

MODULAR EXPONENTIATION VIA THE EXPLICIT CHINESE REMAINDER THEOREM

DANIEL J. BERNSTEIN AND JONATHAN P. SORENSON

ABSTRACT. Fix pairwise coprime positive integers p_1, p_2, \dots, p_s . We propose representing integers u modulo m , where m is any positive integer up to roughly $\sqrt{p_1 p_2 \cdots p_s}$, as vectors $(u \bmod p_1, u \bmod p_2, \dots, u \bmod p_s)$. We use this representation to obtain a new result on the parallel complexity of modular exponentiation: there is an algorithm for the Common CRCW PRAM that, given positive integers x , e , and m in binary, of total bit length n , computes $x^e \bmod m$ in time $O(n/\lg \lg n)$ using $n^{O(1)}$ processors. For comparison, a parallelization of the standard binary algorithm takes superlinear time; Adleman and Kompella gave an $O((\lg n)^3)$ expected time algorithm using $\exp(O(\sqrt{n \lg n}))$ processors; von zur Gathen gave an NC algorithm for the highly special case that m is polynomially smooth.

1. INTRODUCTION

In this paper we consider the problem of computing $x^e \bmod m$ for large integers x , e , and m . This is the bottleneck in Rabin's algorithm for testing primality, the Diffie-Hellman algorithm for exchanging cryptographic keys, and many other common algorithms. See, e.g., [18, Section 4.5.4].

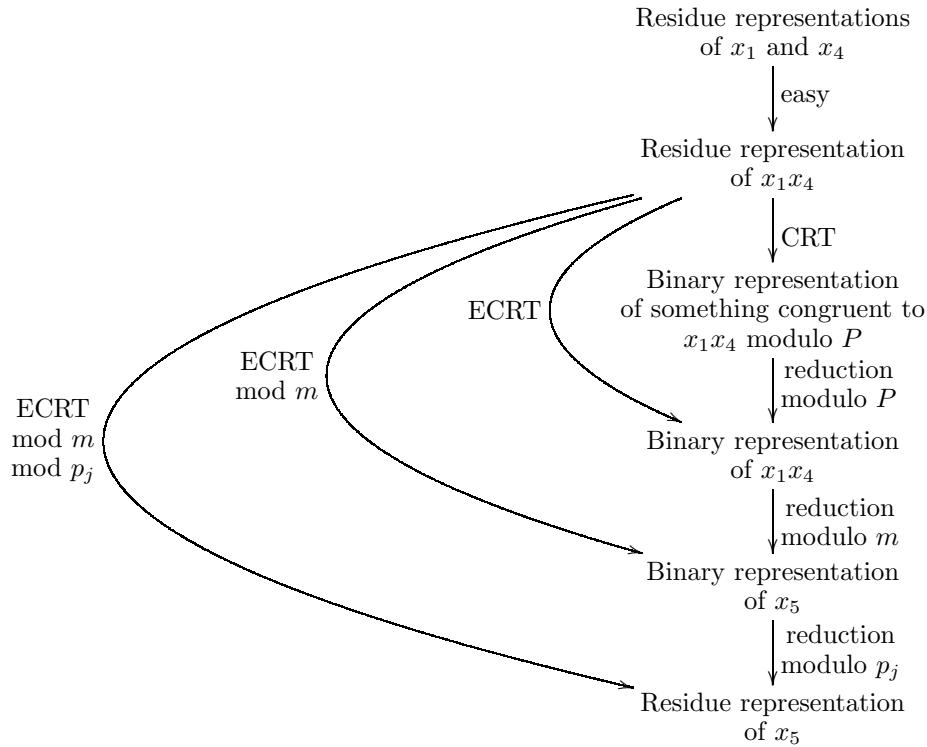
The usual solution for computing $x^e \bmod m$ is to compute small integers that are congruent modulo m to various powers of x . See, e.g., [18, Section 4.6.3], [11, Section 1.2], [16], [20], and [9]. For example, say $e = 10$. One can compute $x_1 = x \bmod m$, then $x_2 = x_1^2 \bmod m$, then $x_4 = x_2^2 \bmod m$, then $x_5 = x_1 x_4 \bmod m$, and finally $x_{10} = x_5^2 \bmod m$. It is often convenient to allow x_1, x_2, x_4, x_5 to be slightly larger than m .

The output $x^e \bmod m$ and inputs x, e, m are conventionally written in binary. Standard practice is to also use the binary representation for x_1, x_1^2, x_2, x_2^2 , etc. We instead use the **residue representation**: x_1 is represented as $(x_1 \bmod p_1, x_1 \bmod p_2, \dots, x_1 \bmod p_s)$, x_1^2 is represented as $(x_1^2 \bmod p_1, x_1^2 \bmod p_2, \dots, x_1^2 \bmod p_s)$, etc. Here p_1, p_2, \dots, p_s are small pairwise coprime positive integers, typically primes, such that $P = p_1 p_2 \cdots p_s$ is sufficiently large. Note that, because p_j is small, it is easy to compute (e.g.) $x_1 x_4 \bmod p_j$ from $x_1 \bmod p_j$ and $x_4 \bmod p_j$.

