmod x^{n-1}$
 - 5: $b_i \leftarrow \frac{1}{i}[x^{i-1}]h^i$
 - 6: $t \leftarrow h^m$
 - 7: **for** $i = m, 2m, 3m, \dots, lm < n$ **do**
 - 8: $b_i \leftarrow \frac{1}{i}[x^{i-1}]t$
 - 9: **for** $1 \leq j < m$ while $i + j < n$ **do**
 - 10: $b_{i+j} \leftarrow \frac{1}{i+j} \sum_{k=0}^{i+j-1} ([x^k]t) \cdot ([x^{i+j-k-1}]h^j)$
 - 11: $t \leftarrow t \times h^m \bmod x^{n-1}$
- return** $b_1 + b_2x + \dots + b_{n-1}x^{n-1}$
-

An improved version

Algorithm 4.5.1 clearly requires $O(n^{1/2}M(n) + n^2)$ operations in R , as many as BK 2.1 with classical matrix multiplication. We can improve the complexity by packing the inner loops into a single matrix product as shown in Algorithm 4.5.2. This allows us to exploit fast matrix multiplication.

In the description of Algorithm 4.5.2, the matrices are indexed from 1 and the pseudocode has been simplified by letting the exponents run out of bounds, using the convention that $[x^k]p = 0$ when $k < 0$ or $k \geq n - 1$. To see that the algorithm is correct, write $[x^{i_1+(i_2-1)m-1}]h^{i_1+(i_2-1)m}$ as

$$\sum_{j=0}^{i_1+(i_2-1)m-1} ([x^j] h^{i_1}) \left([x^{i_1+(i_2-1)m-1-j}] h^{(i_2-1)m} \right)$$

and shift the summation index to obtain

$$\sum_{j=m-i_1+1}^{i_2m} ([x^{i_1+j-m-1}] h^{i_1}) \left([x^{i_2m-j}] h^{(i_2-1)m} \right)$$

which is the inner product of the nonzero part of row i_1 in B with the nonzero part of row i_2 in A .

Algorithm 4.5.2 Fast Lagrange inversion, matrix version

Input: $f = a_1x + a_2x^2 + \dots + a_{n-1}x^{n-1}$ where $n > 1$ and a_1 is a unit in R

Output: $g = b_1x + \dots + b_{n-1}x^{n-1}$ such that $f(g(x)) = g(f(x)) = x \bmod x^n$

- 1: $m \leftarrow \lceil \sqrt{n-1} \rceil$
 - 2: $h \leftarrow x/f \bmod x^{n-1}$
 - 3: Assemble $m \times m^2$ matrices B and A from h, h^2, \dots, h^m and $h^m, h^{2m}, h^{3m}, \dots$
 - 4: **for** $1 \leq i \leq m$, $1 \leq j \leq m^2$ **do**
 - 5: $B_{i,j} \leftarrow [x^{i+j-m-1}] h^i$
 - 6: $A_{i,j} \leftarrow [x^{im-j}] h^{(i-1)m}$
 - 7: $C \leftarrow AB^T$
 - 8: **for** $1 \leq i < n$ **do**
 - 9: $b_i \leftarrow C_i/i$ (C_i is the i th entry of C read rowwise)
 - return** $b_1 + b_2x + \dots + b_{n-1}x^{n-1}$
-

The structure of the matrices is perhaps illustrated more clearly by an example. Taking $n = 8$ and $m = 3$, we need the coefficients of $1, x, \dots, x^6$ in powers of h . Letting h_i^k denote $[x^i]h^k$, the matrices become

$$A = \begin{pmatrix} h_2^0 & h_1^0 & h_0^0 & 0 & 0 & 0 & 0 & 0 & 0 \\ h_5^3 & h_4^3 & h_3^3 & h_2^3 & h_1^3 & h_0^3 & 0 & 0 & 0 \\ (h_8^6) & (h_7^6) & h_6^6 & h_5^6 & h_4^6 & h_3^6 & h_2^6 & h_1^6 & h_0^6 \end{pmatrix}$$
$$B = \begin{pmatrix} 0 & 0 & h_0^1 & h_1^1 & h_2^1 & h_3^1 & h_4^1 & h_5^1 & h_6^1 \\ 0 & h_0^2 & h_1^2 & h_2^2 & h_3^2 & h_4^2 & h_5^2 & h_6^2 & (h_7^2) \\ h_0^3 & h_1^3 & h_2^3 & h_3^3 & h_4^3 & h_5^3 & h_6^3 & (h_7^3) & (h_8^3) \end{pmatrix}$$

where entries in parentheses do not contribute to the final result and may be set to zero. In this example the coefficient of x^4 in h^5 is given by the fifth entry in C , namely $C_{2,2} = h_4^3h_0^2 + h_3^3h_1^2 + h_2^3h_2^2 + h_1^3h_3^2 + h_0^3h_4^2$.

4.5.2 Complexity analysis

We now study the complexity in some more detail. Let $m = \lceil \sqrt{n-1} \rceil$. Then Algorithm 4.5.2 involves at most:

1. $2m + O(1)$ polynomial multiplications, each with cost $M(n)$
2. One $(m \times m^2)$ times $(m^2 \times m)$ matrix multiplication
3. $O(n)$ additional operations

For comparison, BK 2.1 requires at most:

1. m polynomial multiplications, each with cost $M(n)$
2. One $(m \times m)$ times $(m \times m^2)$ matrix multiplication
3. m polynomial multiplications and additions, each with cost $M(n) + n$

Brent and Kung break the matrix multiplication into m products of $m \times m$ matrices, requiring $m\text{MM}(m)$ operations. We can do the same in Algorithm 4.5.2, writing the product as a length- m inner product of $m \times m$ matrices. The extra $O(n^{3/2})$ additions in this matrix operation do not affect the complexity, but it is interesting to note that they match the $O(n^{3/2})$ additions in the last polynomial stage of BK 2.1. To summarize, both Algorithm 4.5.2 and BK 2.1 require at most $(2n^{1/2} + O(1))M(n) + n^{1/2}\text{MM}(n^{1/2}) + O(n^{3/2})$ operations.

The primary drawback of our algorithm as well as BK 2.1 is the requirement to store $O(n^{3/2})$ temporary coefficients in memory, compared to $O(n \log n)$ for BK 2.2 and $O(n)$ for a naive implementation of Lagrange inversion.

Avoiding Newton iteration

In effect, we need the same number of operations to perform a length- n reversion with fast Lagrange inversion as to perform a length- n composition with BK 2.1. However, to perform a reversion with BK 2.1, we must employ Newton iteration. Using the update

$$g_{k+1}(x) = \frac{f(g_k(x)) - x}{f'(g_k(x))},$$

where the chain rule allows us to reuse the composition in the numerator for the denominator, this entails computing a sequence of compositions of length

$$l = 1, \dots, \lceil n/4 \rceil, \lceil n/2 \rceil, n,$$

plus a fixed number of polynomial multiplications of the same length at each stage. If c and r are such that a length- n composition takes $C(n) = cn^r$ operations, Newton iteration asymptotically takes

$$C(n) + C(n/2) + C(n/4) + \dots = cn^r \left(\frac{2^r}{2^r - 1} \right)$$

operations, ignoring additional polynomial multiplications. For example, with classical polynomial multiplication as the dominant cost ($r = 5/2$), the speedup given by the expression in parentheses is $\frac{4}{31}(8 + \sqrt{2}) \approx 1.214$. With FFT polynomial multiplication, and classical matrix multiplication as the dominant cost ($r = 2$), the speedup is $4/3$. We note that a more efficient form of the Newton iteration might exist, in which case the speedup would be smaller.

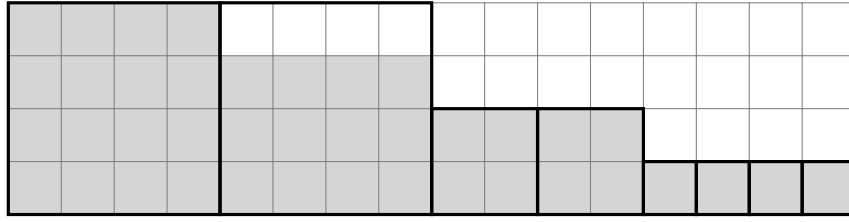


Figure 4.1: Square block decomposition of the matrix A .

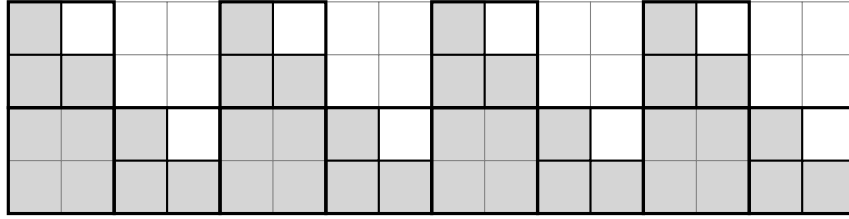


Figure 4.2: Square block decomposition of the matrix AP .

Improving the matrix multiplication

If the matrix multiplication dominates, we can gain an additional speedup from the fact that the i th $m \times m$ block of the matrix A only has $m - i + 1$ nonzero rows, whereas the matrices in BK 2.1 are full. Classically this gives a twofold speedup, reflected in the loop boundaries of Algorithm 4.5.1. We should ideally modify Algorithm 4.5.2 to include this saving.

In fact, a speedup is attainable with any square matrix multiplication algorithm having complexity $\text{MM}(m) \sim m^\omega$ where $\omega > 2$. For simplicity, assume that m is sufficiently composite. Do the first $m/2$ products as full products of size m , the next $(m/2 - m/3)$ in blocks of size $m/2$, the next $(m/3 - m/4)$ in blocks of size $m/3$, and so on (see Figure 4.1). At stage k , only k^2 products of blocks of size m/k are required. The speedup achieved through this procedure is

$$m^{\omega+1} \left(\sum_{k=1}^{\infty} \left(\frac{m}{k} - \frac{m}{k+1} \right) k^2 \left(\frac{m}{k} \right)^\omega \right)^{-1} > \left(\sum_{k=0}^{\infty} \frac{2^{k-1}}{2^{k\omega}} \right)^{-1} = 2 - 2^{2-\omega} > 1$$

where the nontrivial inequality can be obtained by considering the analogous subdivision with blocks of size $m/2^k$ only.

Alternatively, we can write $AB^T = (AP)(P^{-1}B^T)$ where P is a permutation matrix that makes each $m \times m$ block in A lower triangular, and use any algorithm that speeds up multiplication between a full and a triangular matrix. A simple recursive decomposition of size- k blocks into size- $k/2$ blocks (see Figure 4.2) has a proportional cost of $C(k) = 4C(k/2) + 2(k/2)^\omega + O(k^2)$, providing a speedup of $2^{\omega-1} - 2$. This is greater than 1 when

$\omega > \log_2 6 \approx 2.585$, and better than the first method when $\omega > 1 + \log_2(2 + \sqrt{2}) \approx 2.771$. In particular, we recover an optimal factor-two speedup with classical multiplication, and a $3/2$ speedup with the Strassen algorithm.

The asymptotic speedup ratios with both methods are illustrated in Figure 4.3.

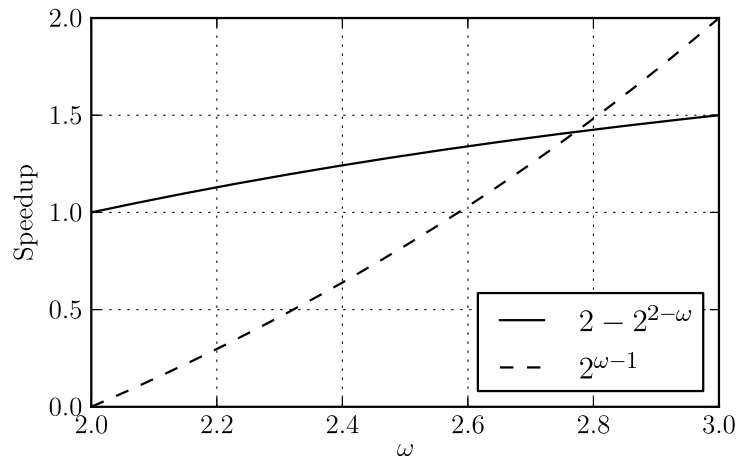


Figure 4.3: Speedup of the block decomposition algorithms for structured matrix multiplication.

Using rectangular multiplication

The complexity improvement of Huang and Pan for multiplication of rectangular matrices [51] also translates to Algorithm 4.5.2. More precisely, given any algorithm for $m \times m$ by $m \times n$ matrix multiplication over a general ring, there is a transposed version for $m \times n$ by $n \times m$ matrix multiplication that uses the same number of scalar multiplications [50] and a number of extra scalar additions bounded by the number of entries [90]. We can therefore take $\text{MM}(m, n, m) = (1 + o(1))\text{MM}(m, m, n)$.

With the Huang-Pan algorithm, we do not know whether a constant factor can be saved by exploiting the zero entries. This problem would be interesting to explore further. In any case, the Huang-Pan algorithm is currently only of theoretical interest, as the advantage probably only can be realized for infeasibly large matrices.

Practical performance

Table 4.1 gives a summary of the theoretical speedup gained by Algorithm 4.5.2 over BK 2.1 with various matrix and polynomial multiplication algorithms. With FFT-based polynomial multiplication, BK 2.2 is asymptotically faster than BK 2.1 and hence also

Dominant operation	Complexity	Newton	Matrix	Total
Polynomial, classical	$O(n^{5/2})$	1.214	1	1.214
Polynomial, Karatsuba	$O(n^{1/2+\log_2 3})$	1.308	1	1.308
Matrix, classical	$O(n^2)$	1.333	2.000	2.666
Matrix, classical, n -bit coeff.	$O(n^{3+o(1)})$	1.142	2.000	2.285
Matrix, Strassen	$O(n^{(1+\log_2 7)/2})$	1.364	1.500	2.047
Matrix, Cop.-Win.	$O(n^{1.688})$	1.450	1.229	1.782
Matrix, Huang-Pan	$O(n^{1.667})$	1.458	1?	1.458?
(Polynomial, FFT)*	$O(n^{3/2+o(1)})$	1.546	1	1.546
(Polynomial, FFT, n -bit coeff.)*	$O(n^{5/2+o(1)})$	1.214	1	1.214

Table 4.1: Theoretical speedup of Algorithm 2 over BK 2.1 due to avoiding Newton iteration and exploiting the matrix structure. *Assuming that matrix multiplication can be ignored.

than Algorithm 4.5.2. In practice, however, the overhead of quasilinear polynomial multiplication compared to matrix multiplication is likely to be large. Fast Lagrange inversion can therefore be expected to be faster than not only BK 2.1 but also BK 2.2 even for quite large n .

Of course, counting ring operations may not accurately reflect actual speed since operations in most interesting rings take a variable amount of time to execute on a physical computer. One consequence of this fact is that Newton iteration is likely to impose a smaller overhead than predicted, since coefficients generally are smaller in earlier steps than in later ones. Newton iteration can also be expected to perform better than generically when the output as a whole has small coefficients.

Over \mathbb{Z} in particular, arithmetic operations with b -bit numbers cost $O(b^{1+o(1)})$ bit operations. In power series arising in applications, we often have $b = O(n^{1\pm\epsilon})$. Two complexity estimates based on this assumption are included in Table 4.1. In practice, the speed will be sensitive to the sizes of the coefficients appearing internally in each algorithm, varying with the structure of $f(x)$.

We note that fast Lagrange inversion becomes faster than generically when the coefficients of $x/f(x)$ grow slowly. This is often the case when $f(x)$ is a rational function. The reversion of a rational function of fixed degree can be computed faster by a dedicated algorithm (Newton iteration takes $O(M(n))$ operations, using polynomial evaluation and series division for the composition), but it is desirable for a general-purpose algorithm to be efficient in this common case, and Lagrange inversion of course also works for nonrational functions having this growth property.

4.5.3 Benchmarks

We have implemented tuned versions of naive Lagrange inversion (“Lagrange”), BK 2.1 with Newton iteration, and Algorithm 4.5.1 (“Fast Lagrange”) over $\mathbb{Z}/p\mathbb{Z}$, \mathbb{Z} and \mathbb{Q} as part of the FLINT library. For each of these rings, FLINT provides fast coefficient and polynomial arithmetic. Matrix multiplication over $\mathbb{Z}/p\mathbb{Z}$ uses the Strassen algorithm when the smallest dimension is at least 256, which in principle helps BK 2.1 for $n > 256^2$ (the speedup is not significant for feasible n , however).

Timings over $\mathbb{Z}/p\mathbb{Z}$ obtained on an Intel Xeon E5-2650 2.0 GHz CPU with 256 GiB of RAM are given in Table 4.2. Algorithm 4.5.1 consistently runs about 1.6 times faster than BK 2.1, agreeing with a predicted speedup of 1.546 with quasilinear polynomial multiplication and negligible cost of matrix multiplication – we see that polynomial multiplication indeed dominates in BK 2.1 for n up to at least 10^6 . We have also implemented BK 2.2 over $\mathbb{Z}/p\mathbb{Z}$, finding it to take about twice as much time as BK 2.1 in the tested range. BK 2.2 might however be preferable for larger n due to memory limits (with $n = 10^6$ and 64-bit coefficients, BK 2.1 uses 15 GiB of memory and fast Lagrange reversion uses 7.5 GiB of memory).

Ring operations in \mathbb{Z} and \mathbb{Q} do not take constant time, as reflected in Tables 4.3 and 4.4. Timings are roughly cubic in n as expected from theory, but sensitive to the inputs. Fast Lagrange inversion is the fastest algorithm for small n in all examples, the fastest in all examples over \mathbb{Z} , and substantially faster for the rational functions f_3 and f_6 (in both cases $x/f(x)$ has small coefficients). For large n , BK 2.1 performs well on f_4 and f_5 , presumably due to generating smaller coefficients internally.

n	Lagrange	BK 2.1	BK 2.2	Fast Lagrange
10	10 μ s	10	18	6.1
10^2	2.8 ms	0.92	1.6	0.54
10^3	690 ms	66	120	45
10^4	110 s	3.3 (8%)	7.1	2.1
10^5	12100 s	144 (20%)	315	85
10^6	$1.9 \cdot 10^6$ s	8251 (28%)	15131	4832

Table 4.2: Time for reversion of a random power series over $\mathbb{Z}/p\mathbb{Z}$, $p = 2^{63} + 29$. The time spent on matrix multiplication in BK 2.1 is shown in parentheses.

With larger coefficients (as seen especially in the case of f_1), matrix multiplication appears to take a larger proportion of the time, suggesting that BK 2.2 becomes competitive for smaller n . We have not implemented BK 2.2 over \mathbb{Z} and \mathbb{Q} , however, and can therefore not provide a direct comparison.

n	Lagrange			BK 2.1			Fast Lagrange		
	f_1	f_2	f_3	f_1	f_2	f_3	f_1	f_2	f_3
10	7.0 μ s	6.5	6.3	10	10	10	4.9	4.3	4.1
10^2	31 ms	7.2	3.2	7.8	2.1	2.1	6.4	0.96	0.65
10^3	106 s	10	4.5	10	1.1	0.96	7.1	0.71	0.22
				(38%)	(31%)	(11%)			
10^4	-	-	-	24356 s	1453	538	8903	426	152
				(81%)	(67%)	(10%)			

Table 4.3: Time for the reversion of $f_1(x) = \sum_{k \geq 1} k!x^k$, $f_2(x) = \frac{x}{\sqrt{1-4x}}$, $f_3(x) = \frac{x+x^2}{1+x+x^2}$ over \mathbb{Z} .

n	Lagrange			BK 2.1			Fast Lagrange		
	f_4	f_5	f_6	f_4	f_5	f_6	f_4	f_5	f_6
10	15 μ s	15	13	31	28	28	11	11	9.1
10^2	40 ms	40	10	12	21	8.8	8.9	8.1	1.9
10^3	145 s	133	9.7	8.8	17	3.1	14	13	0.65
				(28%)	(24%)	(19%)			
10^4	-	-	-	13812 s	27057	1990	35633	27823	784
				(27%)	(27%)	(14%)			

Table 4.4: Time for the reversion of $f_4(x) = \exp(x) - 1$, $f_5(x) = x \exp(x)$, $f_6(x) = \frac{3x(1-x^2)}{2(1-x+x^2)^2}$ over \mathbb{Q} .

Care must be taken to handle denominators efficiently. In our implementation of BK 2.1, we found that naive matrix multiplication over \mathbb{Q} took ten times as long as polynomial multiplications. Clearing denominators and multiplying matrices over \mathbb{Z} resulted in a comparable time being spent on the matrix and polynomial stages. Similar concerns apply when implementing Algorithms 4.5.1 and 4.5.2. On the other hand, translating the *entire* composition or reversion to one over \mathbb{Z} by rescaling the inputs typically results in too much coefficient inflation, and can even run slower than a classical algorithm working directly over \mathbb{Q} . We expect the situation to be similar when working with parametric power series having rational functions as coefficients.

An interesting alternative would be to work modulo small primes and reconstruct the output using the Chinese remainder theorem. We have not investigated this option in detail. It would provide additional memory benefits: for example, if the coefficients have $O(n)$ bits, direct application of BK 2.1 or fast Lagrange reversion uses $O(n^{5/2})$ bits of temporary space, while modular reversions each require $O(n^{3/2+o(1)})$ bits of space – less than the $O(n^2)$ bits required to store the output.

4.5.4 Reversion with numerical coefficients

We have implemented the previously-mentioned algorithms for reversion of power series with real or complex coefficients (represented by balls) in the Arb library. Some experimentation suggests that fast Lagrange inversion generally performs about as well as BK 2.1 both in terms of speed and numerical accuracy, as the example in Table 4.5 and Table 4.6 illustrates.

n	b	Lagrange	BK 2.1	Fast Lagrange
10	100	0.0002	0.00017	0.00007
100	100	0.025	0.0093	0.0059
100	1000	0.10	0.030	0.018
1000	100	9.8	0.82	0.65
1000	1000	19	1.8	1.2
1000	10000	193	30	16
10000	100	-	110	130
10000	1000	-	226	200
10000	10000	-	2690	1349

Table 4.5: Time (in seconds) for the reversion of $\exp(x) - 1$ to order $O(x^n)$ over \mathbb{R} , using ball arithmetic with a precision of b bits.

Lagrange		BK 2.1		Fast Lagrange	
Error	Bound	Error	Bound	Error	Bound
2^{-100}	2^{-59}	2^{-99}	2^{-55}	2^{-99}	2^{-77}
2^{-96}	2^{+564}	2^{-96}	2^{+137}	2^{-96}	2^{-9}
2^{-996}	2^{-336}	2^{-996}	2^{-833}	2^{-996}	2^{-909}
2^{-91}	2^{+9325}	2^{-96}	2^{+3929}	2^{-91}	2^{+227}
2^{-991}	2^{+8425}	2^{-994}	2^{-642}	2^{-991}	2^{-673}
2^{-9991}	2^{-575}	2^{-9994}	2^{-9642}	2^{-9991}	2^{-9673}
-	-	2^{-92}	2^{+72552}	2^{-22}	2^{+1167}
-	-	2^{-990}	2^{-337}	2^{-921}	2^{+209}
-	-	2^{-9990}	2^{-9337}	2^{-9921}	2^{-8791}

Table 4.6: For each corresponding entry in Table 4.5, this table shows the actual error of the floating-point approximation computed for the last coefficient in the output, as well as the error bound computed using ball arithmetic.

On this test problem, we observe that all algorithms are numerically stable in floating-point arithmetic, giving midpoints that are accurate to high precision. The error bounds

produced with our implementation of ball arithmetic are far from tight, however. The poor error bounds are explained by the fact that a large number of successive polynomial multiplications are being performed, together with the characteristics of our multiplication algorithm. We compute a slight overestimate for the error bound in each multiplication, and the successive multiplications cause these overestimates to accumulate.

Since both BK 2.1 and fast Lagrange inversion reduce the number of successive multiplications from $O(n)$ to $O(n^{1/2})$, they are better than naive Lagrange inversion in terms of both execution speed and the tightness of the error bounds (for composition of power series, BK 2.1 represents a similar improvement over Horner's rule).

In either BK 2.1 or fast Lagrange inversion, the error bounds could potentially be improved by computing powers from scratch using binary exponentiation rather than recurrently. This would reduce the maximum number of sequential multiplications to $O(\log n)$, at the cost of increasing the total number of polynomial multiplications by a factor $O(\log n)$.

4.5.5 Remarks

Fast Lagrange inversion is a practical algorithm for reversion of formal power series, having essentially no higher overhead than a naive implementation of Lagrange inversion for small n and requiring fewer coefficient operations than Newton iteration coupled with BK 2.1 for large n . Among currently available general-purpose algorithms, it is likely to be the fastest choice for typical coefficient rings, input series, and values of n , and may thus be a good choice as a default reversion algorithm in a computer algebra system. Newton iteration with BK 2.2 remains faster asymptotically when FFT polynomial multiplication is available, and uses less memory, but may require very large n to become advantageous.

An interesting question is whether a reversion analog of BK 2.2 can be constructed that avoids Newton iteration, or whether we can otherwise save constant factors in BK 2.2. We may also ask whether the close correspondence in complexity between Algorithm 4.5.2 and BK 2.1 can be explained by some underlying duality along the lines of the transposition principle. Further investigation of improvements over particular rings and of the special matrix multiplications arising in Algorithm 4.5.2 and BK 2.1 would also be warranted.

4.6 The gamma function

The gamma function $\Gamma(z)$ appears in many contexts, for example in normalization factors of hypergeometric functions and in the functional equation of the Riemann zeta function (and more general L -functions). Thus the ability to compute series expansions of the gamma function is often useful.

Recall that to evaluate the gamma function with an accuracy of P bits, we need a rising factorial of $r \sim P$ factors, and $N \sim P$ terms of the Stirling series

$$\log \Gamma(z) = \left(z - \frac{1}{2}\right) \log z - z + \frac{\ln 2\pi}{2} + \sum_{k=1}^{N-1} \frac{B_{2k}}{2k(2k-1)z^{2k-1}} + R_N(z). \quad (4.6.1)$$

Assume now that we want to compute D coefficients in the Taylor series of the gamma function at a fixed complex number z . This can be done by substituting $z \rightarrow z + x \in \mathbb{C}[[x]]/\langle x^D \rangle$ and evaluating all expressions in (4.6.1) using power series arithmetic. Note that we can evaluate $\Gamma(Z)$ for any formal power series $Z = z + z_1x + z_2x^2 + \dots$ by first evaluating $\Gamma(z + x)$ and then formally right-composing by $Z - z$ using, for instance, the Brent-Kung algorithm (unlike the elementary functions, we are not aware of a better way to evaluate the gamma function of a whole power series at once).

The remainder in the Stirling series is given exactly (Olver [84], pp. 293–295) by

$$R_N(z) = \int_0^\infty \frac{B_{2N} - \tilde{B}_{2N}(t)}{2N(z+t)^{2N}} dt.$$

where $\tilde{B}_n(t) = B_n(t - [t])$ denotes the n -th periodic Bernoulli polynomial. Since $R_N(z)$ is an analytic function of $z \in \mathbb{C} \setminus (-\infty, +\epsilon)$, we can expand it in a power series and bound the coefficients to get error bounds for derivatives of all orders. By extending Olver's calculation, we obtain the coefficientwise inequality

$$|R_N(z + x)| \leq \sum_{k=0}^{\infty} b_k x^k \in \mathbb{R}[[x]]$$

where

$$b_k = 2|B_{2N}| \frac{\Gamma(2N + k - 1)}{\Gamma(k + 1)\Gamma(2N + 1)} |z| (A(z)/|z|)^{2N+k}$$

and $A(z) = 1/\cos(\arg(z)/2)$.

If we use recurrence relations to generate and add the terms in the Stirling series and in the rising factorial one by one, computing D derivatives to a precision of P bits takes $O^\sim(P^2D)$ time. When D is small (say $D \leq 4$) the rising factorial can be accelerated using the rectangular splitting algorithm of section 3.5 (applied to a power series).

Now assume that $D \sim P$. In this case, we can improve the complexity of computing $\Gamma(z + x)$ to $O^\sim(PD)$, which is softly optimal. This follows by applying binary splitting to both the rising factorial and the tail sum in the Stirling series.

For the rising factorial $(z+x)(z+x+1)\cdots(z+x+r-1)$, the binary splitting algorithm simply amounts to recursively dividing the product in half. Although the bit sizes of the coefficients are the same ($O^\sim(P)$ bits) throughout, binary splitting balances the polynomial degrees, which is sufficient to achieve the complexity reduction.

The binary splitting scheme for the tail sum can be derived by writing down a matrix recurrence for the sequence of partial sums. We observe that the recurrence is linear but not holonomic due to the presence of Bernoulli numbers. As with the rising factorial, the complexity reduction comes from the fact that the Bernoulli numbers have bounded size, where the relevant measure of size is the polynomial degree.

Algorithm 4.6.1 Tail of the Stirling series using binary splitting

Input: $z \in \mathbb{C}$ and $N \in \mathbb{Z}_{>2}, D \in \mathbb{Z}_{\geq 1}$

Output: $\sum_{k=1}^{N-1} \frac{B_{2k}}{2k(2k-1)(z+x)^{2k-1}} \in \mathbb{C}[[x]]/\langle x^D \rangle$

```

1: Let  $x$  denote the generator of  $\mathbb{C}[[x]]/\langle x^D \rangle$ 
2: function BINSPLIT( $j, k$ )
3:   if  $j+1 = k$  then
4:     if  $j = 1$  then
5:       return  $(B_{2j}/(2j(2j-1)), z+x)$ 
6:     else
7:       return  $(B_{2j}/(2j(2j-1)), (z+x)^2)$ 
8:   else
9:      $(T_1, Q_1) \leftarrow$  BINSPLIT( $j, \lfloor (j+k)/2 \rfloor$ )
10:     $(T_2, Q_2) \leftarrow$  BINSPLIT( $\lfloor (j+k)/2 \rfloor, k$ )
11:    return  $(T_1Q_2 + T_2, Q_1Q_2)$  ▷ Polynomial multiplications mod  $x^D$ 
12:  $(T, Q) \leftarrow$  BINSPLIT(1,  $N$ )
13: return  $T/Q$  ▷ Power series division mod  $x^D$ 

```

In some more detail, let the terms be $t_k = c_k/(u_1u_2\dots u_k)$ where $c_k = B_{2k}/(2k(2k-1))$ and $u_1 = z+x$ and $u_k = (z+x)^2$ for $k \geq 2$. Then $t_k = t_{k-1}/u_k$, and the partial sums s_k satisfy $s_k = s_{k-1} + (c_k/u_k)t_k$. In matrix form, this becomes

$$\begin{pmatrix} s_k \\ t_k \end{pmatrix} = \frac{1}{u_k} \begin{pmatrix} u_k & c_k \\ 0 & 1 \end{pmatrix} \begin{pmatrix} s_{k-1} \\ t_{k-1} \end{pmatrix}.$$

Denoting the top left entry in the matrix (and the denominator) by Q and the top right entry by T , we obtain algorithm 4.6.1.

Timings of our implementation within the Arb library are shown in Table 4.7. We observe that the complexity really is quasioptimal in practice, as increasing D by a factor 10 roughly increases the running time by a factor 100. The series expansion of the gamma function in Mathematica, included for scale, seems to be poorly implemented.

D	Mathematica 9.0	Arb
10	0.568	0.00187
30	60.9	0.00346
100	(crashes)	0.023
300		0.178
1000		2.944
3000		37.9
10000		535

Table 4.7: Time in seconds to compute the Taylor series to order $O(x^D)$ of $\Gamma(z+x)$ at $z = \sqrt{2} + \sqrt{3}i$, with a precision of D digits. Computations were done on a 64-bit Intel Xeon X5675 3.07 GHz CPU.

4.7 Integer zeta values

Integer arguments of the gamma function are particularly common. Since

$$\log \Gamma(1+x) = -\gamma x + \sum_{n=2}^{\infty} \frac{(-1)^n \zeta(n)}{n} x^n,$$

computing the Taylor series expansion of the gamma function at an integer is equivalent modulo a power series exponential or logarithm to multi-evaluation of the Riemann zeta function $\zeta(s)$ at positive integer values of s (Euler’s constant $\gamma = \lim_{s \rightarrow 1} \zeta(s) - 1/(s-1)$ can also be computed softly optimally to high precision using binary splitting, by the Brent-McMillan algorithm [26, 24]).

It is sufficient to consider $s \lesssim p$ where p is the numerical precision in bits, as $\zeta(s)$ can be computed easily from the defining series $\zeta(s) = \sum_{k=1}^{\infty} k^{-s}$ or the corresponding Euler product when s is about as large as p . Moreover, we mainly need to worry about odd s , as the even zeta values when $s \ll p$ can be computed quickly from Bernoulli numbers.

Borwein, Bradley and Crandall [15] propose several methods for multi-evaluation of $\zeta(s)$ or Bernoulli numbers, including “hyperbolic series” of Ramanujan and Zagier, and methods based on power series. One method uses the identity

$$\sum_{k=0}^{\infty} \zeta(2k+3)x^k = \sum_{k=1}^{\infty} \frac{1}{k^3(1-x/k^2)} \quad (4.7.1)$$

$$= \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k^3 \binom{2k}{k}} \left(\frac{1}{2} + \frac{2}{1-x/k^2} \right) \prod_{j=1}^{k-1} \left(1 - \frac{x}{j^2} \right) \quad (4.7.2)$$

Each term of the last sum adds roughly two bits of accuracy, for every derivative (it is easy to derive an explicit error bound). Borwein, Bradley and Crandall suggest matching coefficients in (4.7.2) to obtain fast schemes for evaluating single values of $\zeta(s)$, but we

can also apply polynomial binary splitting directly. If we write (4.7.2) as $\sum_{k=1}^{\infty} T(k)$, the terms are hypergeometric, satisfying $T(k+1) = P(k)T(k)/Q(k)$ with

$$\begin{aligned} P(k) &= k(k^2 - x)^2(x - 5(k+1)^2) \\ Q(k) &= 2(k+1)^2(2k+1)((k+1)^2 - x)(5k^2 - x) \\ T(1) &= \frac{x-5}{4(x-1)}. \end{aligned}$$

Adding the terms using binary splitting, we get an algorithm that is softly optimal both when we want a large number of zeta values to high precision and (unlike the Stirling series) when we just want a few zeta values to high precision.

A second method is a generalization of the binary splitting algorithm to evaluate $\Gamma(z)$ quickly when z is a rational number, which we already mentioned in 3.6.2. With power series notation, we have

$$\Gamma(z+x) = \int_0^{\infty} e^{-t} t^{z+x-1} dt \approx \int_N^{\infty} e^{-t} t^{z+x-1} dt \approx N^{z+x} H \quad (4.7.3)$$

where

$$H = \sum_{k=0}^R \frac{(-1)^k N^k}{(k+z+x)k!}. \quad (4.7.4)$$

This generalization is due to Karatsuba [57, 58], who used it to prove that the Hurwitz zeta function $\zeta(s, z)$ can be computed asymptotically fast for fixed positive integers s and algebraic z . Karatsuba also proves explicit error bounds: for a precision of p bits, it is sufficient to choose $N \approx p \log 2$ and $R \approx 4N$ (this can be improved to $R \approx cN$ where $c = 1/W(1/e) \approx 3.59$ and $W(x)$ is the Lambert W -function).

We can apply polynomial arithmetic directly to the sum H . Alternatively, we can expand (4.7.4) as a geometric series with respect to x to obtain D scalar sums to be evaluated separately using binary splitting:

$$H = \sum_{j=0}^{D-1} x^j \sum_{k=0}^R \frac{(-1)^{k+j} N^k}{(k+z)^{j+1} k!}. \quad (4.7.5)$$

Either way, we obtain H which we finally multiply by the power series

$$N^{z+x} = \sum_{k=0}^{D-1} N^z (\log N)^k x^k$$

to obtain the Taylor series for the gamma function.

The scheme using polynomial binary splitting is equivalent to Algorithm 5 given by Borwein, Bradley and Crandall, while the scheme using scalar binary splitting is equivalent to the algorithm described by Karatsuba. Borwein, Bradley and Crandall ask whether “[Karatsuba’s] methods may be used to accelerate even further the series computations

of Algorithm 5”, apparently overlooking the correspondence between (4.7.4) and (4.7.5). In fact, the version given by Borwein, Bradley and Crandall is faster for large D since it balances both the polynomial degrees and the bit sizes of the terms, while the factors $(k+z)^{j+1}$ in Karatsuba’s algorithm become progressively more expensive as j increases. Karatsuba also expresses the final multiplication by N^{z+x} as an explicit convolution sum, which is less efficient to evaluate verbatim.

We finally also mention Borwein’s remarkably simple and efficient formula [17]

$$(1 - 2^{1-s})\zeta(s) = \frac{1}{d_n} \sum_{k=0}^{n-1} \frac{(-1)^k (d_n - d_k)}{(k+1)^s} + \varepsilon_n(s) \quad (4.7.6)$$

where

$$d_k = n \sum_{i=0}^k \frac{(n+i-1)! 4^i}{(n-i)! (2i)!} \quad (4.7.7)$$

and where $|\varepsilon_n(s)| \leq 3/(3 + \sqrt{8})^n$ for real $s \geq 2$. This method can also be used for complex s , but the error bound deteriorates with larger imaginary parts. When s is a fixed positive integer, the partial sums are a holonomic sequence of order 3, so applying binary splitting to (4.7.6)–(4.7.7) then allows computing $\zeta(s)$ with softly optimal complexity (this fact was already pointed out in [43]). A drawback of this binary splitting method is that it rapidly gets slower for larger s , due to the appearance of factors $\prod_k k^s$ in the denominators, and it does not appear possible to save a significant amount of time when evaluating $\zeta(s)$ for several s simultaneously. Nevertheless, the method is of interest for computing isolated values of $\zeta(s)$ to millions of digits.

Algorithm 4.7.1 is an extension of the sequential implementation of Borwein’s formula used in the MPFR library [78]. We have omitted the computation of error bounds for simplicity. Notably, the coefficients $\tilde{d}_k = d_n - d_k$ are computed sequentially (in the end we have $\tilde{d}_0 = d_n$) and the main part of the computation only uses integer arithmetic. We have extended the algorithm by adding the inner loop which computes \tilde{d}_{k-1}/k^s for several s in progression. The algorithm has complexity $O^\sim(p^2)$ for each zeta value whether one calls the function repeatedly with $N = 1$ and different s or computes several values at once with $N > 1$, but the latter is much faster since the divisor $u = k^h$ is small (usually a single machine word) and t moreover shrinks with each iteration (whereas computing the power on line 6 gets more expensive for larger s).

A benchmark comparison of various algorithms as implemented in the Arb library is shown in Table 4.8. Borwein’s algorithm with power recycling is the fastest method up to precisions exceeding 10000 digits. The Stirling series holds up well overall, especially when computing a large number of derivatives to high precision. As expected, the two other power series-based algorithms are superior when computing a small number of gamma function derivatives (equivalently, zeta values) at very high precision. The series (4.7.2) seems consistently faster than (4.7.4). The latter also gives us the even zeta

Algorithm 4.7.1 Evaluation of $\zeta(s + hj), j = 0 \dots N - 1, s, h \in \mathbb{N}, s \geq 2$

```

1: function ZETABORWEINRECYCLED( $s, h, N, p$ )
2:    $n \leftarrow \lceil 2 + p \log(2) / \log(3 + \sqrt{8}) \rceil$             $\triangleright$  Accuracy of approximately  $p$  bits
3:    $c \leftarrow d \leftarrow 2^{2n-1}$                             $\triangleright c, d, z_i, t, u$  will be integer variables
4:    $z_0, z_1, \dots, z_{N-1} \leftarrow 0$ 
5:   for  $k \leftarrow n$  to 1 do
6:      $t \leftarrow (-1)^{k+1} \lfloor d / k^s \rfloor$             $\triangleright$  Approximate division, creating rounding error
7:      $z_0 \leftarrow z_0 + t$ 
8:      $u \leftarrow k^h$ 
9:     for  $j \leftarrow 1$  to  $N - 1$  do
10:       $t \leftarrow \lfloor t / u \rfloor$                     $\triangleright$  Approximate division, creating rounding error
11:       $z_j \leftarrow z_j + t$ 
12:      $c \leftarrow c \cdot k(2k - 1)$ 
13:      $c \leftarrow c / (2(n - k + 1)(n + k - 1))$         $\triangleright$  Exact integer division
14:      $d \leftarrow d + c$ 
15:   for  $j \leftarrow 0$  to  $N - 1$  do
16:      $w_j \leftarrow z_j / (d(1 - 2^{1-(s+hj)}))$           $\triangleright$  Approximate division ( $w_j$  is real)
17:   return  $w_0, w_1, \dots, w_{N-1}$ 

```

p	n	ST	S1	S2	B
100	10	0.00075	0.0015	0.0072	0.000089
100	100	0.0047	0.013	0.026	0.00046
1000	10	0.012	0.023	0.14	0.0031
1000	100	0.077	0.29	0.91	0.024
1000	1000	0.57	1.3	4.6	0.071
10000	10	1.5	0.6	3.44	0.22
10000	100	4.61	9.1	22	2.0
10000	1000	16	62	163	17
10000	10000	95	224	727	34
100000	10	410	17	84	23
100000	100	530	213	535	196
100000	1000	1108	1965	5082	2124
100000	10000	3196	12129	36029	16818

Table 4.8: Time in seconds to compute either the Taylor series of $\Gamma(1 + x)$ to order $O(x^n)$, or to compute the odd zeta values $\zeta(2k + 1)$ with $2k + 1 < n$, at a precision of p digits. ST: the Stirling series with binary splitting; S1: the series (4.7.2) with polynomial binary splitting; S2: the series (4.7.4) with polynomial binary splitting; B: Borwein's algorithm with power recycling. Computations were done on a 64-bit Intel Xeon X5675 3.07 GHz CPU.

values, but this is not a large advantage since we can compute those faster separately. An advantage of (4.7.4) is that we also can use it at rational or algebraic $z \neq 1$.