Received by the editor August 18, 2003 and, in revised form, June 15, 2005.

2000 *Mathematics Subject Classification*. Primary 11Y16; Secondary 68W10.

This paper combines and improves the preliminary papers [7] by Bernstein and [29] by Sorenson. Bernstein was supported by the National Science Foundation under grants DMS-9600083 and DMS-9970409. Sorenson was supported by the National Science Foundation under grant CCR-9626877. Sorenson completed part of this work while on sabbatical at Purdue University in Fall 1998.



The usual Chinese remainder theorem says that (for example) $x_1x_4 \bmod P$ is determined by the residue representation of x_1x_4 . In fact, any integer u is congruent modulo P to a particular linear combination of $u \bmod p_1, u \bmod p_2, \dots, u \bmod p_s$. The *explicit* Chinese remainder theorem says, in a computationally useful form, exactly what multiple of P must be subtracted from the linear combination to obtain u . See §2.

Once we know the binary representation of x_1x_4 , we could reduce it modulo m to obtain the binary representation of x_5 . We actually use the explicit CRT modulo m to obtain the binary representation of x_5 directly from the residue representation of x_1x_4 . See §3.

Once we know the binary representation of x_5 , we reduce it modulo each p_j to obtain the residue representation of x_5 . An alternative is to use the explicit CRT modulo m modulo p_j to obtain the residue representation of x_5 directly from the residue representation of x_1x_4 . See §4.

If P is sufficiently large then, as discussed in §5, one can perform several multiplications in the residue representation before reduction modulo m . This is particularly beneficial for parallel computation: operations in the residue representation are more easily parallelized than reductions modulo m . By optimizing parameters, we obtain the following result:

Theorem 1.1. *There is an algorithm for the Common CRCW PRAM that, given the binary representations of positive integers x , e , and m , of total bit length n ,*

computes the binary representation of $x^e \bmod m$ in time $O(n/\lg \lg n)$ using $n^{O(1)}$ processors.

In §6 we define the Common CRCW PRAM, in §8 we present the algorithm, and in §7 we present a simpler algorithm taking time $O(n)$.

For comparison: A naive parallelization of the standard binary algorithm takes superlinear time; see §7. Adleman and Kompella gave an $O((\lg n)^3)$ expected time algorithm using $\exp(O(\sqrt{n \lg n}))$ processors; see [2]. Von zur Gathen gave an NC algorithm for the highly special case that m is polynomially smooth; see [35]. It is unknown whether general integer modular exponentiation is in NC. For more on parallel exponentiation algorithms, see Gordon’s survey [16].

Our techniques can easily be applied to exponentiation in finite rings more general than \mathbf{Z}/m . For example, the deterministic polynomial-time primality test of Agrawal, Kayal, and Saxena in [4] can be carried out in sublinear time using a polynomial number of processors on the Common CRCW PRAM; the bottleneck is exponentiation in a ring of the form $(\mathbf{Z}/m)[x]/(x^k - 1)$.

2. THE EXPLICIT CHINESE REMAINDER THEOREM

Algorithms for integer arithmetic using the residue representation were introduced in the 1950s by Svoboda, Valach, and Garner, according to [18, Section 4.3.2]. This section discusses conversion from the residue representation to the binary representation.

For each real number α such that $\alpha - 1/2 \notin \mathbf{Z}$, define $\text{round } \alpha$ as the unique integer r with $|r - \alpha| < 1/2$.

Theorem 2.1. *Let p_1, p_2, \dots, p_s be pairwise coprime positive integers. Write $P = p_1 p_2 \cdots p_s$. Let q_1, q_2, \dots, q_s be integers with $q_i P/p_i \equiv 1 \pmod{p_i}$. Let u be an integer with $|u| < P/2$. Let u_1, u_2, \dots, u_s be integers with $u_i \equiv u \pmod{p_i}$. Let t_1, t_2, \dots, t_s be integers with $t_i \equiv u_i q_i \pmod{p_i}$. Then $u = P\alpha - P \text{round } \alpha$ where $\alpha = \sum_i t_i/p_i$.*

Proof. $P\alpha = \sum_i t_i(P/p_i) \equiv \sum_i u_i q_i(P/p_i) \equiv u_j q_j(P/p_j) \equiv u \pmod{p_j}$ for each j , so $P\alpha \equiv u \pmod{P}$. Write $r = \alpha - u/P$. Then r is an integer, and $|r - \alpha| = |u/P| < 1/2$, so $r = \text{round } \alpha$, i.e., $u/P = \alpha - \text{round } \alpha$. \square

The usual Chinese remainder theorem says that $u \equiv \sum_i t_i P/p_i \pmod{P}$; this is the integer version of Lagrange’s interpolation formula. This is a popular way to compute u from the residues u_i : first compute $t_i = u_i q_i \bmod p_i$ or simply $t_i = u_i q_i$, then compute $P\alpha = \sum_i t_i P/p_i$, then reduce $P\alpha$ modulo P to the right range.