4.8 The Hurwitz zeta function

The Hurwitz zeta function $\zeta(s, a)$ is defined for complex numbers s and a by analytic continuation of the sum

$$\zeta(s, a) = \sum_{k=0}^{\infty} \frac{1}{(k+a)^s}.$$

The usual Riemann zeta function is given by $\zeta(s) = \zeta(s, 1)$.

In this section, we consider numerical computation of $\zeta(s, a)$ by the Euler-Maclaurin formula with rigorous error control. Error bounds for $\zeta(s)$ are classical (see for example [36, 15] and numerous references therein), but previous works have restricted to the case $a = 1$ or have not considered derivatives. Our main contribution is to give an efficiently computable error bound for $\zeta(s, a)$ valid for any complex s and a and for an arbitrary number of derivatives with respect to s (equivalently, we allow s to be a formal power series).

We also discuss implementation aspects, such as parallelization and use of fast polynomial arithmetic. An implementation of $\zeta(s, a)$ based on the algorithms described here is available in the Arb library. In the end of this section, we present results from some new record computations done with this implementation.

Our interest is in evaluating $\zeta(s, a)$ to high precision (hundreds or thousands of digits) for a single s of moderate height, say with imaginary part less than 10^6 . Investigations of zeros of large height typically use methods based on the Riemann-Siegel formula and fast multi-evaluation techniques such as the Odlyzko-Schönhage algorithm [82] or the recent algorithm of Hiary [49].

This work is motivated by several applications. For example, recent work of Matiyasevich and Beliakov required values of thousands of nontrivial zeros ρ_n of $\zeta(s)$ to a precision of several thousand digits [73, 74]. Investigations of quantities such as the Stieltjes constants $\gamma_n(a)$ and the Keiper-Li coefficients λ_n also call for high-precision values [60, 67]. The difficulty is not necessarily that the final result needs to be known to very high accuracy, but that intermediate calculations may involve catastrophic cancellation.

More broadly, the Riemann and Hurwitz zeta functions are useful for numerical evaluation of various other special functions such as polygamma functions, polylogarithms, Dirichlet L -functions, generalized hypergeometric functions at singularities [11], and certain number-theoretical constants [39]. High-precision numerical values are of particular interest for guessing algebraic relations among special values of such functions (which

subsequently may be proved rigorously by other means) or ruling out the existence of algebraic relations with small norm [5].

4.8.1 Evaluation using the Euler-Maclaurin formula

Assume that f is analytic on a domain containing $[N, U]$ where $N, U \in \mathbb{Z}$, and let M be a positive integer. Let B_n denote the n -th Bernoulli number and let $\tilde{B}_n(t) = B_n(t - [t])$ denote the n -th periodic Bernoulli polynomial. The Euler-Maclaurin summation formula (described in numerous works, such as [84]) states that

$$\sum_{k=N}^U f(k) = I + T + R \quad (4.8.1)$$

where

$$I = \int_N^U f(t) dt, \quad (4.8.2)$$

$$T = \frac{1}{2} (f(N) + f(U)) + \sum_{k=1}^M \frac{B_{2k}}{(2k)!} \left(f^{(2k-1)}(U) - f^{(2k-1)}(N) \right), \quad (4.8.3)$$

$$R = - \int_N^U \frac{\tilde{B}_{2M}(t)}{(2M)!} f^{(2M)}(t) dt. \quad (4.8.4)$$

If f decreases sufficiently rapidly, (4.8.1)–(4.8.4) remain valid after letting $U \rightarrow \infty$. To evaluate the Hurwitz zeta function, we set

$$f(k) = \frac{1}{(a+k)^s} = \exp(-s \log(a+k))$$

with the conventional logarithm branch cut on $(-\infty, 0)$. The derivatives of $f(k)$ are given by

$$f^{(r)}(k) = \frac{(-1)^r (s)_r}{(a+k)^{s+r}}$$

where $(s)_r = s(s+1) \cdots (s+r-1)$ denotes a rising factorial. The Euler-Maclaurin summation formula now gives, at least for $\Re(s) > 1$ and $a \neq 0, -1, -2, \dots$,

$$\zeta(s, a) = \sum_{k=0}^{N-1} f(k) + \sum_{k=N}^{\infty} f(k) = S + I + T + R \quad (4.8.5)$$

where

$$S = \sum_{k=0}^{N-1} \frac{1}{(a+k)^s}, \quad (4.8.6)$$

$$I = \int_N^{\infty} \frac{1}{(a+t)^s} dt = \frac{(a+N)^{1-s}}{s-1}, \quad (4.8.7)$$

$$T = \frac{1}{(a+N)^s} \left(\frac{1}{2} + \sum_{k=1}^M \frac{B_{2k}}{(2k)!} \frac{(s)_{2k-1}}{(a+N)^{2k-1}} \right), \quad (4.8.8)$$

$$R = - \int_N^{\infty} \frac{\tilde{B}_{2M}(t)}{(2M)!} \frac{(s)_{2M}}{(a+t)^{s+2M}} dt. \quad (4.8.9)$$

If we choose N and M such that $\Re(a + N) > 0$ and $\Re(s + 2M - 1) > 0$, the integrals in I and R are well-defined, giving us the analytic continuation of $\zeta(s, a)$ to $s \in \mathbb{C}$ except for the pole at $s = 1$.

To evaluate derivatives with respect to s of $\zeta(s, a)$, we substitute $s \rightarrow s + x \in \mathbb{C}[[x]]$ and evaluate (4.8.5)–(4.8.9) with the corresponding arithmetic operations done on formal power series (which may be truncated at some arbitrary finite order in an implementation). For example, the summand in (4.8.6) becomes

$$\frac{1}{(a+k)^{s+x}} = \sum_{i=0}^{\infty} \frac{(-1)^i \log(a+k)^i}{(a+k)^s} x^i \in \mathbb{C}[[x]]. \quad (4.8.10)$$

Note that we can evaluate $\zeta(S, a)$ for any formal power series $S = s + s_1x + s_2x^2 + \dots$ by first evaluating $\zeta(s+x, a)$ and then formally right-composing by $S - s$. We can also easily evaluate derivatives of $\zeta(s, a) - 1/(s-1)$ at $s = 1$. The pole of $\zeta(s, a)$ only appears in the term I on the right hand side of (4.8.5), so we can remove the singularity as

$$\begin{aligned} \lim_{s \rightarrow 1} \left[I - \frac{1}{(s+x)-1} = \frac{(a+N)^{1-(s+x)}}{(s+x)-1} - \frac{1}{(s+x)-1} \right] \\ = \sum_{i=0}^{\infty} \frac{(-1)^{i+1} \log(a+N)^{i+1}}{i!} x^i \in \mathbb{C}[[x]]. \end{aligned} \quad (4.8.11)$$

We now wish to bound the coefficientwise error $|R(s+x)|$ (using the notation in section 4.2) where $R(s) = R$ is the remainder integral given in (4.8.9).

To express the error bound in a compact form, we introduce the sequence of integrals defined for integers $k \geq 0$ and real parameters $A > 0, B > 1, C \geq 0$ by

$$J_k(A, B, C) \equiv \int_A^{\infty} t^{-B} (C + \log t)^k dt.$$

Using the binomial theorem, $J_k(A, B, C)$ can be evaluated in closed form for any fixed k . In fact, collecting factors gives

$$J_k(A, B, C) = \frac{L_k}{(B-1)^{k+1} A^{B-1}}$$

where $L_0 = 1$, $L_k = kL_{k-1} + D^k$ and $D = (B-1)(C + \log A)$. This recurrence allows computing J_0, J_1, \dots, J_n easily, using $O(n)$ arithmetic operations.

Theorem 4.8.1. *Given complex numbers $s = \sigma + \tau i$, $a = \alpha + \beta i$ and positive integers N, M such that $\alpha + N > 1$ and $\sigma + 2M > 1$, the error term (4.8.9) in the Euler-Maclaurin summation formula applied to $\zeta(s+x, a) \in \mathbb{C}[[x]]$ satisfies*

$$|R(s+x)| \leq \frac{4|(s+x)_{2M}|}{(2\pi)^{2M}} \left| \sum_{k=0}^{\infty} R_k x^k \right| \in \mathbb{R}[[x]] \quad (4.8.12)$$

where $R_k \leq (K/k!) J_k(N + \alpha, \sigma + 2M, C)$, with

$$C = \frac{1}{2} \log \left(1 + \frac{\beta^2}{(\alpha + N)^2} \right) + \operatorname{atan} \left(\frac{|\beta|}{\alpha + N} \right) \quad (4.8.13)$$

and

$$K = \exp \left(\max \left(0, \tau \operatorname{atan} \left(\frac{\beta}{\alpha + N} \right) \right) \right). \quad (4.8.14)$$

Proof. We have

$$\begin{aligned} |R(s+x)| &= \left| \int_N^\infty \frac{\tilde{B}_{2M}(t)}{(2M)!} \frac{(s+x)_{2M}}{(a+t)^{s+x+2M}} dt \right| \\ &\leq \int_N^\infty \left| \frac{\tilde{B}_{2M}(t)}{(2M)!} \frac{(s+x)_{2M}}{(a+t)^{s+x+2M}} \right| dt \\ &\leq \frac{4|(s+x)_{2M}|}{(2\pi)^{2M}} \int_N^\infty \left| \frac{1}{(a+t)^{s+x+2M}} \right| dt \end{aligned}$$

where the last step invokes the fact that

$$|\tilde{B}_{2M}(t)| < \frac{4(2M)!}{(2\pi)^{2M}}.$$

Thus it remains to bound the coefficients R_k satisfying

$$\int_N^\infty \left| \frac{1}{(a+t)^{s+x+2M}} \right| dt = \sum_k R_k x^k, \quad R_k = \int_N^\infty \frac{1}{k!} \left| \frac{\log(a+t)^k}{(a+t)^{s+2M}} \right| dt.$$

By the assumption that $\alpha + t \geq \alpha + N \geq 1$, we have

$$\begin{aligned} |\log(\alpha + \beta i + t)| &= \left| \log(\alpha + t) + \log \left(1 + \frac{\beta i}{\alpha + t} \right) \right| \\ &\leq \log(\alpha + t) + \left| \log \left(1 + \frac{\beta i}{\alpha + t} \right) \right| \\ &= \log(\alpha + t) + \left| \frac{1}{2} \log \left(1 + \frac{\beta^2}{(\alpha + t)^2} \right) + i \operatorname{atan} \left(\frac{\beta}{\alpha + t} \right) \right| \\ &\leq \log(\alpha + t) + C \end{aligned}$$

where C is defined as in (4.8.13). By the assumption that $\sigma + 2M > 1$, we have

$$\frac{1}{|(\alpha + \beta i + t)^{\sigma + \tau i + 2M}|} = \frac{\exp(\tau \arg(\alpha + \beta i + t))}{|\alpha + \beta i + t|^{\sigma + 2M}} \leq \frac{K}{(\alpha + t)^{\sigma + 2M}}$$

where K is defined as in (4.8.14). Bounding the integrand in R_k in terms of the integrand in the definition of J_k now gives the result. \square

The bound given in Theorem 4.8.1 should generally approximate the exact remainder (4.8.9) quite well, even for derivatives of large order, if $|a|$ is not too large. The quantity K is especially crude, however, as it does not decrease when $|a + t|^{-\tau i}$ decreases

exponentially as a function of τ . We have made this simplification in order to obtain a bound that is easy to evaluate for all s, a . In fact, assuming that a is small, we can simplify the bounds a bit further using

$$C \leq \frac{\beta^2}{2(\alpha + N)^2} + \frac{|\beta|}{(\alpha + N)}.$$

In practice, the Hurwitz zeta function is usually only considered for $0 < a \leq 1$, unless s is an integer greater than 1 in which case it reduces to a polygamma function of a . It is easy to derive error bounds for polygamma functions that are accurate for large $|a|$, and we do not consider this special case further here.

4.8.2 Algorithmic matters

The evaluation of $\zeta(s + x, a)$ can be broken into three stages:

1. Choosing parameters M and N and bounding the remainder R .
2. Evaluating the power sum S .
3. Evaluating the tail T (and the trivial term I).

In this section, we describe some algorithmic techniques that are useful at each stage. We sketch the computational complexities, but do not attempt to prove strict complexity bounds.

We assume that arithmetic on real and complex numbers is done using ball arithmetic [103], which essentially is floating-point arithmetic with the added automatic propagation of error bounds. This is probably the most reasonable approach: *a priori* floating-point error analysis would be overwhelming to do in full generality (an analysis of the floating-point error when evaluating $\zeta(s)$ for real s , with a partial analysis of the complex case, is given in [88]).

Evaluating the error bound

For a precision of P bits, we should choose $N \sim M \sim P$. A simple strategy is to do a binary search for an N that makes the error bound small enough when $M = cN$ where $c \approx 1$. This is sufficient for our present purposes, but more sophisticated approaches are possible. In particular, for evaluation at large heights in the critical strip, N should be larger than M .

Given complex balls for s and a , and integers N and M , we can evaluate the error bound (4.8.12) using ball arithmetic. The output is a power series with ball coefficients. The

absolute value of each coefficient in this series should be added to the radius for the corresponding coefficient in $S + I + T \approx \zeta(s + x, a)$ at the end of the whole computation. If the assumptions that $\Re(a) + N > 1$ and $\Re(s) + 2M > 1$ are not satisfied for all points in the balls s and a , we set the error bounds for all coefficients to $+\infty$.

If we are computing D derivatives and D is large, the rising factorial $|(s + x)_{2M}|$ can be computed using binary splitting and the outer power series product in (4.8.12) can be done using fast polynomial multiplication, so that only $O^\sim(D + M)$ real number operations are required. Or, if D is small and M is large, $|(s + x)_{2M}|$ can be computed via the gamma function in time independent of M .

Evaluating the power sum

As a power series, the power sum S becomes $\sum_{k=0}^{N-1} (\sum_i c_i(k) x^i)$ where the coefficients $c_i(k)$ are given by (4.8.10). For $i \geq 1$, the coefficients can be computed using the recurrence

$$c_{i+1}(k) = -\frac{\log(a + k)}{i + 1} c_i(k).$$

If we are computing D derivatives with a working precision of P bits, the complexity of evaluating the power sum is $O^\sim(DNP)$, or $O^\sim(DN^2)$ if $N \sim P$. The computation is easy to parallelize by assigning a range of k values to each thread (for large D , a more memory-efficient method is to assign a range of i to each thread).

When evaluating the ordinary Riemann zeta function, i.e. when $a = 1$, and we just want to compute a small number of derivatives, we can speed up the power sum a bit. Writing the sum as $\sum_{k=1}^N f(k)$, the terms $f(k) = k^{-(s+x)}$ are completely multiplicative, i.e. $f(k_1 k_2) = f(k_1) f(k_2)$. This means that we only need to evaluate $f(k)$ from scratch when k is prime; when k is composite, a single multiplication is sufficient.

This method has two drawbacks: we have to store previously computed terms, which requires $O(DNP)$ space, and the power series multiplication $f(k_1) f(k_2)$ becomes more expensive than evaluating $f(k_1 k_2)$ from scratch for large D . For both reasons, this method is only useful when D is quite small (say $D \leq 4$).

We can avoid some redundant work by collecting multiples of small primes. For example, if we extract all powers of two, $\sum_{k=1}^{10} f(k)$ can be written as

$$\begin{aligned} & [f(1) + f(3) + f(5) + f(7) + f(9)] \\ & + f(2) [f(1) + f(3) + f(5)] \\ & + f(4) [f(1)] \\ & + f(8) [f(1)]. \end{aligned}$$

Algorithm 4.8.1 Sieved summation of a completely multiplicative function

Input: A function f such that $f(jk) = f(j)f(k)$ for $j, k \in \mathbb{Z}_{\geq 1}$, and an integer $N \geq 1$

Output: $\sum_{k=1}^N f(k)$

```
1:  $p \leftarrow 2^{\lceil \log_2 N \rceil}$  (largest power of two such that  $p \leq N$ )
2:  $h \leftarrow 1, z \leftarrow 0, u \leftarrow 0$ 
3:  $D = []$  ▷ Build table of divisors
4: for  $k \leftarrow 1; k \leq N; k \leftarrow k + 2$  do
5:    $D[k] \leftarrow 0$ 
6: for  $k \leftarrow 3; k \leq \lfloor \sqrt{N} \rfloor; k \leftarrow k + 2$  do
7:   if  $D[k] = 0$  then
8:     for  $j \leftarrow k^2; j \leq N; j \leftarrow j + 2k$  do
9:        $D[j] \leftarrow k$ 
10:  $F = []$  ▷ Create initially empty cache of  $f(k)$  values
11:  $F[2] \leftarrow f(2)$ 
12: for  $k \leftarrow 1; k \leq N; k \leftarrow k + 2$  do
13:   if  $D[k] = 0$  then ▷  $k$  is prime (or 1)
14:      $t \leftarrow f(k)$ 
15:   else
16:      $t \leftarrow F[D[k]]F[k/D[k]]$  ▷  $k$  is composite
17:   if  $3k \leq N$  then
18:      $F[k] \leftarrow t$  ▷ Store  $f(k)$  for future use
19:    $u \leftarrow u + t$ 
20:   while  $k = h$  and  $p \neq 1$  do ▷ Horner's rule
21:      $z \leftarrow u + F[2]z$ 
22:      $p \leftarrow p/2$ 
23:      $h \leftarrow \lfloor N/p \rfloor$ 
24:     if  $h$  is even then
25:        $h \leftarrow h - 1$ 
26: return  $u + F[2]z$ 
```

This is a polynomial in $f(2)$ and can be evaluated from bottom to top using Horner's rule while progressively adding the terms in the brackets. Asymptotically, this reduces the number of multiplications and the size of the tables by half. Algorithm 4.8.1 implements this trick, and requires about $\pi(N) \approx N/\log N$ evaluations of $f(k)$ and $N/2$ multiplications, at the expense of having to store about $N/6$ function values plus a table of divisors of about $N/2$ integers. Constructing the table of divisors using the sieve of Eratosthenes requires $O(N \log \log N)$ integer operations, but this cost is negligible when multiplications and $f(k)$ evaluations are expensive. One could also extract other powers besides 2 (for example powers of 3 and 5), but this gives diminishing returns.

Another trick that can save time at high precision is to avoid computing the logarithms of integers from scratch. If q and p are nearby integers (such as two consecutive primes) and we already know $\log(p)$, we can use the identity

$$\log(q) = \log(p) + 2 \operatorname{atanh} \left(\frac{q-p}{q+p} \right)$$

and evaluate the inverse hyperbolic tangent by applying binary splitting to its Taylor series. This is not an asymptotic improvement over the best known algorithm for computing the logarithm (which uses the arithmetic-geometric mean), but likely faster in practice.

If $D \sim N$, we can improve the asymptotic complexity of computing S to $O^\sim(DP)$, which is softly optimal in the bit size of the output (the author thanks David Harvey for sharing this observation). The vector of coefficients $((-1)^k k! [x^k] S)_{k=0}^{D-1}$ is given by $V^T Y$ where

$$V = \begin{bmatrix} 1 & \log a & \cdots & \log^{D-1} a \\ 1 & \log(a+1) & \cdots & \log^{D-1}(a+1) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \log(a+N-1) & \cdots & \log^{D-1}(a+N-1) \end{bmatrix}, \quad Y = \begin{bmatrix} a^{-s} \\ (a+1)^{-s} \\ \vdots \\ (a+N-1)^{-s} \end{bmatrix}.$$

It is well known that multiplying a vector from the left by the Vandermonde matrix V can be done in $O^\sim(N)$ coefficient operations in what amounts to *fast multipoint evaluation*. Multiplying a vector from the left by V^T when $D \sim N$ then has essentially the same complexity according to the transposition principle (this problem is discussed, for example, in [38]).

Evaluating the tail

Except for the multiplication by Bernoulli numbers, the terms of the tail sum T satisfy a simple (hypergeometric) recurrence relation. When we are computing D derivatives with a working precision of P bits, the complexity of evaluating the tail by repeated

application of the recurrence relation is $O^\sim(DMP)$, or $O^\sim(DP^2)$ if $M \sim P$. We can do better if D is large, using binary splitting (Algorithm 4.8.2).

Algorithm 4.8.2 Evaluation of the tail T using binary splitting

Input: $s, a \in \mathbb{C}$ and $N, M, D \in \mathbb{Z}_{\geq 1}$

Output: $T = \frac{1}{(a+N)^{s+x}} \left(\frac{1}{2} + \sum_{k=1}^M \frac{B_{2k}}{(2k)!} \frac{(s+x)_{2k-1}}{(a+N)^{2k-1}} \right) \in \mathbb{C}[[x]]/\langle x^D \rangle$

```

1: Let  $x$  denote the generator of  $\mathbb{C}[[x]]/\langle x^D \rangle$ 
2: function BINSPLIT( $j, k$ )
3:   if  $j + 1 = k$  then
4:     if  $j = 0$  then
5:        $P \leftarrow (s + x)/(2(a + N))$ 
6:     else
7:        $P \leftarrow \frac{(s + 2j - 1 + x)(s + 2j + x)}{(2j + 1)(2j + 2)(a + N)^2}$ 
8:     return  $(P, B_{2j+2}P)$ 
9:   else
10:     $(P_1, R_1) \leftarrow \text{BINSPLIT}(j, \lfloor (j + k)/2 \rfloor)$ 
11:     $(P_2, R_2) \leftarrow \text{BINSPLIT}(\lfloor (j + k)/2 \rfloor, k)$ 
12:    return  $(P_1P_2, R_1 + P_1R_2)$  ▷ Polynomial multiplications mod  $x^D$ 
13:  $(P, R) \leftarrow \text{BINSPLIT}(0, M)$ 
14:  $T \leftarrow (a + N)^{-(s+x)}(1/2 + R)$  ▷ Polynomial multiplication mod  $x^D$ 
15: return  $T$ 

```

If $D \sim M$, the complexity with binary splitting is only $O^\sim(PD)$, or softly optimal in the bit size of the output. A drawback is that the intermediate products increase the memory consumption.

4.8.3 Implementation and benchmarks

We have implemented the Hurwitz zeta function for $s \in \mathbb{C}[[x]]$ and $a \in \mathbb{C}$ with rigorous error bounds as part of the Arb library. Our implementation incorporates most of the techniques discussed in the previous section, including optional parallelization of the power sum (we have not implemented the fast algorithm based on multiplication by a transposed Vandermonde matrix). Bernoulli numbers are computed using the algorithm of Bloemen [10]. In the remainder of this section, we present the results of some computations done with our implementation.

Computing zeros to high precision

For $n \geq 1$, let ρ_n denote the n -th smallest zero of $\zeta(s)$ with positive imaginary part. We assume that ρ_n is simple and has real part $1/2$. Using Newton's method, we can evaluate ρ_n to high precision nearly as fast as we can evaluate $\zeta(s)$ for s near ρ_n .

It is convenient to work with real numbers. The ordinate $t_n = \Im(\rho_n)$ is a simple zero of the real-valued function $Z(t) = e^{i\theta(t)}\zeta(1/2 + it)$ where

$$\theta(t) = \frac{\log \Gamma\left(\frac{2it+1}{4}\right) - \log \Gamma\left(\frac{-2it+1}{4}\right)}{2i} - \frac{\log \pi}{2}t.$$

We assume that we are given an isolating ball $B_0 = [m_0 - \varepsilon_0, m_0 + \varepsilon_0]$ such that $t_n \in B_0$ and $t_m \notin B_0, m \neq n$, and wish to compute t_n to high precision (finding such a ball for a given n is an interesting problem, but we do not consider it here).

Newton's method maps an approximation z_n of a root of a real analytic function $f(z)$ to a new approximation z_{n+1} via $z_{n+1} = z_n - f(z_n)/f'(z_n)$. Using Taylor's theorem, the error can be shown to satisfy

$$|\epsilon_{n+1}| = \frac{|f''(\xi_n)|}{2|f'(z_n)|} |\epsilon_n|^2$$

for some ξ_n between z_n and the root.

As a setup step, we evaluate $Z(s), Z'(s), Z''(s)$ (simultaneously using power series arithmetic) at $s = B_0$, and compute

$$C = \frac{\max |Z''(B_0)|}{2 \min |Z'(B_0)|}.$$

This only needs to be done at low precision.

Starting from an input ball $B_k = [m_k - \varepsilon_k, m_k + \varepsilon_k]$, one step with Newton's method gives an output ball $B_{k+1} = [m_{k+1} - \varepsilon_{k+1}, m_{k+1} + \varepsilon_{k+1}]$. The updated midpoint is given by

$$m_{k+1} = m_k - \frac{Z(m_k)}{Z'(m_k)} \tag{4.8.15}$$

where we evaluate $Z(m_k)$ and $Z'(m_k)$ simultaneously using power series arithmetic. The updated radius is given by $\varepsilon_{k+1} = \varepsilon'_{k+1} + C\varepsilon_k^2$ where ε'_{k+1} is the numerical error (or a bound thereof) resulting from evaluating (4.8.15) using finite-precision arithmetic. The new ball is valid as long as $B_{k+1} \subseteq B_k$ (if this does not hold, the algorithm fails and we need to start with a better B_0 or increase the working precision).

For best performance, the evaluation precision should be chosen so that $\varepsilon'_{k+1} \approx C\varepsilon_k^2$. In other words, for a target accuracy of p bits, the evaluations should be done at $\dots, p/4, p/2, p$ bits, plus some guard bits.

Digits	mpmath		Mathematica		Arb	
	$\tilde{\rho}_1$	$\zeta(\tilde{\rho}_1)$	$\tilde{\rho}_1$	$\zeta(\tilde{\rho}_1)$	$\tilde{\rho}_1$	$\zeta(\tilde{\rho}_1)$
100	0.080	0.0031	0.044	0.012	0.012	0.0011
1000	7.1	0.24	11	1.6	0.18	0.05
10000	7035	252	5127	779	29	15
100000	-	-	-	-	6930	3476
303000	-	-	-	-	73225	31772

Table 4.9: Time in seconds to compute an approximation $\tilde{\rho}_1$ of the first nontrivial zero ρ_1 accurate to the specified number of decimal digits, and then to evaluate $\zeta(\tilde{\rho}_1)$ at the same precision. Computations were done on a 64-bit Intel Xeon E5-2650 2.00 GHz CPU.

As a benchmark problem, we compute an approximation $\tilde{\rho}_1$ of the first nontrivial zero $\rho_1 \approx 1/2 + 14.1347251417i$ and then evaluate $\zeta(\tilde{\rho}_1)$ to the same precision. We compare our implementation of the zeta function and the root-refinement algorithm described above (starting from a double-precision isolating ball) with the `zetazero` and `zeta` functions provided in mpmath version 0.17 in Sage 5.10 [97] and the `ZetaZero` and `Zeta` functions provided in Mathematica 9.0. The results of this benchmark are shown in Table 4.9. At 10000 digits, our code for computing the zero is about two orders of magnitude faster than the other systems, and the subsequent single zeta evaluation is about one order of magnitude faster.

We have computed ρ_1 to 303000 digits, or slightly more than one million bits, which appears to be a record (a 20000-digit value is given in [74]). The computation used up to 62 GiB of memory for the sieved power sum and the storage of Bernoulli numbers up to B_{325328} (to attain even higher precision, the memory usage could be reduced by evaluating the power sum without sieving, perhaps using several CPUs in parallel, and not caching Bernoulli numbers).

Computing the Keiper-Li coefficients

Riemann's function $\xi(s) = \frac{1}{2}s(s-1)\pi^{-s/2}\Gamma(s/2)\zeta(s)$ satisfies the symmetric functional equation $\xi(s) = \xi(1-s)$. The coefficients $\{\lambda_n\}_{n=1}^{\infty}$ defined by

$$\log \xi \left(\frac{1}{1-x} \right) = \log \xi \left(\frac{x}{x-1} \right) = -\log 2 + \sum_{n=1}^{\infty} \lambda_n x^n$$

were introduced by Keiper [60], who noted that the truth of the Riemann hypothesis would imply that $\lambda_n > 0$ for all $n > 0$. In fact, Keiper observed that if one makes an assumption about the distribution of the zeros of $\zeta(s)$ that is even stronger than the

	$n = 1000$	$n = 10000$	$n = 100000$
1: Error bound	0.017	1.0	97
1: Power sum	0.048	47	65402
(1: Power sum, CPU time)	(0.65)	(693)	(1042210)
1: Bernoulli numbers	0.0020	0.19	59
1: Tail	0.058	11	1972
2: Series logarithm	0.047	8.5	1126
3: $\log \Gamma(1 + x)$ series	0.019	3.0	1610
4: Composition	0.022	4.1	593
Total wall time	0.23	84	71051
Peak RAM usage (MiB)	8	730	48700

Table 4.10: Elapsed time in seconds to evaluate the Keiper-Li coefficients $\lambda_0 \dots \lambda_n$ with a working precision of $1.1n + 50$ bits, giving roughly $0.1n$ accurate bits. The computations were done on a multicore system with 64-bit Intel Xeon E7-8837 2.67 GHz CPUs (16 threads were used for the power sum, and all other parts were computed sequentially on a single core).

Riemann hypothesis, the coefficients λ_n should behave as

$$\lambda_n \approx (1/2) (\log n - \log(2\pi) + \gamma - 1). \quad (4.8.16)$$

Keiper presented numerical evidence for this conjecture by computing λ_n up to $n = 7000$, showing that the approximation error appears to fluctuate increasingly close to zero. Some years later, Li proved [72] that the Riemann hypothesis actually is equivalent to the positivity of λ_n for all $n > 0$ (this reformulation of the Riemann hypothesis is known as Li's criterion). Recently, Arias de Reyna has proved that a certain precise statement of (4.8.16) also is equivalent to the Riemann hypothesis [4].

A computation of the Keiper-Li coefficients up to $n = 100000$ shows agreement with Keiper's conjecture (and the Riemann hypothesis), as illustrated in Figure 4.4. We obtain $\lambda_{100000} = 4.62580782406902231409416038\dots$ (plus about 2900 more accurate digits), whereas (4.8.16) gives $\lambda_{100000} \approx 4.626132$. Empirically, we need a working precision of about n bits to determine λ_n accurately. A breakdown of the computation time to determine the signs of λ_n up to $n = 1000$, 10000 and 100000 is shown in Table 4.10.

Our computation of the Keiper-Li coefficients uses the formula

$$\log \xi(s) = \log(-\zeta(s)) + \log \Gamma\left(1 + \frac{s}{2}\right) + \log(1 - s) - \frac{s \log \pi}{2}$$

which we evaluate at $s = x \in \mathbb{R}[[x]]$. This arrangement of the terms avoids singularities and branch cuts at the expansion point. We carry out the following steps (plus some more trivial operations):

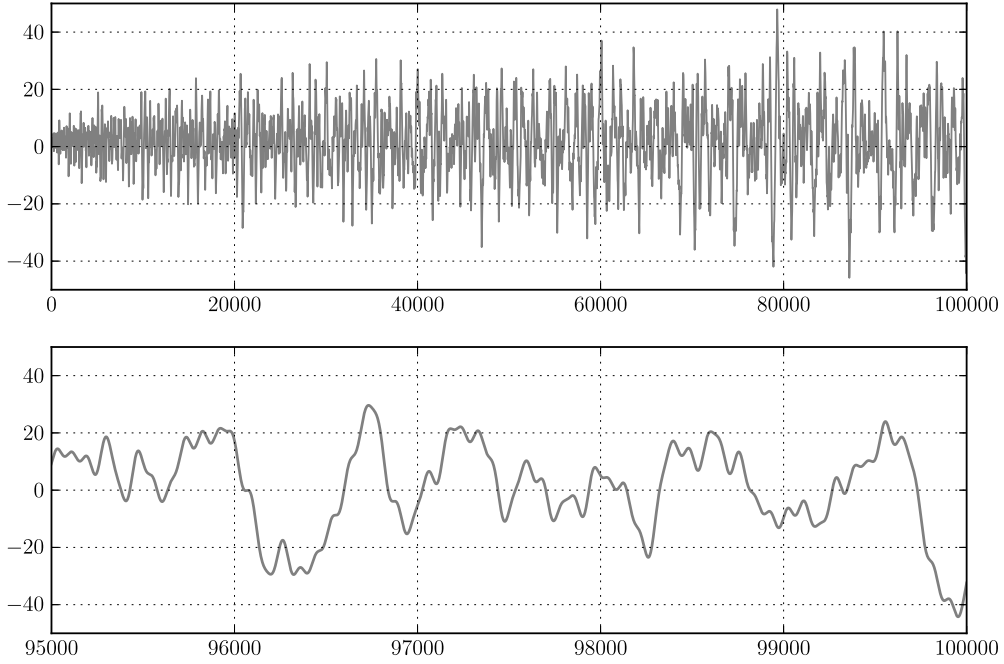


Figure 4.4: Plot of $n(\lambda_n - (\log n - \log(2\pi) + \gamma - 1)/2)$.

1. Computing the series expansion of $\zeta(s)$ at $s = 0$.
2. Computing the logarithm of a power series, i.e. $\log f(x) = \int f'(x)/f(x)dx$.
3. Computing the series expansion of $\log \Gamma(s)$ at $s = 1$, i.e. computing the sequence of values $\gamma, \zeta(2), \zeta(3), \zeta(4), \dots$
4. Finally, right-composing by $x/(x - 1)$ to obtain the Keiper-Li coefficients.

Step 2 requires $O(M(n))$ arithmetic operations on real numbers. For step 3, see the remarks in section 4.7. There is a very fast way to perform step 4. For $f = \sum_{k=0}^{\infty} a_k x^k \in \mathbb{C}[[x]]$, the binomial (or Euler) transform $T: \mathbb{C}[[x]] \rightarrow \mathbb{C}[[x]]$ is defined by

$$T[f(x)] = \frac{1}{1-x} f\left(\frac{x}{x-1}\right) = \sum_{n=0}^{\infty} \left(\sum_{k=0}^n (-1)^k \binom{n}{k} a_k \right) x^n.$$

We have

$$f\left(\frac{x}{x-1}\right) = a_0 + xT\left[\frac{a_0 - f}{x}\right].$$

If $B: \mathbb{C}[[x]] \rightarrow \mathbb{C}[[x]]$ denotes the Borel transform

$$B\left[\sum_{k=0}^{\infty} a_k x^k\right] = \sum_{k=0}^{\infty} \frac{a_k}{k!} x^k,$$

then (see [42]) $T[f(x)] = B^{-1}[e^x B[f(-x)]]$. This identity gives an algorithm for evaluating the composition which requires only $M(n) + O(n)$ coefficient operations. Moreover,

this algorithm is numerically stable (in the sense that it does not significantly increase errors from the input when using ball arithmetic), provided that a numerically stable polynomial multiplication algorithm is used.

The composition could also be carried out using various generic algorithms for composition of power series. We tested three other algorithms, and found them to perform much worse:

- Horner’s rule is slow (requiring about $nM(n)$ operations) and is numerically unsatisfactory in the sense that it gives extremely poor error bounds with ball arithmetic.
- The Brent-Kung algorithm BK 2.1 turns out to give adequate error bounds, but uses about $O(n^{1/2}M(n) + n^2)$ operations which still is expensive for large n .
- We also tried binary splitting: to evaluate $f(p/q)$ where f is a power series and p and q are polynomials, we recursively split the evaluation in half and keep numerator and denominator polynomials separated. In the end, we perform a single power series division. This only costs $O(M(n) \log n)$ operations, but turns out to be numerically unstable. It would be of independent interest to investigate whether this algorithm can be modified to avoid the stability problem.

Computing the Stieltjes constants

The generalized Stieltjes constants $\gamma_n(a)$ are defined by

$$\zeta(s, a) = \frac{1}{s-1} + \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} \gamma_n(a) (s-1)^n.$$

The “usual” Stieltjes constants are $\gamma_n(1) = \gamma_n$, and $\gamma_0 = \gamma \approx 0.577216$ is Euler’s constant. The Stieltjes constants were first studied over a century ago. Some historical notes and numerical values of γ_n for $n \leq 20$ are given in [12]. Keiper [60] provides a method for computing the Stieltjes constants based on numerical integration and recurrence relations, and lists various γ_n up to $n = 150$. This algorithm is implemented in Mathematica [111].

More recently, Kreminski [67] has given an algorithm for the Stieltjes constants, also based on numerical integration but different from Keiper’s. He reports having computed γ_n to a few thousand digits for all $n \leq 10000$, and provides further isolated values up to γ_{50000} (accurate to 1000 digits) as well as tables of $\gamma_n(a)$ with various $a \neq 1$.

The best proven bounds for the Stieltjes constants appear to be very pessimistic. In a recent paper, Knessl and Coffey [62] give an asymptotic approximation formula for the

Stieltjes constants that seems to be very accurate even for small n . Based on numerical computations done with Mathematica, they note that their approximation correctly predicts the sign of γ_n up to at least $n = 35000$ with the single exception of $n = 137$.

Our implementation immediately gives the generalized Stieltjes constants by computing the series expansion of $\zeta(s, a) - 1/(s - 1)$ at $s = 1$ using (4.8.11). The costs are similar to those for computing the Keiper-Li coefficients: due to ill-conditioning, it appears that we need about $n + p$ bits of precision to determine γ_n with p bits of accuracy. This makes our method somewhat unattractive for computing just a few digits of γ_n when n is large, but reasonably good if we want a large number of digits. Our method is also useful if we want to compute a table of all the values $\gamma_0, \dots, \gamma_n$ simultaneously.

For example, we can compute γ_n for all $n \leq 1000$ to 1000-digit accuracy in just over 10 seconds on a single CPU. Computing the single coefficient γ_{1000} to 1000-digit accuracy with Mathematica 9.0 takes 80 seconds, with an estimated 20 hours required for all $n \leq 1000$. Thus our implementation is nearly four orders of magnitude faster. We can compute a table of accurate values of γ_n for all $n \leq 10000$ in a few minutes on an ordinary workstation with around one GiB of memory.

We have computed all γ_n up to $n = 100000$ using a working precision of 125050 bits, resulting in an accuracy from about 37640 decimal digits for γ_0 to about 10860 accurate digits for γ_{100000} . The computation took 26 hours on a multicore system with 16 threads utilized for the power sum, with a peak memory consumption of about 80 GiB during the binary splitting evaluation of the tail. As shown in Figure 4.5, the accuracy of the Knessl-Coffey approximation approaches six digits on average. Our computation gives $\gamma_{100000} = 1.991927306312541095658 \dots \times 10^{83432}$, while the Knessl-Coffey approximation gives $\gamma_n \approx 1.9919333 \times 10^{83432}$. We are able to confirm that $n = 137$ is the only instance for $n \leq 100000$ where the Knessl-Coffey approximation has the wrong sign.

In [61], Knessl and Coffey extend their asymptotic approximation formula to $a \neq 1$. Their approximation gives, for example, the estimate

$$\gamma_{50000}(1 + i) \approx (1.0324943 - 1.4419586i) \times 10^{39732}$$

while we compute (rounded to 15 decimal digits)

$$\gamma_{50000}(1 + i) = (1.03250208743188 - 1.44196255284053i) \times 10^{39732}.$$

We emphasize that our implementation computes $\gamma_n(a)$ with proved error bounds, while the other cited works and implementations (to our knowledge) depend on heuristic error estimates.