The explicit Chinese remainder theorem, Theorem 2.1, suggests another way to divide $P\alpha$ by P . Use t_i and p_i directly to compute a low-precision approximation to $\alpha = \sum_i t_i/p_i$ with sufficient accuracy to determine $\text{round } \alpha$; see, for example, Theorem 2.2 below. Then subtract $P \text{round } \alpha$ from $P\alpha$ to obtain u .

As far as we know, the first use of the explicit Chinese remainder theorem was by Montgomery and Silverman in [23, Section 4]. It has also appeared in [22], [13, Section 2.1], [7], [29], and [3, Section 5].

Theorem 2.2. *Let $\beta_1, \beta_2, \dots, \beta_s$ be real numbers. Let r and a be integers. If $|r - \sum_i \beta_i| < 1/4$ and $2^a \geq 2s$ then $r = \lfloor 3/4 + 2^{-a} \sum_i \lfloor 2^a \beta_i \rfloor \rfloor$.*

Proof. $r < 1/4 + \sum_i \beta_i = 1/4 + 2^{-a} \sum_i 2^a \beta_i \leq 1/4 + 2^{-a}(s + \sum_i \lfloor 2^a \beta_i \rfloor) \leq 1/4 + 2^{-a}(2^a/2 + \sum_i \lfloor 2^a \beta_i \rfloor) = 3/4 + 2^{-a} \sum_i \lfloor 2^a \beta_i \rfloor \leq 3/4 + \sum_i \beta_i < r + 1$. \square

In the situation of Theorem 2.1, assume further that $|u| < P/4$. We can then use Theorem 2.2 to quickly compute $r = \text{round } \alpha$. We choose a with $2^a \geq 2s$, then compute the fixed-point approximations $2^{-a} \lfloor 2^a t_i / p_i \rfloor$ to t_i / p_i , then compute $r = \lfloor 3/4 + 2^{-a} \sum_i \lfloor 2^a t_i / p_i \rfloor \rfloor$.

3. THE EXPLICIT CRT MOD m

Theorem 3.1. *Let p_1, p_2, \dots, p_s be pairwise coprime positive integers. Write $P = p_1 p_2 \cdots p_s$. Let q_1, q_2, \dots, q_s be integers with $q_i P / p_i \equiv 1 \pmod{p_i}$. Let u be an integer with $|u| < P/2$. Write $t_i = u q_i \pmod{p_i}$. Let m be a positive integer. Write $v = \sum_i t_i (P/p_i \pmod{m}) - (P \pmod{m}) \text{round} \sum_i t_i / p_i$. Then $u \equiv v \pmod{m}$, and $|v| \leq m \sum_i p_i$.*

The hypotheses on p, P, q, u are the same as in Theorem 2.1. The hypothesis on t is more restrictive: Theorem 2.1 allowed any integer $t_i \equiv u q_i \pmod{p_i}$, whereas Theorem 3.1 insists that $0 \leq t_i < p_i$.

One can allow a slightly wider range of t_i , often making t_i easier to compute, at the expense of a larger bound on $|v|$. In the other direction, one can reduce the bound on $|v|$ by taking t_i between $-p_i/2$ and $p_i/2$, and one can reduce the bound further by similarly adjusting $P/p_i \pmod{m}, P \pmod{m}$.

Proof. By Theorem 2.1, $u = \sum_i t_i P / p_i - P \text{round} \sum_i t_i / p_i \equiv v$. Furthermore, $v \leq \sum_i t_i (P/p_i \pmod{m}) \leq \sum_i p_i m$ since $t_i \leq p_i$, and $-v \leq (P \pmod{m}) \text{round} \sum_i t_i / p_i \leq m \sum_i p_i$ since $t_i \leq p_i^2$. \square

Theorem 3.1 suggests the following representation of integers modulo m . Select moduli p_1, \dots, p_s whose product $P = p_1 \cdots p_s$ exceeds $4(m \sum_i p_i)^2$. Use the vector $(x \pmod{p_1}, \dots, x \pmod{p_s})$, where x is any integer between $-m \sum_i p_i$ and $m \sum_i p_i$, to represent $x \pmod{m}$. Note that each element of \mathbf{Z}/m has many representations.

The following procedure, given two such vectors $(x \pmod{p_1}, \dots, x \pmod{p_s})$ and $(y \pmod{p_1}, \dots, y \pmod{p_s})$, computes another such vector $(v \pmod{p_1}, \dots, v \pmod{p_s})$ with $v \equiv xy \pmod{m}$:

- (1) Precompute $q_i, P/p_i \pmod{m}$, and $P \pmod{m}$.
- (2) Compute $t_i = ((x \pmod{p_i})(y \pmod{p_i})q_i) \pmod{p_i}$, so that $t_i = (xy)q_i \pmod{p_i}$.
- (3) Compute $\text{round} \sum_i t_i / p_i$ by Theorem 2.2.
- (4) Compute $v = \sum_i t_i (P/p_i \pmod{m}) - (P \pmod{m}) \text{round} \sum_i t_i / p_i$. Then v is between $-m \sum_i p_i$ and $m \sum_i p_i$, and $v \equiv xy \pmod{m}$, by Theorem 3.1.
- (5) Compute the residues $(v \pmod{p_1}, \dots, v \pmod{p_s})$.

The output can then be used in subsequent multiplications. One can carry out more componentwise operations before applying Theorem 3.1 if P is chosen larger; see §5 for further discussion.

Montgomery and Silverman in [23, Section 4] suggested computing $u \pmod{m}$ by first computing v . (Another way to compute $u \pmod{m}$, well suited for FFT-based arithmetic, is to perform the computation of [8, Section 13] modulo m .) The idea of using v for subsequent operations, and not bothering to compute $u \pmod{m}$, was introduced in [7].

4. THE EXPLICIT CRT MOD $m \bmod p_j$

Theorem 4.1. *Under the assumptions of Theorem 3.1,*

$$v \equiv \sum_i t_i(P/p_i \bmod m \bmod p_j) - (P \bmod m \bmod p_j) \text{ round } \sum_i t_i/p_i \pmod{p_j}.$$

Proof. Reduce the definition of v modulo p_j . □

In §3 we discussed precomputing $P/p_i \bmod m$, precomputing $P \bmod m$, computing v , and reducing v modulo p_j . Theorem 4.1 suggests a different approach: precompute $P/p_i \bmod m \bmod p_j$; precompute $P \bmod m \bmod p_j$; compute $v \bmod p_j$ as $(\sum_i t_i(P/p_i \bmod m \bmod p_j) - (P \bmod m \bmod p_j) \text{ round } \sum_i t_i/p_i) \bmod p_j$. This idea was introduced in [7].

Theorem 4.1 is particularly well suited to parallel computation, as briefly observed in [7]. By using additional ideas described in the remaining sections of this paper, one can exponentiate on a polynomial-size parallel computer in sublinear time, as shown in [29]. This paper includes all the results from [7] and [29].

5. HIGHER POWERS

We begin by reviewing the left-to-right base-2 algorithm for computing $x^e \bmod m$:

```

Let  $l$  denote the number of bits in  $e$ ;
Write  $e = \sum_{k=0}^{l-1} e_k 2^k$ , where  $e_k \in \{0, 1\}$ ;
 $y \leftarrow 1$ ;
For( $k \leftarrow l - 1$ ;  $k \geq 0$ ;  $k \leftarrow k - 1$ ) do:
     $y \leftarrow y^2 x^{e_k} \bmod m$ ;
Output( $y$ );
    
```

More generally, let b be a power of 2. The base- b algorithm is as follows:

```

Let  $l$  denote the number of base- $b$  digits in  $e$ ;
Write  $e = \sum_{k=0}^{l-1} e_k b^k$ , where  $0 \leq e_k < b$ ;
 $y \leftarrow 1$ ;
For( $k \leftarrow l - 1$ ;  $k \geq 0$ ;  $k \leftarrow k - 1$ ) do:
     $y \leftarrow y^b x^{e_k} \bmod m$ ;
Output( $y$ );
    
```

Sane people take every opportunity to reduce intermediate results modulo m : for example, they compute $y^b \bmod m$ by squaring y , reducing modulo m , squaring again, reducing again, etc. We call this quite reasonable practice into question. It may be better to first compute $y^b x^{e_k}$, then reduce the result modulo m . The benefit of performing fewer reductions may outweigh the cost of performing arithmetic on larger numbers. (Similar comments apply to other exponentiation algorithms.)

When we use the explicit Chinese remainder theorem to perform multiplications modulo m , the cost of handling larger numbers is in using more primes p_1, p_2, \dots . In the context of parallel computation, increasing the number of primes mainly affects the number of processors, not the run time. In fact, as explained in the next three sections, we can save an asymptotically non-constant factor in the run time in various models of parallel computation, at a reasonable expense in the number of processors.

Perhaps the same idea can also save a small constant factor in the serial case.

6. DEFINITION OF THE CREW PRAM AND THE COMMON CRCW PRAM

We use two models of parallel computation in this paper. This section defines the models.

In both models, computers are “parallel random-access machines” (PRAMs) in which many processors can access variable locations in a shared memory. The models differ in how they handle memory conflicts: the “Common CRCW PRAM” allows many processors to simultaneously write a common value to a single memory location, while the “CREW PRAM” does not. Our exponentiation algorithm in §8 relies on the extra power of the Common CRCW PRAM; the simplified algorithm in §7 works in either model.

We caution the reader that phrases such as “CREW PRAM” are not sufficient to pin down a model of computation in the literature. For example, [25, Sections 2.6 and 15.2] defines unrealistic PRAMs in which the processor word size is permitted to grow linearly as a computation progresses; at the other extreme, [14] presents computations that fit into logarithmic-size index registers and constant-size data registers. The complexity of a computation can depend heavily on details of the model. See [15] for a comparison of several models.