We have not yet implemented a function for computing isolated Stieltjes constants of large index; this would have roughly the same running time as the evaluation of the

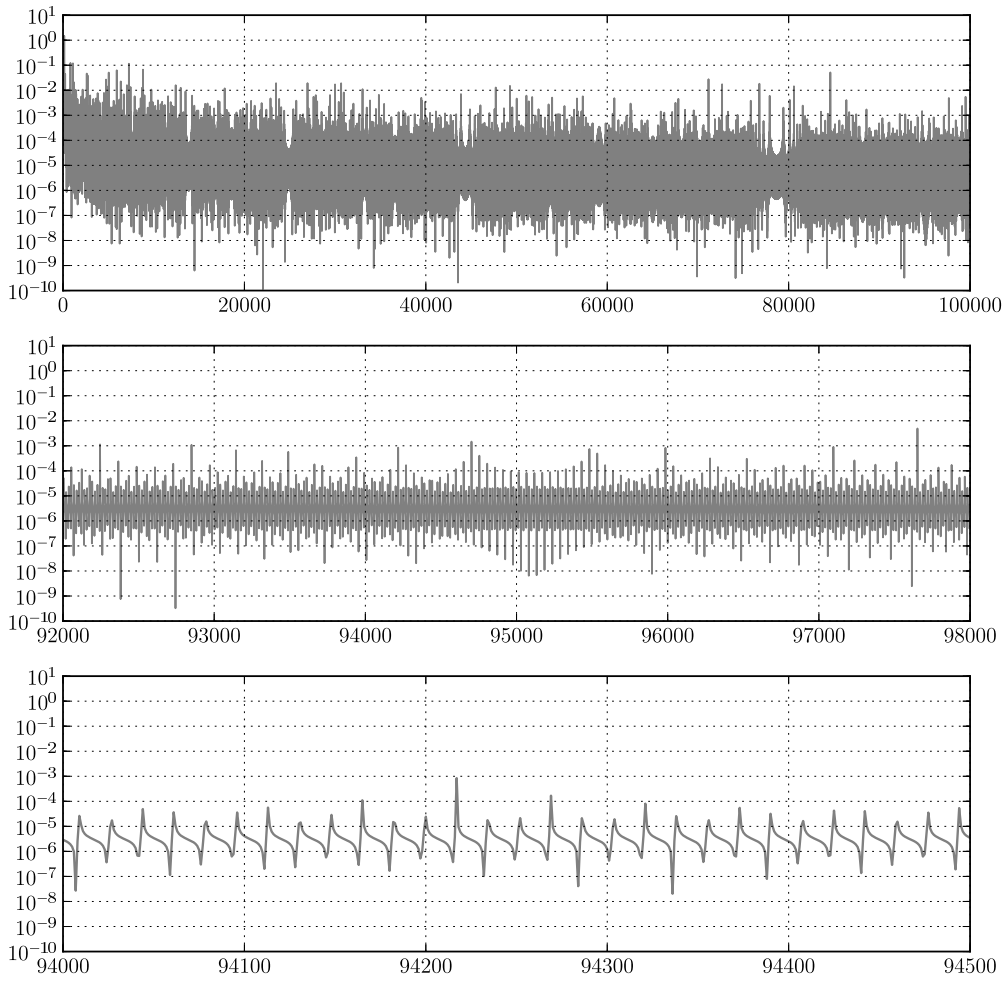


Figure 4.5: Plot of the relative error $|\gamma_n - \tilde{\gamma}_n|/|\gamma_n|$ of the Knessl-Coffey approximation for the Stieltjes constants. The error exhibits a complex oscillation pattern.

tail (since only a single derivative of the power sum would have to be computed). The memory consumption is highest when evaluating the tail, and would therefore remain the same.

4.8.4 Remarks

One direction for further work would be to improve the error bounds for large $|a|$ and to investigate strategies for selecting N and M optimally, particularly when the number of derivatives is large. It would also be interesting to investigate parallelization of the tail sum, or look for ways to evaluate a single derivative of high order of the tail in a memory-efficient way. Further constant-factor improvements are possible in an implementation, for example by reducing the precision of terms that have small magnitude (rather than

naively performing all operations at the same precision). It would also be interesting to implement the asymptotically fast algorithm for the power sum when computing many derivatives.

Finally, it would be interesting to compare the efficiency of the Euler-Maclaurin formula with other approaches to evaluating the Hurwitz zeta function such as the algorithms of Borwein [17], Vepštas [105] and Coffey [34].

Chapter 5

Computing the partition function

Let $p(n)$ denote the number of partitions of n , i.e. the number of ways that n can be written as a sum of positive integers without regard to the order of the terms (A000041 in [83]). For instance $5 = 4 + 1 = 3 + 2 = 3 + 1 + 1 = 2 + 2 + 1 = 2 + 1 + 1 + 1 = 1 + 1 + 1 + 1 + 1$, so $p(5) = 7$. A few more values are $p(10) = 42$, $p(100) = 190569292$, $p(1000) = 24061467864032622473692149727991$. Clearly, the last value was not found by explicitly listing all the partitions of 1000.

5.1 The pentagonal number theorem

It is a classical result of Euler, known as the pentagonal number theorem, that the generating function of $p(n)$ satisfies

$$\sum_{n=0}^{\infty} p(n)x^n = \frac{1}{\phi(x)} \quad (5.1.1)$$

where

$$\phi(x) = \prod_{k=1}^{\infty} (1 - x^k) = \sum_{k=-\infty}^{\infty} (-1)^k x^{k(3k-1)/2}. \quad (5.1.2)$$

The function $\eta(\tau) = q^{1/24}\phi(q)$, where $q = e^{2\pi i\tau}$, is called the Dedekind eta function, and it is an example of a modular form. From (5.1.2), Euler also derived the recursive formula

$$p(n) = \sum_{k=1}^n (-1)^{k+1} \left(p\left(n - \frac{k(3k-1)}{2}\right) + p\left(n - \frac{k(3k+1)}{2}\right) \right). \quad (5.1.3)$$

MacMahon famously used (5.1.3) to compute $p(n)$ up to $n = 200$ by hand. When

examining MacMahon's table, Ramanujan discovered the congruences

$$\begin{aligned} p(5k + 4) &\equiv 0 \pmod{5} \\ p(7k + 5) &\equiv 0 \pmod{7} \\ p(11k + 6) &\equiv 0 \pmod{11} \end{aligned}$$

which he later also proved. Ramanujan's congruences have inspired a vast amount of research, providing an excellent example of the usefulness of the computational approach to number theory.

Determining the list of values $p(0), p(1), \dots, p(n-1), p(n)$ using (5.1.3) requires $O(n^{1.5})$ arithmetic operations. Applying power series inversion to the right-hand side of (5.1.2) with an FFT-based multiplication algorithm requires $O(M(n))$ arithmetic operations, giving a nearly optimal procedure for multi-evaluation of the partition function.

An attractive feature of the pentagonal number theorem, in both the recursive and FFT incarnations, is that the values can be computed more efficiently modulo a small prime number. This is useful for investigating partition function congruences, such as in a recent large-scale computation of $p(n)$ modulo small primes for n up to 10^9 [80].

5.2 The Hardy-Ramanujan-Rademacher formula

While efficient for computing $p(n)$ for all n up to some limit, the pentagonal number theorem is impractical for evaluating $p(n)$ for an isolated, large n . One of the most astonishing number-theoretical discoveries of the 20th century is the Hardy-Ramanujan-Rademacher formula, first given as an asymptotic expansion by Hardy and Ramanujan in 1917 [46] and subsequently refined to an exact representation by Rademacher in 1936 [91], which provides a direct and computationally efficient expression for the single value $p(n)$.

Simplified to a first-order estimate, the Hardy-Ramanujan-Rademacher (HRR) formula states that

$$p(n) \sim \frac{1}{4n\sqrt{3}} e^{\pi\sqrt{2n/3}}, \quad (5.2.1)$$

from which one gathers that $p(n)$ is a number with roughly $n^{1/2}$ decimal digits. The

full version can be stated as

$$p(n) = \left(\sum_{k=1}^N T(n, k) \right) + R(n, N), \quad (5.2.2)$$

$$T(n, k) = \left(\sqrt{\frac{3}{k}} \frac{4}{24n-1} \right) A_k(n) U \left(\frac{C(n)}{k} \right), \quad (5.2.3)$$

$$U(x) = \cosh(x) - \frac{\sinh(x)}{x}, \quad (5.2.4)$$

$$C(n) = \frac{\pi}{6} \sqrt{24n-1}, \quad (5.2.5)$$

$$A_k(n) = \sum_{h=0}^{k-1} \delta_{\gcd(h,k),1} \exp \left(\pi i \left[s(h, k) - \frac{2hn}{k} \right] \right) \quad (5.2.6)$$

where $s(h, k)$ is the Dedekind sum

$$s(h, k) = \sum_{i=1}^{k-1} \frac{i}{k} \left(\frac{hi}{k} - \left\lfloor \frac{hi}{k} \right\rfloor - \frac{1}{2} \right) \quad (5.2.7)$$

and where the remainder satisfies $|R(n, N)| < M(n, N)$ with

$$M(n, N) = \frac{44\pi^2}{225\sqrt{3}} N^{-1/2} + \frac{\pi\sqrt{2}}{75} \left(\frac{N}{n-1} \right)^{1/2} \sinh \left(\frac{\pi}{N} \sqrt{\frac{2n}{3}} \right). \quad (5.2.8)$$

It is easily shown that $M(n, cn^{1/2}) \sim n^{-1/4}$ for every positive c . Rademacher's bound (5.2.8) therefore implies that $O(n^{1/2})$ terms in (5.2.2) suffice to compute $p(n)$ exactly by forcing $|R(n, N)| < 1/2$ and rounding to the nearest integer. For example, we can take $N = \lceil n^{1/2} \rceil$ when $n \geq 65$.

The main result which we prove in this chapter is the following:

Theorem 5.2.1. *The value $p(n)$ can be computed using $O(n^{1/2} \log^{4+o(1)} n)$ bit operations.*

Since $p(n)$ has $\Theta(n^{1/2})$ bits, we show that the partition function can be computed in softly optimal time. The possibility to compute $p(n)$ in softly optimal time was already suggested by Odlyzko [81], but he did not give a proof or a complete algorithm. In fact, even with the use of fast arithmetic, the steps required to achieve a softly optimal complexity bound are nontrivial.

The existence of a softly optimal algorithm to compute $p(n)$ is interesting since the partition function is not known to belong to any general class of sequences for which this property holds "automatically". For example, Fibonacci numbers and factorials are P-finite and can therefore be computed fast, but $p(n)$ is not P-finite (since its generating function has a natural boundary of analyticity on the unit circle).

In contrast to the case of integer partitions, no softly optimal algorithm is known for counting *set partitions*, i.e. isolated values in the sequence of Bell numbers B_n defined by

$$\sum_{n=0}^{\infty} \frac{B_n}{n!} x^n = e^{e^x - 1}.$$

Multi-evaluation of Bell numbers can nonetheless be done nearly optimally using fast power series arithmetic.

The computational utility of the Hardy-Ramanujan-Rademacher formula was realized before the availability of electronic computers. For instance, Lehmer [68] used it to verify Ramanujan’s conjectures $p(599) \equiv 0 \pmod{5^3}$ and $p(721) \equiv 0 \pmod{11^2}$.

Implementations are now available in numerous mathematical software systems, including Pari/GP, Maple, Mathematica and Sage. However, apart from Odlyzko’s remark, we find few algorithmic accounts of the Hardy-Ramanujan-Rademacher formula in the literature, nor any investigation into the optimality of the available implementations.

We have implemented the partition function in the FLINT and Arb libraries, and observe that our implementations behave quasioptimally in practice, improving on the speed of previously published software by more than two orders of magnitude.

We benchmark our code by computing some extremely large isolated values of $p(n)$. We also investigate efficiency compared to power series methods for evaluation of multiple values, and finally apply our implementation in FLINT to the problem of computing congruences for $p(n)$.

5.3 Evaluating the exponential sums

A naive implementation of formulas (5.2.2)–(5.2.7) requires $O(n^{3/2})$ integer operations to evaluate Dedekind sums, and $O(n)$ numerical evaluations of complex exponentials (or cosines, since the imaginary parts ultimately cancel out). In the following section, we describe how the number of integer operations and cosine evaluations can be reduced, for the moment ignoring numerical evaluation.

A first improvement, used for instance in Bober’s implementation of $p(n)$ for the Sage computer algebra system, is to recognize that Dedekind sums can be evaluated in $O(\log k)$ steps using a GCD-style algorithm, as described by Apostol [3], or with Knuth’s fraction-free algorithm [63] which avoids the overhead of rational arithmetic. This reduces the total number of integer operations to $O(n \log n)$, which is a dramatic improvement but still leaves the cost of computing $p(n)$ quadratic in the size of the final result.

Fortunately, the $A_k(n)$ sums have additional structure as discussed in [69, 70, 92, 110, 86], allowing the computational complexity to be reduced. Since numerous implementers of the Hardy-Ramanujan-Rademacher formula until now appear to have overlooked these results, it seems appropriate that we reproduce the main formulas and assess the computational issues in more detail.

5.3.1 A simple algorithm

Using properties of the Dedekind eta function, one can derive the formula (which Whiteman [110] attributes to Selberg)

$$A_k(n) = \left(\frac{k}{3}\right)^{1/2} \sum_{(3l^2+l)/2 \equiv -n \pmod k} (-1)^l \cos\left(\frac{6l+1}{6k}\pi\right) \quad (5.3.1)$$

in which the summation ranges over $0 \leq l < 2k$ and only $O(k^{1/2})$ terms are nonzero. With a simple brute force search for solutions of the quadratic equation, this representation provides a way to compute $A_k(n)$ that is both simpler and more efficient than the usual definition (5.2.6).

Although a brute force search requires $O(k)$ loop iterations, the successive quadratic terms can be generated without multiplications or divisions using two coupled linear recurrences. This only costs a few processor cycles per loop iteration, which is a substantial improvement over computing Dedekind sums, and means that the cost up to fairly large k effectively will be dominated by evaluating $O(k^{1/2})$ cosines, adding up to $O(n^{3/4})$ function evaluations for computing $p(n)$.

Algorithm 5.3.1 Simple algorithm for evaluating $A_k(n)$

Input: Integers $k, n \geq 0$

Output: $s = A_k(n)$, where $A_k(n)$ is defined as in (5.2.6)

- 1: **if** $k \leq 1$ **then return** k
 - 2: **else if** $k = 2$ **then return** $(-1)^n$
 - 3: $(s, r, m) \leftarrow (0, 2, (n \bmod k))$
 - 4: **for** $0 \leq l < 2k$ **do**
 - 5: **if** $m = 0$ **then**
 - 6: $s \leftarrow s + (-1)^l \cos(\pi(6l+1)/(6k))$
 - 7: $m \leftarrow m + r$
 - 8: **if** $m \geq k$ **then** $m \leftarrow m - k$ $\triangleright m \leftarrow m \bmod k$
 - 9: $r \leftarrow r + 3$
 - 10: **if** $r \geq k$ **then** $r \leftarrow r - k$ $\triangleright r \leftarrow r \bmod k$
 - 11: **return** $(k/3)^{1/2} s$
-

A basic implementation of (5.3.1) is given as Algorithm 5.3.1. Here the variable m runs over the successive values of $(3l^2 + l)/2$, and r runs over the differences between consecutive m . Various improvements are possible: a modification of the equation allows cutting the loop range in half when k is odd, and the number of cosine evaluations can be reduced by counting the multiplicities of unique angles after reduction to $[0, \pi/4)$, evaluating a weighted sum $\sum w_i \cos(\theta_i)$ at the end – possibly using trigonometric addition theorems to exploit the fact that the differences $\theta_{i+1} - \theta_i$ between successive angles tend to repeat for many different i .

5.3.2 A fast algorithm

From Selberg’s formula (5.3.1), Whiteman [110] obtained a more efficient but considerably more complicated multiplicative decomposition of $A_k(n)$. The essence of Whiteman’s factorization is the following.

Theorem 5.3.1. *$A_k(n)$ can be written as $\sqrt{a} \prod_{i=1}^r \cos(\pi b)$ where $r = O(\log k)$ and a and b are rational numbers with $O(\log n + \log k)$ bits in the numerator and denominator.*

Since Whiteman’s factorization only involves $O(\log k)$ cosine factors per term in the HRR series, the total number for $p(n)$ is brought down to $O(n^{1/2} \log n)$. His factorization also reveals exactly when $A_k(n) = 0$ (which is about half the time).

We now reproduce the details of Whiteman’s factorization theorem. We omit the proofs, which are given in full detail in [110].

First consider the case when k is a power of a prime. Clearly $A_1(n) = 1$ and $A_2(n) = (-1)^n$. Otherwise let $k = p^\lambda$ and $v = 1 - 24n$. Then, using the notation $(a|m)$ for Jacobi symbols to avoid confusion with fractions, we have

$$A_k(n) = \begin{cases} (-1)^\lambda (-1|m_2) k^{1/2} \sin(4\pi m_2/8k) & \text{if } p = 2 \\ 2(-1)^{\lambda+1} (m_3|3) (k/3)^{1/2} \sin(4\pi m_3/3k) & \text{if } p = 3 \\ 2(3|k) k^{1/2} \cos(4\pi m_p/k) & \text{if } p > 3 \end{cases} \quad (5.3.2)$$

where m_2 , m_3 and m_p respectively are any solutions of

$$(3m_2)^2 \equiv v \pmod{8k} \quad (5.3.3)$$

$$(8m_3)^2 \equiv v \pmod{3k} \quad (5.3.4)$$

$$(24m_p)^2 \equiv v \pmod{k} \quad (5.3.5)$$

provided, when $p > 3$, that such an m_p exists and that $\gcd(v, k) = 1$. If, on the other

hand, $p > 3$ and either of these two conditions do not hold, we have

$$A_k(n) = \begin{cases} 0 & \text{if } v \text{ is not a quadratic residue modulo } k \\ (3|k)k^{1/2} & \text{if } v \equiv 0 \pmod{p}, \lambda = 0 \\ 0 & \text{if } v \equiv 0 \pmod{p}, \lambda > 1. \end{cases} \quad (5.3.6)$$

If k is not a prime power, assume that $k = k_1 k_2$ where $\gcd(k_1, k_2) = 1$. Then we can factor $A_k(n)$ as $A_k(n) = A_{k_1}(n_1)A_{k_2}(n_2)$, where n_1, n_2 are any solutions of the following equations. If $k_1 = 2$, then

$$\begin{cases} 32n_2 \equiv 8n + 1 \pmod{k_2} \\ n_1 \equiv n - (k_2^2 - 1)/8 \pmod{2}, \end{cases} \quad (5.3.7)$$

if $k_1 = 4$, then

$$\begin{cases} 128n_2 \equiv 8n + 5 \pmod{k_2} \\ k_2^2 n_1 \equiv n - 2 - (k_2^2 - 1)/8 \pmod{4}, \end{cases} \quad (5.3.8)$$

and if k_1 is odd or divisible by 8, then

$$\begin{cases} k_2^2 d_2 e n_1 \equiv d_2 e n + (k_2^2 - 1)/d_1 \pmod{k_1} \\ k_1^2 d_1 e n_2 \equiv d_1 e n + (k_1^2 - 1)/d_2 \pmod{k_2} \end{cases} \quad (5.3.9)$$

where $d_1 = \gcd(24, k_1)$, $d_2 = \gcd(24, k_2)$, $24 = d_1 d_2 e$.

Here $(k^2 - 1)/d$ denotes an operation done on integers, rather than a modular division. All other solving steps in (5.3.2)–(5.3.9) amount to computing greatest common divisors, carrying out modular ring operations, finding modular inverses, and computing modular square roots.

Repeated application of these formulas results in Algorithm 5.3.2, where we omit the detailed arithmetic for brevity.

We now analyze the complexity of Algorithm 5.3.2.

Lemma 5.3.2. *Assume that the prime factorization k is given, and that we know a quadratic nonresidue modulo any prime up to k . Then the exponential sum $A_k(n)$ can be factored using $O(\log^{4+o(1)} k)$ bit operations by use of Algorithm 5.3.2.*

Proof. A fixed index k is a product of at most $O(\log k)$ prime powers with exponents bounded by $O(\log k)$. For each prime power, Algorithm 5.3.2 performs $O(1)$ operations on integers with $O(\log k)$ bits, where an operation is an addition, multiplication, a division, a modular inversion, a greatest common divisor, or a Jacobi symbol (which can all be done with bit complexity $O(\log^{1+o(1)} k)$), or a square root modulo p^λ .

Algorithm 5.3.2 Fast algorithm for evaluating $A_k(n)$

Input: Integers $k \geq 1, n \geq 0$

Output: $s = A_k(n)$, where $A_k(n)$ is defined as in (5.2.6)

```
1: Compute the prime factorization  $k = p_1^{\lambda_1} p_2^{\lambda_2} \dots p_j^{\lambda_j}$ 
2:  $s \leftarrow 1$ 
3: for  $1 \leq i \leq j$  and while  $s \neq 0$  do
4:   if  $i < j$  then
5:      $(k_1, k_2) \leftarrow (p_i^{\lambda_i}, k/p_i^{\lambda_i})$ 
6:     Compute  $n_1, n_2$  by solving the respective case of (5.3.7)–(5.3.9)
7:      $s \leftarrow s \times A_{k_1}(n_1)$   $\triangleright$  Handle the prime power case using (5.3.2)–(5.3.6)
8:      $(k, n) \leftarrow (k_2, n_2)$ 
9:   else
10:     $s \leftarrow s \times A_k(n)$   $\triangleright$  Prime power case
11: return  $s$ 
```

To compute square roots modulo p^λ , we can use the Tonelli-Shanks algorithm [99, 94] modulo p followed by Hensel lifting up to p^λ . Assuming that we know a quadratic nonresidue modulo p , the Tonelli-Shanks algorithm requires $O(\log^3 k)$ multiplications in the worst case and $O(\log^2 k)$ multiplications on average. Hensel lifting costs $O(\log^{o(1)} k)$ multiplications. We thus perform at most $O(\log^2 k)$ multiplications for each of the $O(\log k)$ prime powers, each multiplication having bit complexity $O(\log^{1+o(1)} k)$. \square

At first sight, the need for the prime factorization of k might seem to pose a problem for Algorithm 5.3.2, since no $\log^{O(1)} k$ algorithm for integer factorization is known. More subtly, finding quadratic nonresidues poses a similar difficulty. Fortunately, to compute $p(n)$, these operations can be batched for all the required indices.