Computers. A parallel computer is parametrized by positive integers p, s, w, r, i such that $p \leq 2^w$ and $s \leq 2^w$. The computer has p processors, numbered from 0 through $p - 1$, operating on s words of shared memory, labelled $\text{mem}[0]$ through $\text{mem}[s-1]$. Each word is a w -bit string, often interpreted as the binary representation of an integer between 0 and $2^w - 1$. Each processor has r registers, labelled $\text{reg}[0]$ through $\text{reg}[r - 1]$; each register contains one word. The computer also has space for i instructions; each processor has an instruction pointer.

Algorithms. A parallel algorithm is a sequence of instructions. (The algorithm fails on a computer of size p, s, w, r, i if it has more than i instructions.) Here are the possible instructions:

- (1) Clear j : Set $\text{reg}[j] \leftarrow 0$. (The algorithm fails if $j \geq r$.)
- (2) Increment j, k : Set $\text{reg}[j] \leftarrow (\text{reg}[k] + 1) \bmod 2^w$.
- (3) Add j, k, ℓ : Set $\text{reg}[j] \leftarrow (\text{reg}[k] + \text{reg}[\ell]) \bmod 2^w$.
- (4) Subtract j, k, ℓ : Set $\text{reg}[j] \leftarrow (\text{reg}[k] - \text{reg}[\ell]) \bmod 2^w$.
- (5) Shift left j, k, ℓ : Set $\text{reg}[j] \leftarrow (\text{reg}[k] 2^{\text{reg}[\ell]}) \bmod 2^w$.
- (6) Shift right j, k, ℓ : Set $\text{reg}[j] \leftarrow \lfloor \text{reg}[k] / 2^{\text{reg}[\ell]} \rfloor$.
- (7) Word size j : Set $\text{reg}[j] \leftarrow w$.
- (8) Identify j : Set $\text{reg}[j] \leftarrow$ this processor’s number.
- (9) Read j, k : Set $\text{reg}[j] \leftarrow \text{mem}[\text{reg}[k]]$. (The algorithm fails if $\text{reg}[k] \geq s$.)
- (10) Write j, k : Set $\text{mem}[\text{reg}[k]] \leftarrow \text{reg}[j]$.
- (11) Jump j, k, ℓ : If $\text{reg}[k] \geq \text{reg}[\ell]$, go to instruction j , rather than proceeding to the next instruction as usual.

An n -bit input to the algorithm is placed into memory at time 0, with n in the first word of memory, and the n bits packed into the next $\lceil n/w \rceil$ words of memory. (The algorithm fails if $2^w \leq n$, or if $s < 1 + \lceil n/w \rceil$.) All other words of memory, registers, etc. are set to 0. Each processor performs one instruction at time 1; each processor performs one instruction at time 2; and so on until all the processors have halted. The output is then in memory, encoded in the same way as the input.

Memory conflicts. During one time step, a single memory location might be accessed by more than one processor.

The **concurrent-read exclusive-write PRAM**, or **CREW PRAM**, allows any number of processors to simultaneously read the same memory location, but it does not allow two processors to simultaneously write to a single memory location.

The **concurrent-read common-concurrent-write PRAM**, or **Common CRCW PRAM**, allows any number of processors to simultaneously read the same memory location, and allows any number of processors to simultaneously write the same memory location, if they all write the same value. (If two processors attempt to write different values, the algorithm fails.)

A memory location cannot be read and written simultaneously.

Asymptotics. When we say that a parallel algorithm handles an n -bit input using (for example) time $O(n)$ on $O(n^3)$ processors with a word size of $O(\lg n)$ using $O(n^4)$ memory locations, we mean that it works on any sufficiently large computer: there are functions $p(n) \in O(n^3)$, $s(n) \in O(n^4)$, $w(n) \in O(\lg n)$, and $t(n) \in O(n)$ such that, for every n , every $p \geq p(n)$, every $s \geq s(n)$, every $w \geq w(n)$, every r larger than the highest register number used in the algorithm, every i larger than the length of the algorithm, and every input string of length n , the algorithm runs without failure on that input in time at most $t(n)$ on a parallel computer of size p, s, w, r, i .

When we do not mention the word size (and number of memory locations), we always mean that the required word size is logarithmic in the time-processor product (and, consequently, the required number of memory locations is polynomial in the time-processor product).

7. MODULAR EXPONENTIATION ON THE CREW PRAM

In this section we present three successively better parallel algorithms, using a polynomial number of processors, for modular exponentiation on the CREW PRAM.

Time $O(n(\lg n)^2)$. Given the binary representations of n -bit integers x and y , one can compute xy in time $O(\lg n)$ using $n^{O(1)}$ processors. One can compute x/y and $x \bmod y$ by Newton's method in time $O((\lg n)^2)$ using $n^{O(1)}$ processors. See [32, Theorem 12.1].