Lemma 5.3.3. *We can compute the prime factorizations of all k up to $n^{1/2}$, and find a quadratic nonresidue modulo p for all primes p up to $n^{1/2}$, using $O(n^{1/2} \log^{1+o(1)} n)$ bit operations.*

Proof. Using the sieve of Eratosthenes, we can precompute a list of length $n^{1/2}$ where entry k is the largest prime dividing k using $O(n^{1/2} \log^{1+o(1)} n)$ bit operations.

If $n_2(p_k)$ denotes the least quadratic nonresidue modulo the k th prime number, it is a theorem of Erdős [37, 89] that as $x \rightarrow \infty$,

$$\frac{1}{\pi(x)} \sum_{p_k \leq x} n_2(p_k) \rightarrow \sum_{k=1}^{\infty} \frac{p_k}{2^k} = C < 3.675. \quad (5.3.10)$$

Given the primes up to $x = n^{1/2}$, we can therefore build a table of nonresidues by testing no more than $(C + o(1))\pi(n^{1/2})$ candidates. Since $\pi(n^{1/2}) = O(n^{1/2}/\log n)$ and

a quadratic residue test takes $O(\log^{1+o(1)} p)$ time, the total precomputation time for quadratic residues is $O(n^{1/2} \log^{o(1)} n)$. \square

Combining Lemmas 5.3.2 and 5.3.3 gives the following result.

Theorem 5.3.4. *The factorizations of the exponential sums $A_k(n)$ for all k up to $O(n^{1/2})$ can be computed using a total of $O(n^{1/2} \log^{4+o(1)} n)$ bit operations.*

The precomputation of Lemma 5.3.3 is mostly of theoretical interest. In practice, the k are word-sized and can be factored on the fly in time that is negligible compared to all the other operations involved in evaluating $p(n)$. Likewise, it is sufficient in practice to generate nonresidues on the fly since $O(1)$ candidates need to be tested on average, but we can only prove an $O(\log^c k)$ bound for factoring an isolated $A_k(n)$ (where the prime factorization of k also is given) by assuming the Extended Riemann Hypothesis which gives $n_2(p) = O(\log^2 p)$ [2].

In our original paper [52], we erroneously claimed an $O(n^{1/2} \log^{3+o(1)} n)$ complexity bound for factoring all the $A_k(n)$. This analysis was based on using Cipolla's algorithm [33] which requires $O(\log^2 k)$ multiplications in the worst case to compute a modular square root (as discussed in [35], pp. 99–103). However, while the smallest quadratic nonresidue modulo a given p can be used in Tonelli-Shanks algorithm, Cipolla's algorithm requires knowing a quadratic nonresidue of special form. The original complexity bound is still valid as a heuristic estimate, and could possibly be proved rigorously by a stronger argument.

5.4 Numerical evaluation

To evaluate $p(n)$ using the HRR formula, the following strategy can be used: we first select $N = O(n^{1/2})$ such that the remainder satisfies $|R(n, N)| < 0.25$, and then compute a floating-point number $s = u2^v$, $u, v \in \mathbb{Z}$ such that

$$\left| \sum_{k=1}^N T(n, k) - s \right| < 0.25. \quad (5.4.1)$$

Since $|p(n) - s| < 0.5$, rounding s to the nearest integer guarantees a correct result.

Clearly, s needs to be computed with a precision of $O(n^{1/2})$ bits. Even with the aid of the factorization of the previous section, evaluating all terms to the same precision would cost $O(n^{1/2}) \times O(n^{1/2+o(1)}) = O(n^{1+o(1)})$ bit operations. The key step to get a quasioptimal algorithm is to choose a tight working precision for each term.

The bound (5.4.1) holds if each term $T(n, k)$ is approximated with an absolute error of

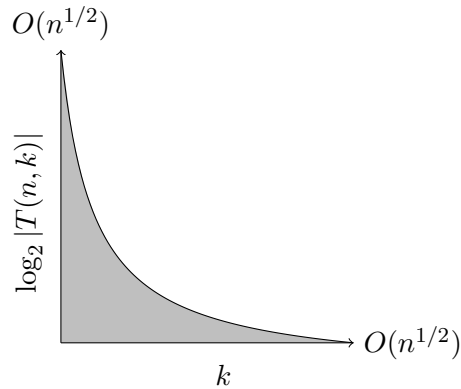


Figure 5.1: Distribution of the term magnitudes in the Hardy-Ramanujan-Rademacher series: we need a few terms with high precision and many terms with low precision.

at most $0.125/N$ and each addition in the main summation is done with an absolute error of at most $0.125/N$.

Since $\log_2 |T(n, k)| = O(n^{1/2}/k + \log k)$, it is sufficient to evaluate $T(n, k)$ from its factorization using floating-point arithmetic with a precision of $\log_2 |T(n, k)| + O(\log n) = O(n^{1/2}/k + \log n)$ bits, where the extra $O(\log n)$ term accounts for intermediate rounding errors.

A more detailed error analysis, providing an explicit numerical bound for the number of guard bits under the assumption that the operations are done in a certain order, is given in our original paper [52].

Up to additional logarithmic factors, the preceding argument shows that the complexity of numerical evaluation required to compute $p(n)$ is the area below the hyperbola $O(n^{1/2}/k)$ as k goes up to $O(n^{1/2})$, which amounts to $O(n^{1/2} \log n)$. This is illustrated in Figure 5.1. If a fixed precision were used for all terms, the complexity would be the area of the enveloping rectangle.

Stated more precisely, we have the following result:

Theorem 5.4.1. *Assume that the factorizations of $A_k(n)$ for all k up to $N = O(n^{1/2})$ are given. Then we can compute a floating-point number s such that*

$$\left| \sum_{k=1}^N T(n, k) - s \right| < 0.25$$

using $O(n^{1/2} \log^{4+o(1)} n)$ bit operations.

Proof. By Theorem 5.3.1, each term $T(n, k)$ can be evaluated using $O(\log k)$ elementary operations (arithmetic operations, real exponentials, and real cosines), each of which needs to be carried out at a precision of $O(n^{1/2}/k + \log n)$ bits.

The elementary functions can be evaluated to a precision of r bits using $E(r) \equiv O(M(r) \log^\alpha r) = O(r \log^{1+\alpha+o(1)} r)$ bit operations, where $\alpha = 1$ using the arithmetic-geometric mean iteration and $\alpha = 2$ using the bit-burst (binary splitting) algorithm.

This estimate assumes that the function argument is bounded. For $\cos(x)$ where $x = (p/q)\pi$, we can enforce $0 < x < \pi/2$ by adjusting the integers p and q . For $\exp(x)$ where $x = O(n^{1/2})$, we can reduce the argument to a standard interval such as $(-1, 1)$ via $\exp(x) = \exp(x - m \log 2) 2^m$ where a single floating-point division determines the integer m . This argument reduction step requires no more than $O(\log n^{1/2}) = O(\log n)$ guard bits. Obviously, the constants π and $\log 2$ only need to be computed once.

Numerically evaluating all terms thus costs

$$B = O\left(\sum_{k=1}^{n^{1/2}} (\log k) E(n^{1/2}/k + \log n)\right)$$

bit operations. To clean up the nested logarithms that appear when expanding the E term, we note that $\log(n^{1/2}/k + \log n) = O(\log n^{1/2}) = O(\log n)$. This gives

$$\begin{aligned} B &= O\left(\sum_{k=1}^{n^{1/2}} (\log k) (n^{1/2}/k + \log n) (\log^{1+\alpha+o(1)} n)\right) \\ &= O(\log^{2+\alpha+o(1)} n) \left(\sum_{k=1}^{n^{1/2}} n^{1/2}/k + \log n\right) \\ &= O(\log^{2+\alpha+o(1)} n) O\left(\int_1^{n^{1/2}} n^{1/2}/k + \log n \, dk\right) \\ &= O(\log^{2+\alpha+o(1)}) O(n^{1/2} \log n) \\ &= O(n^{1/2} \log^{3+\alpha+o(1)}). \end{aligned}$$

Finally, we need to compute $s = \sum_{k=1}^N T(n, k)$. If the terms are added in reverse order $k = N, \dots, 2, 1$, then addition k costs $O(n^{1/2}/k + \log n)$ bit operations, which adds up to less than the cost of evaluating all the terms. \square

The main result (Theorem 5.2.1) now follows from Theorem 5.3.4 and Theorem 5.4.1.

The implementation of additions is a subtle but crucial point. If the additions are done in forward order, the total complexity would be $O(n)$. We can still get the best complexity when generating terms in forward order by amortizing the additions: we keep separate summation variables for the partial sums of terms not exceeding $r_1, r_1/2, r_1/4, r_1/8, \dots$ bits where r_1 is the precision of the first term.

Alternatively, if the additions are performed in-place in memory, we can sum in the forward order and rely on carry propagation terminating in an expected $O(1)$ steps, but

many implementations of arbitrary-precision floating-point arithmetic do not provide this optimization.

5.5 Rapid computation of roots of unity

To determine $p(n)$ using the Hardy-Ramanujan-Rademacher formula, we require numerical approximations of a large number of algebraic numbers, namely the special trigonometric values

$$\alpha_{p,q} = \cos\left(\frac{p\pi}{q}\right)$$

where p, q are coprime integers, $q > 0$, and where we can assume that $0 < p < 2q$. Computing $\alpha(p, q)$ is equivalent to computing the root of unity

$$\beta_{p,q} = \exp\left(\frac{p\pi i}{q}\right)$$

since $\alpha_{p,q} = \operatorname{Re}(\beta_{p,q}) = (\beta_{p,q} + \overline{\beta_{p,q}})/2$ and $\beta_{p,q} = \alpha_{p,q} \pm i\sqrt{1 - \alpha_{p,q}^2}$.

The best way to compute $\alpha_{p,q}$ or $\beta_{p,q}$ depends on q and the precision. When q is large and the precision is small (which occurs in the majority of the terms in the HRR series), a satisfactory algorithm is to just compute a floating-point approximation of $x = p\pi/q$ and then evaluate $\cos(x)$ the normal way, for instance using Taylor series (to improve multi-evaluation speed, precomputed lookup tables and polynomial approximations can be useful).

When q is small and the precision is large (which occurs in the first few terms of the HRR series), it is more efficient to exploit the fact that $\alpha_{p,q}$ (or $\beta_{p,q}$) is an algebraic number, meaning that we can find an annihilating polynomial $A_{p,q} \in \mathbb{Q}[x]$ such that $A_{p,q}(\alpha_{p,q}) = 0$ (or $B_{p,q} \in \mathbb{Q}[x]$ such that $B_{p,q}(\beta_{p,q}) = 0$).

Knowing an annihilating polynomial $A_{p,q}$, we can compute a low-precision approximation x_0 of $\alpha_{p,q}$ via the cosine function, and refine it to high precision using the Newton iteration $x_{k+1} = x_k - A_{p,q}(x_k)/A'_{p,q}(x_k)$, and analogously for $\beta_{p,q}$. The low-precision approximation only roughly needs to be accurate enough to isolate the desired root from the other roots of the annihilating polynomial.

If p and q are fixed, computing $\alpha_{p,q}$ or $\beta_{p,q}$ to a precision of b bits using Newton iteration costs $O(M(b))$, provided that the Newton iteration is performed with a precision that doubles in each iteration. In other words, we gain a factor $\log b$ compared to evaluating the cosine or exponential using the arithmetic-geometric mean, or a factor $\log^2 b$ compared to evaluating the cosine or exponential using binary splitting.

Since the complexity grows with q , this trick alone does not appear to give a way to

reduce the asymptotic complexity for computing $p(n)$, but it does reduce the running time by a significant factor in practice.

5.5.1 Choice of annihilating polynomial

If we choose to compute $\beta_{p,q}$, then $B_{p,q}(x) = x^q + (-1)^{p+1}$ is a natural choice as annihilating polynomial. This polynomial can be evaluated using $O(\log q)$ multiplications with the binary exponentiation algorithm, so the cost grows quite slowly with q .

We could alternatively use the minimal polynomial of $\beta_{p,q}$, which is a cyclotomic polynomial. Although this lowers the polynomial degree, it appears to be less efficient since sparsity is lost.

Working with $\beta_{p,q}$ has the drawback that complex arithmetic is more costly than real arithmetic. If we compute $\alpha_{p,q}$, a possible choice as annihilating polynomial is the Chebyshev polynomial $U_{q-1}(x)$. This polynomial has degree $q - 1$ and coefficients with $O(q \log q)$ bits, so the cost to evaluate the polynomial in expanded form grows rapidly with q . However, $U_{q-1}(x)$ can also be evaluated using $O(\log q)$ multiplications with a binary exponentiation-like scheme based on the composition theorem for Chebyshev polynomials [65]. This algorithm could be competitive with complex exponentiation for large q .

Finally, we might take $A_{p,q}$ to be the minimal polynomial of $\alpha_{p,q}$. Although this polynomial is dense, its evaluation can be accelerated using the rectangular Paterson-Stockmeyer scheme. In fact, since Newton iteration requires the values $A_{p,q}(x_k)$ and $A'_{p,q}(x_k)$, we can save time by reusing the table of powers of x_k for both.

For the range of q and precisions encountered when computing the partition function, we experimentally tested both the real minimal polynomial and complex exponentiation, finding the former to be more efficient.

5.5.2 Generating trigonometric minimal polynomials

The minimal polynomial of $\alpha_{p,q}$ (which is a factor of $U_{q-1}(x)$) can be generated quite easily thanks to the results of Watkins and Zeitlin [108]. Let $\Psi_n(x) \in \mathbb{Q}[x]$ denote the monic minimal polynomial of $\cos(2\pi/n)$. It is clear that the monic minimal polynomial of $\alpha_{p,q}$ equals $\Psi_q(x)$ if p is even and $\Psi_{2q}(x)$ if p is odd. Watkins and Zeitlin show that

$$d(n) = \deg \Psi_n(x) = \begin{cases} 1, & \text{if } n = 1, 2 \\ \varphi(n)/2, & \text{if } n \geq 3 \end{cases}$$

and that the roots $(\alpha_i)_{i=1}^{d(n)}$ of $\Psi_n(x)$ when $n \geq 3$ precisely are $\cos(2\pi k/n)$ where $0 \leq k \leq \lfloor n/2 \rfloor$ and $\gcd(k, n) = 1$.

Since $2 \cos(2\pi/n)$ is an algebraic integer [71], we also note that $2^{d(n)}\Psi_n(x) \in \mathbb{Z}[x]$. We can therefore construct $\Psi_n(x)$ efficiently by computing numerical approximations of $(\alpha_i)_{i=1}^{d(n)}$, expanding $2^{d(n)} \prod_{i=1}^d (x - \alpha_i)$ using binary splitting, and rounding each coefficient to the nearest integer. Since $|\alpha_i| \leq 1$, binomial expansion shows that the integer coefficients are bounded by

$$2^{d(n)} \binom{n}{\lfloor n/2 \rfloor} \leq 2^{d(n)+n},$$

which allows us to estimate the required numerical precision.

An algebraic method to compute Ψ_n is also given in [108]. If $T_n(x)$ denotes a Chebyshev polynomial of the first kind and

$$F(n) = \begin{cases} 2^{-s} (T_{s+1}(x) - T_s(x)) & \text{if } n = 2s + 1 \text{ is odd} \\ 2^{-s} (T_{s+1}(x) - T_{s-1}(x)) & \text{if } n = 2s \text{ is even,} \end{cases}$$

then $\Psi_n(x) = \prod_{d|n} F(n/d)^{\mu(d)}$. We note that $T_n(x)$ can be generated very cheaply as a hypergeometric series. In fact, it is more efficient to compute the minimal polynomial of $2 \cos(2\pi/n)$. Then one computes the corresponding product of the rescaled polynomials

$$F(n) = \begin{cases} 2(T_{s+1}(x/2) - T_s(x/2)) & \text{if } n = 2s + 1 \text{ is odd} \\ 2(T_{s+1}(x/2) - T_{s-1}(x/2)) & \text{if } n = 2s \text{ is even.} \end{cases}$$

which are monic and have integer coefficients. In particular, polynomials with integer coefficients can be used throughout. Experiments suggest that the algebraic method with integer polynomial arithmetic usually is faster than the numerical binary splitting algorithm.

5.5.3 Performance evaluation

Figure 5.2 compares the running time of four different algorithms for computing $\cos(\pi/q)$ at a precision of 10^5 digits. The first algorithm is to evaluate the cosine function using the MPFR library (which internally uses the bit-burst binary splitting algorithm at this precision). The cost is essentially constant, and in fact decreases slightly with larger q since convergence of the Taylor series becomes more rapid.

The next two algorithms compute a root of the minimal polynomial $\Psi_n(x)$ using real Newton iteration. When the polynomial evaluations are done using Horner's rule, the cost grows roughly as $O(q)$ (as long as q is much smaller than the precision so that the coefficients of the minimal polynomial can be considered to have constant bit size).

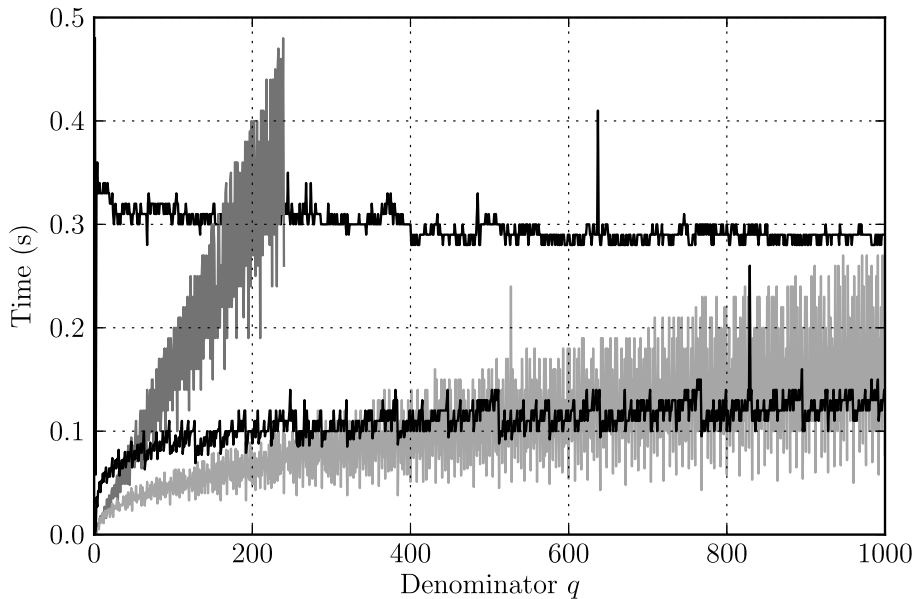


Figure 5.2: Time to evaluate $\cos(\pi/q)$ to a precision of 10^5 digits, shown as a function of the denominator q , using: (black, top) the MPFR cosine function; (dark gray) real Newton iteration with the minimal polynomial and Horner's rule; (light gray) real Newton iteration with the minimal polynomial and rectangular splitting; (black, bottom) complex Newton iteration with the polynomial $z^q + 1$.

This algorithm becomes slower than the MPFR cosine function around $q \approx 200$. The Paterson-Stockmeyer rectangular splitting algorithm is much more efficient, with a cost growing like $O(q^{1/2})$ for small q . It becomes slower than the MPFR cosine function only for $q > 1000$. The pseudorandom local variation in the running times is explained by the variation of $\deg \Psi_n$ (the three visible extra large spikes are runtime noise).

The fourth algorithm, complex Newton iteration, remains faster than direct cosine evaluation for q even larger than the previous two algorithms. Up to $q \approx 500$, however, the overhead of the complex arithmetic makes this algorithm slower than the minimal polynomial algorithm with rectangular splitting. The local variation in the running time is due to the fact that the cost of computing z^q depends on the sparseness of the binary representation of q (powers of two are cheapest).

As a last remark, we note that the computation of $\alpha_{p,q}$ or $\beta_{p,q}$ can be optimized for special values. If $q = m \times 2^k$ where m is odd, we can apply the trigonometric double-angle formulas k times to reduce the denominator to m . This reduces the degree of the minimal polynomial at the expense of performing k square roots. When $m = 1$, $m = 3$ or $m = 5$ (among other values) we can reduce the evaluation entirely to square roots.