The base-2 exponentiation algorithm shown in §5, using these subroutines, computes $x^e \bmod m$ in time $O(n(\lg n)^2)$ if x, e, m have n bits. There are $O(n)$ iterations in the algorithm, each iteration involving a constant number of multiplications and divisions of $O(n)$ -bit integers. The number of processors is polynomial in n .

Time $O(n \lg n)$. A division algorithm of Beame, Cook, and Hoover takes time only $O(\lg n)$ using $n^{O(1)}$ processors, after a precomputation taking time $O((\lg n)^2)$. See [5].

The main subroutine in the Beame-Cook-Hoover algorithm computes powers x^v , where x has n bits and $v \in \{0, 1, \dots, n\}$, in time $O(\lg n)$. The idea is to use the CRT to recover x^v from the remainders $x^v \bmod p$ for enough small primes p . Beame, Cook, and Hoover precompute the primes p and the powers $u^v \bmod p$ for all $u \in \{0, 1, \dots, p-1\}$ and all $v \in \{0, 1, \dots, n\}$; then they can compute $x^v \bmod p$ as $(x \bmod p)^v \bmod p$. This is one of the ideas that we use in §8.

Time $O(n)$. To save another $\lg n$ factor, we choose an integer $b \geq 2$ so that $\lg b \in \Theta(\lg n)$, and we use the base- b exponentiation algorithm.

The number of iterations drops to $O(\log_b e) = O(n/\lg n)$. Each iteration involves computing $y^b x^v \bmod m$ for some $v \in \{0, 1, \dots, b-1\}$; we compute y^b and x^v in time $O(\lg n)$ by the Beame-Cook-Hoover subroutine, then multiply, then use the Beame-Cook-Hoover algorithm again to divide by m .

8. MODULAR EXPONENTIATION ON THE COMMON CRCW PRAM

In this section, we present our sublinear-time algorithm, which uses a polynomial number of processors, for modular exponentiation on the Common CRCW PRAM. We begin by informally explaining how we save an additional factor of $\lg \lg n$ from the running time of our linear-time algorithm from §7. Next we review some FFT-based subroutines that we use in the algorithm. We then present our algorithm along with a proof of its complexity, followed by some discussion.

How we save a $\lg \lg n$ factor. The algorithm in this section, like the simplified linear-time algorithm of §7, performs several multiplications before each reduction modulo m . It achieves better parallelization than the algorithm of §7 as follows:

- (1) It avoids the binary representation in the main loop. It uses the explicit Chinese remainder theorem modulo m modulo p_j , as described in §4, to work consistently in the residue representation.
- (2) It uses the Cole-Vishkin parallel-addition algorithm to add $O(n)$ integers, each having $O(\lg n)$ bits, in time $O((\lg n)/\lg \lg n)$ using $n^{1+o(1)}$ processors. See [12] and [34]. This is where we are making use of the increased power of the CRCW PRAM over the CREW PRAM. Computing such sums is the bottleneck in using the explicit Chinese remainder theorem modulo m modulo p_j .

Perhaps the same $\lg \lg n$ speedup can be obtained with other redundant representations of integers modulo m ; we leave this exploration to the reader.

FFT-based subroutines. We take advantage of several standard FFT-based tools here:

- (1) Given n bits representing two integers x, y (in binary), one can use the Schönhage-Strassen multiplication algorithm to compute the product xy in time $(\lg n)^{O(1)}$ using $n^{1+o(1)}$ processors. See [28].
- (2) Given n bits representing two integers x, y with $y \neq 0$, one can compute the quotient $\lfloor x/y \rfloor$ in time $(\lg n)^{O(1)}$ using $n^{1+o(1)}$ processors. See [32, Theorem 12.1].
- (3) Given n bits representing s integers, one can compute the product of all the integers in time $(\lg n)^{O(1)}$ using $n^{1+o(1)}$ processors, by multiplying pairs in parallel.
- (4) Given n bits representing integers u, p_1, p_2, \dots, p_s , one can use the Borodin-Moenck remainder-tree algorithm to compute $u \bmod p_1, \dots, u \bmod p_s$ in time $(\lg n)^{O(1)}$ using $n^{1+o(1)}$ processors. See [21], [10, Sections 4–6], and [8, Section 18].

Without the FFT-based tools, Step 4 below requires roughly n^3 processors to be carried out in polylogarithmic time. This is (after some serialization) adequate for Theorem 8.1, but it becomes a bottleneck when e is much shorter than m .

The algorithm. Fix $\epsilon > 0$. Some of the O constants below depend on ϵ ; we indicate when this happens by adding ϵ as a subscript.

Theorem 8.1. *There is an algorithm for the Common CRCW PRAM that, given the binary representations of positive integers x , e , and m , of total bit length n , computes the binary representation of $x^e \bmod m$ in time $O_\epsilon(n/\lg \lg n)$ using $O(n^{2+\epsilon})$ processors.*

Outside the inner loop (Step 6 below), the algorithm takes time only $(\lg n)^{O(1)}$ using $O(n^{2+\epsilon})$ processors.

Step 1: Build multiplication tables. Fix a positive rational number $\delta \leq \epsilon/3$. Compute an integer $a \geq 1$ within 1 of $\delta \lg n$. Note that an a -bit integer fits into $O_\epsilon(1)$ words, since $a \in O_\epsilon(\lg n)$.

Compute xy and (for $y \neq 0$) $\lfloor x/y \rfloor$ for each pair x, y of a -bit integers in parallel. Store the results in a table of $O(2^{2a})$ words. This takes time $O(a)$ using 2^{2a} processors; i.e., time $O_\epsilon(\lg n)$ using $O(n^{2\delta})$ processors.

Note that, if x and y are $O(\lg n)$ -bit integers, then x and y are also $O_\epsilon(a)$ -bit integers, so one processor can compute $x + y$, $x - y$, xy , and $\lfloor x/y \rfloor$ in time $O_\epsilon(1)$ with the help of this table.

Step 2: Find primes. Define $b = 2^a$. Note that $b \in \Theta(n^\delta)$.

Find the smallest integer $k \geq 6$ such that $2^k \geq 4 + 4b(n + k)$. Use a parallel sieve, as described in [30], to find the primes p_1, p_2, \dots, p_s between 2 and 2^k . This takes time $O((1 + \delta) \lg n)$ using $O(n^{1+\delta})$ processors; note that $s \leq p_s \in O(n^{1+\delta})$.

Define $P = p_1 p_2 \cdots p_s$. Note that $\lg P \in O(n^{1+\delta})$ by, e.g., [27, Theorem 9]. Note also that $P \geq 4(m^2 \sum p_i)^b$. (Indeed, $2^k \geq 41$, so $\log P = \log p_1 + \log p_2 + \cdots + \log p_s \geq 2^k(1 - 1/\log 2^k)$ by [27, Theorem 4]; also $p_1 + p_2 + \cdots + p_s \leq 2^{2k}$. Thus $\lg P \geq \log P \geq 2^{k-1} \geq 2 + b(2n + 2k) \geq \lg 4 + b(\lg m^2 + \lg \sum p_i)$.)

Multiply p_1, \dots, p_s to compute P . This takes time $(\lg \lg P)^{O(1)}$ using $(\lg P)^{1+o(1)}$ processors; i.e., time $(\lg n)^{O(1)}$ using $n^{1+\delta+o(1)}$ processors.

Step 3: Build power tables. For each $i \in \{1, 2, \dots, s\}$ in parallel, for each $u \in \{0, 1, \dots, p_i - 1\}$ in parallel, for each $v \in \{0, 1, \dots, b\}$ in parallel (or serial), compute $u^v \bmod p_i$. Store the results in a table. This takes time $O_\epsilon(\lg b)$ using $O(sp_s(b + 1))$ processors; i.e., time $O_\epsilon(\lg n)$ using $O(n^{2+2\delta} \lg n)$ processors.

(An alternative approach, with smaller tables, is to find a generator g_i for the multiplicative group $(\mathbf{Z}/p_i)^*$, build a table of powers of g_i , and build a table of discrete logarithms base g_i . See [30].)

Step 4: Compute ECRT coefficients. For each i in parallel, compute $P_i = P/p_i$; $P_i \bmod p_i$; $P_i \bmod m$; and $P_i \bmod m \bmod p_1, P_i \bmod m \bmod p_2, \dots, P_i \bmod m \bmod p_s$. This takes time $(\lg \lg P)^{O(1)}$ using $s(\lg P)^{1+o(1)}$ processors; i.e., time $(\lg n)^{O(1)}$ using $n^{2+2\delta+o(1)}$ processors.

Next, for each i in parallel, compute $q_i = (P_i \bmod p_i)^{p_i-2} \bmod p_i$, so that q_i is the inverse of P_i modulo p_i . This takes time $O(\lg p_s)$ using s processors; i.e., time $O_\epsilon(\lg n)$ using $O(n^{1+\delta})$ processors.

Step 5: Convert to the residue representation. Set $x \leftarrow x \bmod m$. This takes time $(\lg n)^{O(1)}$ using $n^{1+o(1)}$ processors.

Set $x_1 \leftarrow x \bmod p_1, x_2 \leftarrow x \bmod p_2, \dots, x_s \leftarrow x \bmod p_s$. This takes time $(\lg n)^{O(1)}$ using $n^{1+\delta+o(1)}$ processors. Also set $y_1 \leftarrow 1, y_2 \leftarrow 1, \dots, y_s \leftarrow 1$.

Let l be the number of base- b digits in e . Write $e = \sum_{k=0}^{l-1} e_k b^k$ with $0 \leq e_k < b$. Note that $l \in O_\epsilon(n/\lg n)$.

Step 6: Inner loop. Repeat the following substeps for $k \leftarrow l-1, k \leftarrow l-2, \dots, k \leftarrow 0$. Each substep will take time $O_\epsilon((\lg n)/\lg \lg n)$, so the total time here is $O_\epsilon(n/\lg \lg n)$.