5.6 Implementation

As reported in our original paper [52], we have implemented computation of $p(n)$ via the Hardy-Ramanujan-Rademacher formula in the FLINT library. Our implementation uses the MPFR library for high-precision floating-point arithmetic, hardware double-precision floating-point arithmetic for low-precision numerical evaluation, and FLINT functions for operations on word-size integers (to factor the $A_k(n)$ sums).

For efficiency, it is crucial to set the numerical precision tightly throughout the algorithm. It is easy to implement the HRR formula incorrectly. Past versions of both the Pari/GP and Maple computer algebra systems have computed incorrect values of $p(n)$ for some n (sequence A110375 in [83] lists n where versions 9.5 to 12 of Maple returns the wrong result, starting with $n = 11269, 11566, \dots$).

Based on a detailed (but incomplete) floating-point error analysis, and extensive testing, we concluded that the implementation in FLINT *probably* is correct for all n .

We have subsequently written a new implementation as part of the Arb library. This implementation uses ball arithmetic for all numerical calculations to guarantee that the final value of $p(n)$ is correct. The Arb implementation is slightly slower for small n mainly as a result of not using hardware double arithmetic (this optimization could be re-enabled on platforms where error bounds for the floating-point arithmetic, in particular the evaluation of transcendental functions, can be guaranteed). For very large n , the Arb implementation is slightly faster, mainly as a result of improved code for high-precision roots of unity.

Using a system with an AMD Opteron 6174 processor and 256 GiB RAM, we computed $p(10^{17})$, $p(10^{18})$ and $p(10^{19})$ with the FLINT implementation. The last computation took just less than 100 hours and used more than 150 GiB of memory, producing a result with over 11 billion bits. Some large values of $p(n)$ are listed in Table 5.1.

A comparison of the partition function implementations in several systems is shown in Figure 5.3. The FLINT and Arb implementations exhibit a time complexity only slightly higher than $O(n^{1/2})$, with a comparatively small constant factor. The Sage implementation written by J. Bober is fairly efficient for small n but has a complexity closer to $O(n)$, and is limited to arguments $n < 2^{32} \approx 4 \times 10^9$.

The partition function in Mathematica appears to have complexity slightly higher than $O(n^{1/2})$ as well, but consistently runs more than 100 times slower than our implementation. Based on extrapolation, computing $p(10^{19})$ would take several years. It is unclear whether Mathematica is actually using a nearly-optimal algorithm or whether the slow growth is just the manifestation of various overheads dwarfing the true asymptotic behavior. The ratio compared to our implementation appears too large to be explained by

n	Decimal expansion	Number of digits	Terms	Error
10^{12}	6129000962... 6867626906	1,113,996	264,526	2×10^{-7}
10^{13}	5714414687... 4630811575	3,522,791	787,010	3×10^{-8}
10^{14}	2750960597... 5564896497	11,140,072	2,350,465	-1×10^{-8}
10^{15}	1365537729... 3764670692	35,228,031	7,043,140	-3×10^{-9}
10^{16}	9129131390... 3100706231	111,400,846	21,166,305	-9×10^{-10}
10^{17}	8291300791... 3197824756	352,280,442	63,775,038	5×10^{-10}
10^{18}	1478700310... 1701612189	1,114,008,610	192,605,341	4×10^{-10}
10^{19}	5646928403... 3674631046	3,522,804,578	582,909,398	4×10^{-11}

Table 5.1: Large values of $p(n)$. The table also lists the number of terms N in the Hardy-Ramanujan-Rademacher formula used by FLINT (theoretically bounding the error by 0.25) and the difference between the floating-point sum and the rounded integer.

differences in performance of the underlying arithmetic alone; for example, evaluating the first term in the series for $p(10^{10})$ to required precision in Mathematica only takes about one second.

We get one external benchmark from ([14], p. 250), where it is reported that R. Crandall computed $p(10^9)$ in three seconds on a laptop in December 2008, “using the Hardy-Ramanujan-Rademacher ‘finite’ series for $p(n)$ along with FFT methods”. Even accounting for possible hardware differences, this appears to be an order of magnitude slower than our implementation.

In the FLINT and Arb implementations, about 30% to 50% of the total time is spent evaluating the first term in the Hardy-Ramanujan-Series for large n . Our implementations are therefore nearly optimal in a practical sense, since the first term in the Hardy-Ramanujan-Rademacher expansion hardly can be avoided and at most a factor around two can be gained by improving the tail evaluation.

Naturally, there is some potential to implement a faster version of the exponential function than the one provided by MPFR, reducing the cost of the first term. Improvements on the level of bignum multiplication would, on the other hand, presumably have a comparatively uniform effect.

By similar reasoning, at most a factor two can be gained through parallelization of our implementation by assigning terms in the Hardy-Ramanujan-Rademacher sum to separate threads. Further speedup on a multicore system requires parallelized versions of lower level routines, such as the exponential function or bignum multiplication. Fortunately, it is likely to be more interesting in practice to be able to evaluate $p(n)$ for a range of large values than just for a single value, and this task naturally parallelizes well.

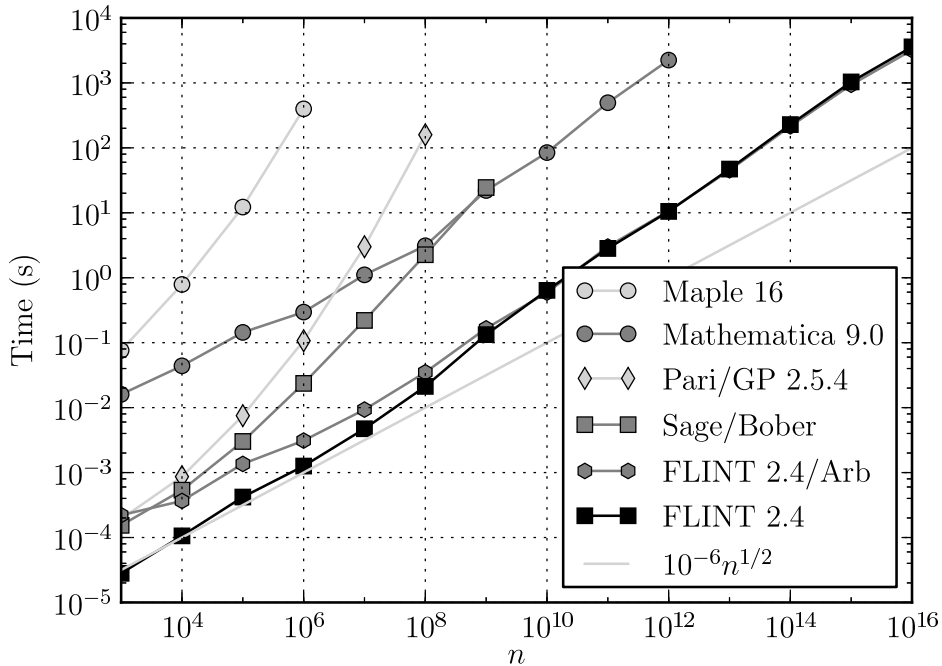


Figure 5.3: CPU time in seconds for computing $p(n)$ using various implementations. The graph of $t = 10^{-6}n^{1/2}$ is included to indicate the slope of an idealized algorithm satisfying the trivial lower complexity bound $\Omega(n^{1/2})$ (the offset 10^{-6} is arbitrary).

5.7 Multi-evaluation and congruence generation

One of the central problems concerning the partition function is the distribution of values of $p(n) \bmod m$. In 2000, Ono [85] proved that for every prime $m \geq 5$, there exist infinitely many congruences of the type

$$p(Ak + B) \equiv 0 \pmod{m} \tag{5.7.1}$$

where A, B are fixed and k ranges over all nonnegative integers. Ono's proof is nonconstructive, but Weaver [109] subsequently gave an algorithm for finding congruences of this type when $m \in \{13, 17, 19, 23, 29, 31\}$, and used the algorithm to compute 76,065 explicit congruences.

Weaver's congruences are specified by a tuple (m, ℓ, ε) where ℓ is a prime and $\varepsilon \in \{-1, 0, 1\}$, where we unify the notation by writing $(m, \ell, 0)$ in place of Weaver's (m, ℓ) .

Such a tuple corresponds to a family of congruences of the form (5.7.1) with coefficients

$$A = m\ell^{4-|\varepsilon|} \quad (5.7.2)$$

$$B = \frac{m\ell^{3-|\varepsilon|}\alpha + 1}{24} + m\ell^{3-|\varepsilon|}\delta, \quad (5.7.3)$$

where α is the unique solution of $m\ell^{3-|\varepsilon|}\alpha \equiv -1 \pmod{24}$ with $1 \leq \alpha < 24$, and where $0 \leq \delta < \ell$ is any solution of

$$\begin{cases} 24\delta \not\equiv -\alpha \pmod{\ell} & \text{if } \varepsilon = 0 \\ (24\delta + \alpha \mid \ell) = \varepsilon & \text{if } \varepsilon = \pm 1. \end{cases} \quad (5.7.4)$$

The free choice of δ gives $\ell - 1$ distinct congruences for a given tuple (m, ℓ, ε) if $\varepsilon = 0$, and $(\ell - 1)/2$ congruences if $\varepsilon = \pm 1$.

5.7.1 Weaver's algorithm

Weaver's test for congruence, described by Theorems 7 and 8 in [109], essentially amounts to a single evaluation of $p(n)$ at a special point n . Namely, for given m, ℓ , we compute the smallest solutions of $\delta_m \equiv 24^{-1} \pmod{m}$, $r_m \equiv -m \pmod{24}$, and check whether $p(mr_m(\ell^2 - 1)/24 + \delta_m)$ is congruent mod m to one of three values corresponding to the parameter $\varepsilon \in \{-1, 0, 1\}$. We give a compact statement of this procedure as Algorithm 5.7.1. To find new congruences, we simply perform a brute force search over a set of candidate primes ℓ , calling Algorithm 5.7.1 repeatedly.

Algorithm 5.7.1 Weaver's congruence test

Input: A pair of prime numbers $13 \leq m \leq 31$ and $\ell \geq 5$, $m \neq \ell$

Output: (m, ℓ, ε) defining a congruence, and Not-a-congruence otherwise

```

 $\delta_m \leftarrow 24^{-1} \pmod{m}$  ▷ Reduced to  $0 \leq \delta_m < m$ 
 $r_m \leftarrow (-m) \pmod{24}$  ▷ Reduced to  $0 \leq m < 24$ 
 $v \leftarrow \frac{m-3}{2}$ 
 $x \leftarrow p(\delta_m)$  ▷ We have  $x \not\equiv 0 \pmod{m}$ 
 $y \leftarrow p(m(r_m(\ell^2 - 1)/24) + \delta_m)$ 
 $f \leftarrow (3 \mid \ell) ((-1)^v r_m \mid \ell)$  ▷ Jacobi symbols
 $t \leftarrow y + fx\ell^{v-1}$ 
if  $t \equiv \omega \pmod{m}$  where  $\omega \in \{-1, 0, 1\}$  then
    return  $(m, \ell, \omega(3(-1)^v \mid \ell))$ 
else
    return Not-a-congruence

```

n	Series ($\mathbb{Z}/13\mathbb{Z}$)	Series (\mathbb{Z})	HRR (all)	HRR (sparse)
10^4	0.01 s	0.1 s	1.4 s	0.001 s
10^5	0.13 s	4.1 s	41 s	0.008 s
10^6	1.4 s	183 s	1430 s	0.08 s
10^7	14 s			0.7 s
10^8	173 s			8 s
10^9	2507 s			85 s

Table 5.2: Comparison of time needed to compute multiple values of $p(n)$ up to the given bound, using power series inversion and the Hardy-Ramanujan-Rademacher formula. The rightmost column gives the time when only computing the subset of terms that are searched with Weaver’s algorithm in the $m = 13$ case.

5.7.2 Comparison of algorithms for multi-evaluation

A timing comparison between various methods for multi-evaluation of $p(n)$ is shown in Table 5.2. Power series division, which we implemented using FLINT over both \mathbb{Z} and $\mathbb{Z}/m\mathbb{Z}$, is clearly the best choice for computing all values up to n modulo a fixed prime, having a complexity of $O(n^{1+o(1)})$. For computing the full integer values, the power series and HRR methods both have complexity $O(n^{3/2+o(1)})$, with the power series method expectedly winning.

Ignoring logarithmic factors, we can expect the HRR formula to be better than the power series for multi-evaluation of $p(n)$ up to some bound n when n/c values are needed. The factor $c \approx 10$ in the FLINT implementation is a remarkable improvement over $c \approx 1000$ attainable with previous implementations of the partition function. For evaluation mod m , the HRR formula is competitive when $O(n^{1/2})$ values are needed; in this case, the constant is highly sensitive to m .

For the sparse subset of $O(n^{1/2})$ terms searched with Weaver’s algorithm, the HRR formula has the same complexity as the modular power series method, but as seen in Table 5.2 runs more than an order of magnitude faster. On top of this, it has the advantage of parallelizing trivially, being resumable from any point, and requiring very little memory (the power series evaluation mod $m = 13$ up to $n = 10^9$ required over 40 GiB memory, compared to a few megabytes with the HRR formula). The recursive version of Euler’s pentagonal number theorem is, of course, also resumable from an arbitrary point, but this requires computing and storing all previous values.

We mention that the authors of [80] use a parallel version of the recursive Euler method. This is not as efficient as power series inversion, but allows the computation to be split across multiple processors more easily.

m	$(m, \ell, 0)$	$(m, \ell, +1)$	$(m, \ell, -1)$	Congruences	CPU	Max n
13	6,189	6,000	6,132	5,857,728,831	448 h	5.9×10^{12}
17	4,611	4,611	4,615	4,443,031,844	391 h	4.9×10^{12}
19	4,114	4,153	4,152	3,966,125,921	370 h	3.9×10^{12}
23	3,354	3,342	3,461	3,241,703,585	125 h	9.5×10^{11}
29	2,680	2,777	2,734	2,629,279,740	1,155 h	2.2×10^{13}
31	2,428	2,484	2,522	2,336,738,093	972 h	2.1×10^{13}
All	23,376	23,367	23,616	22,474,608,014	3,461 h	

Table 5.3: The number of tuples of the given type with $\ell < 10^6$, the total number of congruences defined by these tuples, the total CPU time, and the approximate bound up to which $p(n)$ was evaluated.

5.7.3 Results

Weaver gives 167 tuples, or 76,065 congruences, containing all ℓ up to approximately 1,000–3,000 (depending on m). This table was generated by computing all values of $p(n)$ with $n < 7.5 \times 10^6$ using the recursive version of Euler’s pentagonal theorem. Computing Weaver’s table from scratch with our implementation of the HRR formula, evaluating only the necessary n , takes just a few seconds. We are also able to numerically verify instances of all entries in Weaver’s table for small k .

As a more substantial exercise, we extend Weaver’s table by determining all ℓ up to 10^6 for each prime m . Statistics are listed in Table 5.3. The computation was performed using the FLINT partition function by assigning subsets of the search space to separate processes, running on between 40 and 48 active cores for a period of four days, evaluating $p(n)$ at $6(\pi(10^6) - 3) = 470,970$ distinct n ranging up to 2×10^{13} .

We find a total of 70,359 tuples, corresponding to slightly more than 2.2×10^{10} new congruences. To pick an arbitrary, concrete example, one “small” new congruence is $(13, 3797, -1)$ with $\delta = 2588$, giving

$$p(711647853449k + 485138482133) \equiv 0 \pmod{13}$$

which we easily evaluate for all $k \leq 100$, providing a sanity check on the identity as well as the partition function implementation. As a larger example, $(29, 999959, 0)$ with $\delta = 999958$ gives

$$p(28995244292486005245947069k + 28995221336976431135321047) \equiv 0 \pmod{29}$$

which, despite our efforts, presently is out of reach for direct evaluation.

5.8 Remarks

The techniques we have presented could potentially be applied to other HRR-type series, such as the formula for the number of partitions into distinct parts

$$Q(n) = \frac{\pi^2 \sqrt{2}}{24} \sum_{k=1}^{\infty} \frac{\tilde{A}_{2k-1}(n)}{(1-2k)^2} {}_0F_1 \left(2, \frac{(n + \frac{1}{24})\pi^2}{12(1-2k)^2} \right). \quad (5.8.1)$$

However, the exponential sum $\tilde{A}_{2k-1}(n)$ appearing here is not the same as the exponential sum for the ordinary partition function, and it is not obvious that a similar factorization exists.

It remains an open problem to find a fast way to compute the isolated value $p(n)$ using purely algebraic methods. The recent discovery of an “algebraic formula” for $p(n)$ by Bruinier and Ono [28] could perhaps lead to such an algorithm. An effort in this direction is made by Bruinier, Ono and Sutherland in the paper [29].

Appendix A

Implementations

The Fast Library for Number Theory (FLINT) [47] is a C library for computational number theory, with emphasis on fast arithmetic in various base rings. In particular, FLINT provides asymptotically fast polynomial arithmetic in $\mathbb{Z}[x]$, $\mathbb{Q}[x]$ and $(\mathbb{Z}/n\mathbb{Z})[x]$. Polynomial multiplication uses a Schönhage-Strassen FFT implementation by William Hart. The underlying integer arithmetic is handled by GMP [41] or MPIR [79].

FLINT was initially developed around 2007 by William Hart and David Harvey. The library was later rewritten from scratch and substantially extended by Hart, Sebastian Pancratz, and the present author (whose contributions include many of the methods for power series and linear algebra). Several other authors have also contributed to the project.

For the computations over \mathbb{R} and \mathbb{C} described in this thesis, the author developed the Arb library [53] as an extension of the FLINT project. Both FLINT and Arb are free software, available under version 2 or later of the GNU General Public License. Complementary to this thesis which describes the mathematical and algorithmic background, the Arb library comes with extensive documentation (114 pages in printable format as of Arb version 1.0.0) describing the interface and implementation-level details.

Arb uses ball arithmetic [103] (or mid-rad interval arithmetic) to do numerical computations with rigorous propagation of error bounds. In contrast to endpoint-based interval arithmetic, this representation has negligible space and time overhead at high precision.

The `fmprb_t` type represents a real ball as a pair of floating-point numbers of type `fmpr_t`. Unlike some other implementations of arbitrary-precision floating-point arithmetic, the `fmpr_t` type allocates space for the mantissa dynamically, instead of zero-padding up to the working precision. This feature is especially convenient for binary splitting, where full-precision numbers are grown gradually from smaller integers.

Evaluation of elementary functions uses correctly-rounded floating-point functions provided by the MPFR library [40], along with error propagation based on derivatives. In some cases, custom implementations are used to improve performance. For instance, constants such as π and $\log 2$ are computed using binary splitting code which works for generic hypergeometric series and automatically computes error bounds. Higher transcendental functions (including the gamma and zeta functions) are implemented using algorithms described in this thesis.

The `fmpcb_t` type represents complex numbers in Cartesian form as pairs of `fmpcb_t` values. Thus complex numbers are not strictly speaking represented as “complex balls” but as rectangular boxes. Similarly, the `fmpcb_poly_t` and `fmpcb_mat_t` types represent polynomials (or power series) over the real or complex numbers, and the `fmpcb_mat_t` and `fmpcb_mat_t` types represent matrices, implemented as arrays of `fmpcb_t` or `fmpcb_t` coefficients. This representation is sometimes less efficient than using complex, polynomial or matricial balls with a single common error bound (tradeoffs are discussed in [103]), but we generally prefer it since it often is convenient to track error bounds accurately for individual coefficients.

Fast and numerically stable multiplication in $\mathbb{R}[x]$ and $\mathbb{C}[x]$ is implemented by breaking polynomials into segments with similarly-sized coefficients and using FLINT to compute the subproducts exactly in $\mathbb{Z}[x]$ (a simplified version of van der Hoeven’s block multiplication algorithm [102]). Asymptotically fast polynomial and power series division is implemented using Newton iteration. Elementary functions of power series switch between basecase algorithms and asymptotically fast Newton-based algorithms.

Higher-level features provided in the Arb library include isolation and high-precision polishing of roots of real analytic functions, isolation of complex roots of polynomials, and numerical integration of complex-valued functions using Taylor series (all operations are done with rigorous error bounds).

The following C program demonstrates use of the Arb library. It exactly computes the Bell number B_n by numerically calculating the series expansion of $\exp(\exp(x) - 1) = \sum_{n=0}^{\infty} B_n x^n / n!$ to order $O(x^{n+1})$ and reading off the last coefficient.

At too low precision, rounding a floating-point approximation of B_n to the nearest integer would give the wrong result. The program therefore doubles the working precision until the ball for B_n contains a single integer, which necessarily must be the correct one.

With a simple modification, the program could output all the numbers B_0, \dots, B_{n-1} at no significant extra cost. It could also be modified to output numerical approximations accurate to a requested number of digits, rather than the exact values.

```

#include "fmprb_poly.h"

int main()
{
    fmprb_poly_t f;
    fmprb_t t, u;
    fmpz_t z;
    long i, prec, n = 1000;

    fmprb_poly_init(f);
    fmprb_init(t);
    fmprb_init(u);
    fmpz_init(z);

    for (prec = 100; ; prec *= 2)
    {
        fmprb_poly_zero(f);
        fmprb_poly_set_coeff_si(f, 1, 1); /* f = x */
        fmprb_poly_exp_series(f, f, n + 1, prec); /* f = e^x */
        fmprb_poly_set_coeff_si(f, 0, 0); /* f = e^x - 1 */
        fmprb_poly_exp_series(f, f, n + 1, prec); /* f = e^(e^x-1) */

        fmprb_poly_get_coeff_fmprb(t, f, n); /* t = B_n / n! */
        fmprb_fac_ui(u, n, prec);
        fmprb_mul(t, t, u, prec); /* t = B_n */

        fmprb_printd(t, 15); printf("\n");

        if (fmprb_get_unique_fmpz(z, t)) /* we get the right integer */
            break;
    }

    fmpz_print(z); printf("\n");

    fmprb_poly_clear(f);
    fmprb_clear(t);
    fmprb_clear(u);
    fmpz_clear(z);
}

```


Bibliography

- [1] A. V. Aho, K. Steiglitz, and J. D. Ullman. Evaluating polynomials at fixed sets of points. *SIAM Journal on Computing*, 4(4):533–539, 1975.
- [2] N. Ankeny. The least quadratic non residue. *Annals of Mathematics*, 55(1):65–72, 1952.
- [3] T. Apostol. *Modular functions and Dirichlet series in number theory*. Springer, New York, second edition, 1997.
- [4] J. Arias de Reyna. Asymptotics of Keiper-Li coefficients. *Functiones et Approximatio Commentarii Mathematici*, 45(1):7–21, 2011.
- [5] D. H. Bailey and J. M. Borwein. Experimental mathematics: recent developments and future outlook. In B. Engquist, W. Schmid, and P. W. Michor, editors, *Mathematics Unlimited – 2001 and Beyond*, pages 51–66. Springer, 2000.
- [6] D. J. Bernstein. Composing power series over a finite ring in essentially linear time. *Journal of Symbolic Computation*, 26(3):339–342, 1998.
- [7] D. J. Bernstein. Multidigit multiplication for mathematicians. <http://cr.yp.to/papers.html#m3>, 2001.
- [8] D. J. Bernstein. Removing redundancy in high-precision Newton iteration. <http://cr.yp.to/papers.html#fastnewton>, 2004.
- [9] D. J. Bernstein. Fast multiplication and its applications. *Algorithmic Number Theory*, 44:325–384, 2008.
- [10] R. Bloemen. Even faster $\zeta(2n)$ calculation!, 2009. <http://remcobloemen.nl/2009/11/even-faster-zeta-calculation.html>.
- [11] A. I. Bogolubsky and S. L. Skorokhodov. Fast evaluation of the hypergeometric function ${}_pF_{p-1}(a; b; z)$ at the singular point $z = 1$ by means of the Hurwitz zeta function $\zeta(\alpha, s)$. *Programming and Computer Software*, 32(3):145–153, 2006.

- [12] J. Bohman and C-E. Fröberg. The Stieltjes function – definition and properties. *Mathematics of Computation*, 51(183):281–289, 1988.
- [13] J. M. Borwein and P. B. Borwein. *Pi and the AGM: A study in analytic number theory and computational complexity*. Wiley, Canada, 1987.
- [14] J. M. Borwein and P. B. Borwein. *Experimental and computational mathematics: selected writings*. Perfectly Scientific Press, Portland, OR, 2010.
- [15] J. M. Borwein, D. M. Bradley, and R. E. Crandall. Computational strategies for the Riemann zeta function. *Journal of Computational and Applied Mathematics*, 121:247–296, 2000.
- [16] P. B. Borwein. Reduced complexity evaluation of hypergeometric functions. *Journal of Approximation Theory*, 50(3), July 1987.
- [17] P. B. Borwein. An efficient algorithm for the Riemann zeta function. *Canadian Mathematical Society Conference Proceedings*, 27:29–34, 2000.
- [18] A. Bostan, P. Gaudry, and É. Schost. Linear recurrences with polynomial coefficients and application to integer factorization and Cartier-Manin operator. *SIAM Journal on Computing*, 36(6):1777–1806, 2007.
- [19] A. Bostan and M. Kauers. Automatic classification of restricted lattice walks. *DMTCS Proceedings, FPSAC 2009*, (01):201–215, 2009.
- [20] A. Bostan, B. Salvy, and É. Schost. Power series composition and change of basis. In *Proceedings of the twenty-first international symposium on symbolic and algebraic computation*, pages 269–276. ACM, 2008.
- [21] R. P. Brent. The complexity of multiple-precision arithmetic. *The Complexity of Computational Problem Solving*, pages 126–165, 1976.
- [22] R. P. Brent. Fast multiple-precision evaluation of elementary functions. *Journal of the ACM*, 23(2):242–251, 1976.
- [23] R. P. Brent and D. Harvey. Fast computation of Bernoulli, tangent and secant numbers, 2011. <http://arxiv.org/abs/1108.0286>.
- [24] R. P. Brent and F. Johansson. A bound for the error term in the Brent-McMillan algorithm. *Mathematics of Computation*, to appear.
- [25] R. P. Brent and H. T. Kung. Fast algorithms for manipulating formal power series. *Journal of the ACM*, 25(4):581–595, 1978.
- [26] R. P. Brent and E. M. McMillan. Some new algorithms for high-precision computation of Euler’s constant. *Mathematics of Computation*, 34(149):305–312, 1980.

- [27] R. P. Brent and P. Zimmermann. *Modern Computer Arithmetic*. Cambridge University Press, 2011.
- [28] J. H. Bruinier and K. Ono. Algebraic formulas for the coefficients of half-integral weight harmonic weak Maass forms. <http://arxiv.org/abs/1104.1182>, 2011.
- [29] J. H. Bruinier, K. Ono, and A. V. Sutherland. Class polynomials for nonholomorphic modular functions. <http://arxiv.org/abs/1301.5672>, 2013.
- [30] D. G. Cantor and E. Kaltofen. On fast multiplication of polynomials over arbitrary algebras. *Acta Informatica*, 28(7):693–701, 1991.
- [31] D. V. Chudnovsky and G. V. Chudnovsky. Approximations and complex multiplication according to Ramanujan. In *Ramanujan Revisited*, pages 375–472. Academic Press, 1988.
- [32] D. V. Chudnovsky and G. V. Chudnovsky. Computer algebra in the service of mathematical physics and number theory. In *Computers in Mathematics*, pages 109–232. Dekker, New York, 1990.
- [33] M. Cipolla. Un metodo per la risoluzione della congruenza di secondo grado. *Napoli Rend.*, 9:153–163, 1903.
- [34] M. W. Coffey. An efficient algorithm for the Hurwitz zeta and related functions. *Journal of Computational and Applied Mathematics*, 225(2):338–346, 2009.
- [35] R. Crandall and C. Pomerance. *Prime Numbers: A Computational Perspective*. Springer, New York, 2005.
- [36] H. M. Edwards. *Riemann’s zeta function*. Academic Press, 1974.
- [37] P. Erdős. Remarks on number theory. I. *Mat. Lapok*, 12:10–17, 1961.
- [38] T. Finck, G. Heinig, and K. Rost. An inversion formula and fast algorithms for Cauchy-Vandermonde matrices. *Linear algebra and its applications*, 183:179–191, 1993.
- [39] P. Flajolet and I. Vardi. Zeta function expansions of classical constants. Unpublished manuscript, <http://algo.inria.fr/flajolet/Publications/landau.ps>, 1996.
- [40] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélicier, and P. Zimmermann. MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Transactions on Mathematical Software*, 33(2):13:1–13:15, June 2007. <http://mpfr.org>.

- [41] The GMP development team. GMP: The GNU multiple precision arithmetic library. <http://www.gmpilib.org>.
- [42] H. Gould. Series transformations for finding recurrences for sequences. *Fibonacci Quarterly*, 28:166–171, 1990.
- [43] X. Gourdon and P. Sebah. The Riemann Zeta-function $\zeta(s)$. <http://numbers.computation.free.fr/Constants/Miscellaneous/zeta.html>, 2003.
- [44] B. Haible and T. Papanikolaou. Fast multiprecision evaluation of series of rational numbers. In J. P. Buhler, editor, *Algorithmic Number Theory: Third International Symposium*, volume 1423, pages 338–350. Springer, 1998.
- [45] G. Hanrot and P. Zimmermann. Newton iteration revisited. <http://www.loria.fr/~zimmerma/papers/fastnewton.ps.gz>, 2004.
- [46] G. H. Hardy and S. Ramanujan. Asymptotic formulae in combinatory analysis. *Proceedings of the London Mathematical Society*, 17:75–115, 1918.
- [47] W. B. Hart. Fast Library for Number Theory: An Introduction. In *Proceedings of the Third international congress conference on Mathematical software*, ICMS’10, pages 88–91, Berlin, Heidelberg, 2010. Springer-Verlag. <http://flintlib.org>.
- [48] D. Harvey. Faster algorithms for the square root and reciprocal of power series. *Mathematics of Computation*, 80(273):387–394, 2011.
- [49] G. Hiary. Fast methods to compute the Riemann zeta function. *Annals of Mathematics*, 174:891–946, 2011.
- [50] J. Hopcroft and J. Musinski. Duality applied to the complexity of matrix multiplication and other bilinear forms. *SIAM Journal on Computing*, 2(3):159–173, 1973.
- [51] X. Huang and V. Pan. Fast rectangular matrix multiplication and applications. *Journal of Complexity*, 14(2):257–299, 1998.
- [52] F. Johansson. Efficient implementation of the Hardy-Ramanujan-Rademacher formula. *LMS Journal of Computation of Mathematics*, 15:341–359, 2012.
- [53] F. Johansson. Arb: a C library for ball arithmetic. *ACM Communications in Computer Algebra*, 47(4):166–169, 2013.
- [54] F. Johansson. Evaluating parametric holonomic sequences using rectangular splitting, submitted.
- [55] F. Johansson. Rigorous high-precision computation of the Hurwitz zeta function and its derivatives, submitted.

- [56] F. Johansson. A fast algorithm for reversion of power series. *Mathematics of Computation*, to appear.
- [57] E. A. Karatsuba. Fast computation of the Riemann zeta-function $\zeta(s)$ for integer values of the argument s . *Problems of Information Transmission*, 31(4):353–362, 1995.
- [58] E. A. Karatsuba. Fast evaluation of the Hurwitz zeta function and Dirichlet L -series. *Problems of Information Transmission*, 34(4):62–75, 1998.
- [59] K. S. Kedlaya and C. Umans. Fast polynomial factorization and modular composition. *SIAM Journal on Computing*, 40(6):1767–1802, 2011.
- [60] J. B. Keiper. Power series expansions of Riemann’s ξ function. *Mathematics of Computation*, 58(198):765–773, 1992.
- [61] C. Knessl and M. Coffey. An asymptotic form for the Stieltjes constants $\gamma_k(a)$ and for a sum $S_\gamma(n)$ appearing under the Li criterion. *Mathematics of Computation*, 80(276):2197–2217, 2011.
- [62] C. Knessl and M. Coffey. An effective asymptotic formula for the Stieltjes constants. *Mathematics of Computation*, 80(273):379–386, 2011.
- [63] D. Knuth. Notes on generalized Dedekind sums. *Acta Arithmetica*, 33:297–325, 1977.
- [64] D. Knuth. *The Art of Computer Programming, volume 2: Seminumerical Algorithms*. Addison-Wesley, 3rd edition, 1998.
- [65] W. Koepf. Efficient computation of Chebyshev polynomials in computer algebra. In M. J. Wester, editor, *Computer Algebra System: A Practical Guide*, pages 79–98. Wiley, 1999.
- [66] S. Köhler and M. Ziegler. On the stability of fast polynomial arithmetic. In *Proceedings of the 8th Conference on Real Numbers and Computers*, Santiago de Compostela, Spain, 2008.
- [67] R. Kreminski. Newton-Cotes integration for approximating Stieltjes (generalized Euler) constants. *Mathematics of Computation*, 72(243):1379–1397, 2003.
- [68] D. Lehmer. On a conjecture of Ramanujan. *Journal of the London Mathematical Society*, 11:114–118, 1936.
- [69] D. Lehmer. On the Hardy-Ramanujan series for the partition function. *Journal of the London Mathematical Society*, 3:171–176, 1937.

- [70] D. Lehmer. On the series for the partition function. *Transactions of the American Mathematical Society*, 43(2):271–295, 1938.
- [71] D. H. Lehmer. A note on trigonometric algebraic numbers. *The American Mathematical Monthly*, 40(3):165–166, 1933.
- [72] Xian-Jin Li. The positivity of a sequence of numbers and the Riemann Hypothesis. *Journal of Number Theory*, 65(2):325–333, 1997.
- [73] Y. Matiyasevich. An artless method for calculating approximate values of zeros of Riemann’s zeta function, 2012. <http://logic.pdmi.ras.ru/~yumat/personaljournal/artlessmethod/>.
- [74] Y. Matiyasevich and G. Beliakov. Zeroes of Riemann’s zeta function on the critical line with 20000 decimal digits accuracy, 2011. http://dro.deakin.edu.au/view/DU:30051725?print_friendly=true.
- [75] M. Mezzarobba. NumGfun: a package for numerical and analytic computation with D-finite functions. In W. Koepf, editor, *ISSAC ’10*, pages 139–146. ACM, 2010.
- [76] M. Mezzarobba. *Autour de l’évaluation numérique des fonctions D-finies*. PhD thesis, École polytechnique, 2011.
- [77] M. Mezzarobba. A note on the space complexity of fast D-finite function evaluation. In V. P. Gerdt, W. Koepf, E. W. Mayr, and E. V. Vorozhtsov, editors, *CASC 2012*, number 7442 in Lecture Notes in Computer Science, pages 212–223. Springer, 2012.
- [78] The MPFR team. The MPFR library: algorithms and proofs. <http://www.mpfr.org/algo.html>. retrieved 2013.
- [79] The MPIR development team. MPIR: Multiple Precision Integers and Rationals. <http://www.mpir.org>.
- [80] K. James E. Perez N. Calkin, J. Davis and C. Swannack. Computing the integer partition function. *Mathematics of Computation*, 76(259):1619–1638, 2007.
- [81] A. Odlyzko. Asymptotic enumeration methods. In R. Graham, M. Grötschel, and L. Lovász, editors, *Handbook of combinatorics 2*, pages 1063–1229. Elsevier, The Netherlands, 1995. <http://www.dtc.umn.edu/~odlyzko/doc/asymptotic.enum.pdf>.
- [82] A. M. Odlyzko and A. Schönhage. Fast algorithms for multiple evaluations of the Riemann zeta function. *Transactions of the American Mathematical Society*, 309(2):797–809, 1988.

- [83] OEIS Foundation. The On-Line Encyclopedia of Integer Sequences. <http://oeis.org>.
- [84] F. W. J. Olver. *Asymptotics and Special Functions*. A K Peters, Wellesley, MA, 1997.
- [85] K. Ono. The distribution of the partition function modulo m . *Annals of Mathematics*, 151:293–307, 2000.
- [86] Jr. P. Hagsis. A root of unity occurring in partition theory. *Proceedings of the American Mathematical Society*, 26(4):579–582, 1970.
- [87] M. S. Paterson and L. J. Stockmeyer. On the number of nonscalar multiplications necessary to evaluate polynomials. *SIAM Journal on Computing*, 2(1), March 1973.
- [88] Y.-F.S Pétermann and J-L. Rémy. Arbitrary precision error analysis for computing $\zeta(s)$ with the Cohen-Olivier algorithm: complete description of the real case and preliminary report on the general case. Rapport de recherche RR-5852, INRIA, 2006.
- [89] P. Pollack. The average least quadratic nonresidue modulo m and other variations on a theme of Erdős. *Journal of Number Theory*, 132(6):1185–1202, 2012.
- [90] R. L. Probert. On the additive complexity of matrix multiplication. *SIAM Journal on Computing*, 5(2):187–203, 1976.
- [91] H. Rademacher. On the partition function $p(n)$. *Proceedings of the London Mathematical Society*, 43:241–254, 1938.
- [92] H. Rademacher and A. Whiteman. Theorems on Dedekind sums. *American Journal of Mathematics*, 63(2):377–407, 1941.
- [93] P. Ritzmann. A fast numerical algorithm for the composition of power series with complex coefficients. *Theoretical Computer Science*, 44(1):1–16, 1986.
- [94] D. Shanks. Five number-theoretic algorithms. *Proceedings of the Second Manitoba Conference on Numerical Mathematics*, pages 51–70, 1972.
- [95] D. M. Smith. Efficient multiple-precision evaluation of elementary functions. *Mathematics of Computation*, 52:131–134, 1989.
- [96] D. M. Smith. Algorithm: Fortran 90 software for floating-point multiple precision arithmetic, gamma and related functions. *Transactions on Mathematical Software*, 27:377–387, 2001.

- [97] W. A. Stein et al. *Sage Mathematics Software*. The Sage Development Team, 2010. <http://www.sagemath.org>.
- [98] A. J. Stothers. *On the complexity of matrix multiplication*. PhD thesis, University of Edinburgh, 2010.
- [99] A. Tonelli. Bemerkung über die Auflösung quadratischer Congruenzen. *Göttinger Nachrichten*, pages 344–346, 1891.
- [100] J. van der Hoeven. Fast evaluation of holonomic functions. *Theoretical Computer Science*, 210:199–215, 1999.
- [101] J. van der Hoeven. Relax, but don't be too lazy. *Journal of Symbolic Computation*, 34(6):479–542, 2002.
- [102] J. van der Hoeven. Making fast multiplication of polynomials numerically stable. Technical Report 2008-02, Université Paris-Sud, Orsay, France, 2008.
- [103] J. van der Hoeven. Ball arithmetic. Technical report, HAL, 2009. <http://hal.archives-ouvertes.fr/hal-00432152/fr/>.
- [104] V. Vassilevska Williams. Breaking the Coppersmith-Winograd barrier. <http://cs.berkeley.edu/~virgi/matrixmult.pdf>, 2011.
- [105] L. Vepštas. An efficient algorithm for accelerating the convergence of oscillatory series, useful for computing the polylogarithm and Hurwitz zeta functions. *Numerical Algorithms*, 47(3):211–252, 2008.
- [106] J. von zur Gathen and J. Gerhard. Fast algorithms for Taylor shifts and certain difference equations. In *Proceedings of ISSAC'97*, pages 40–47, 1997.
- [107] J. von zur Gathen and J. Gerhard. *Modern Computer Algebra*. Cambridge University Press, 2nd edition, 2003.
- [108] W. Watkins and J. Zeitlin. The minimal polynomial of $\cos(2\pi/n)$. *The American Mathematical Monthly*, 100(5):471–474, 1993.
- [109] R. Weaver. New congruences for the partition function. *The Ramanujan Journal*, 5:53–63, 2001.
- [110] A. L. Whiteman. A sum connected with the series for the partition function. *Pacific Journal of Mathematics*, 6(1):159–176, 1956.
- [111] Wolfram Research. Some notes on internal implementation (section of the online documentation for Mathematica 9.0). <http://reference.wolfram.com/mathematica/tutorial/SomeNotesOnInternalImplementation.html>, 2013.

- [112] M. Ziegler. Fast (multi-)evaluation of linearly recurrent sequences: Improvements and applications. 2005. <http://arxiv.org/abs/cs/0511033>.

Eidesstattliche Erklärung

Ich erkläre an Eides statt, dass ich die vorliegende Dissertation selbstständig und ohne fremde Hilfe verfasst, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt bzw. die wörtlich oder sinngemäß entnommenen Stellen als solche kenntlich gemacht habe.

Die vorliegende Dissertation ist mit dem elektronisch übermittelten Textdokument identisch.

Linz, February 17, 2014

Fredrik Johansson

Curriculum vitae

Name Fredrik Johansson

Nationality Swedish

Date of birth October 23, 1985

Email fredrik.johansson@gmail.com

Website <http://fredrikj.net>

Education

2010–2014 Doctoral studies in Symbolic Computation, RISC, Johannes Kepler University, Linz.

2005–2010 M.Sc. and B.Sc. in Engineering Physics, Chalmers University of Technology, Gothenburg.

2004–2005 Studies in Software Engineering (first year of B.Sc. program), Chalmers University of Technology, Gothenburg.