Invariant: The integer y represented by (y_1, y_2, \dots, y_s) is between $-m \sum_i p_i$ and $m \sum_i p_i$, and is congruent to $x^{e_{l-1}b^{l-k-2} + \dots + e_{k+2}b^{e_{k+1}}}$ modulo m . The plan is to compute the residue representation of $u = y^b x^{e_k}$, and then change y to the integer v identified in Theorem 3.1. Note that $|u| < (m^2 \sum_i p_i)^b \leq P/4$; hence $v \equiv u \equiv x^{e_{l-1}b^{l-k-1} + \dots + e_{k+2}b^{e_{k+1}} + e_k} \pmod{m}$.

For each i in parallel, compute $t_i \leftarrow (y_i^b \bmod p_i)(x_i^{e_k} \bmod p_i)q_i \bmod p_i$. This takes time $O(1)$ using s processors, thanks to the precomputed power table.

Compute $r \leftarrow \text{round} \sum_i t_i/p_i$ with the help of Theorem 2.2. The fixed-point divisions take time $O(1)$ using s processors. The addition takes time $O((\lg s)/\lg \lg s)$ using $s^{1+o(1)}$ processors using the Cole-Vishkin addition algorithm.

For each j in parallel, compute $y_j \leftarrow \sum_i t_i(P_i \bmod m \bmod p_j)$. This takes time $O((\lg s)/\lg \lg s)$ using $s^{2+o(1)}$ processors, again using the Cole-Vishkin algorithm; i.e., time $O((\lg n)/\lg \lg n)$ using $n^{2+2\delta+o(1)}$ processors.

For each j in parallel, compute $y_j \leftarrow (y_j - r(P \bmod m \bmod p_j)) \bmod p_j$. This takes time $O(1)$ using s processors.

Step 7: Convert to the binary representation. Use Theorem 2.1 to compute the binary representation of the integer y whose residue representation is (y_1, y_2, \dots, y_s) . This takes time $(\lg n)^{O(1)}$ using $n^{2+2\delta+o(1)}$ processors.

Now $x^e = x^{e_{l-1}b^{l-1} + \dots + e_1b + e_0} \equiv y \pmod{m}$. The output of the algorithm is $y \bmod m$.

Discussion. The algorithm can be converted into a polynomial-size unbounded-fan-in circuit of depth $O_\epsilon(n/\lg \lg n)$, using the techniques explained in [31]. The last algorithm discussed in §7 can be converted into a polynomial-size bounded-fan-in circuit of depth $O(n)$ by the same techniques.

The only previous sublinear-time algorithms were algorithms that use many more processors or that drastically restrict m . See [16, Section 6] for a survey.

We do not know an exponentiation algorithm that takes time $O(n/\lg \lg n)$ using $n^{2+o(1)}$ processors. Allowing ϵ to vary with n would hurt the time bound: for example, choosing ϵ as $1/\lg \lg \lg n$ would produce an algorithm that takes time $O(n \lg \lg \lg n / \lg \lg n)$ using $O(n^{2+1/\lg \lg \lg n})$ processors.

In a few applications, the exponent e is only a small part of the input, perhaps \sqrt{n} or $\sqrt[3]{n}$ bits. The running time depends primarily on the length of e : the algorithm takes time $O_\epsilon((\lg e)/\lg \lg n) + (\lg n)^{O(1)}$ using $O(n^{2+\epsilon})$ processors.

REFERENCES

- [1] Leonard M. Adleman, Kireeti Kompella, *Proceedings of the 20th ACM symposium on theory of computing*, Association for Computing Machinery, New York, 1988.
- [2] ———, *Using smoothness to achieve parallelism*, in [1] (1988), 528–538.
- [3] Amod Agashe, Kristin Lauter, Ramarathnam Venkatesan, *Constructing elliptic curves with a known number of points over a prime field*, in [33] (2004), 1–17. URL: <http://research.microsoft.com/~klauter/>. MR2005m:11112

- [4] Manindra Agrawal, Neeraj Kayal, Nitin Saxena, *PRIMES is in P*, *Annals of Mathematics* **160** (2004), 781–793. URL: <http://ProjectEuclid.org/Dienst/UI/1.0/Summarize/euclid.anm/1111770735>. MR2123939 (2006a:11170)
- [5] Paul W. Beame, Stephen A. Cook, H. James Hoover, *Log depth circuits for division and related problems*, *SIAM Journal on Computing* **15** (1986), 994–1003. ISSN 0097-5397. MR0861365 (88b:68059)
- [6] Daniel J. Bernstein, *Detecting perfect powers in essentially linear time, and other studies in computational number theory*, Ph.D. thesis, University of California at Berkeley, 1995.
- [7] Daniel J. Bernstein, *Multidigit modular multiplication with the explicit Chinese remainder theorem*, in [6] (1995). URL: <http://cr.yp.to/papers.html#mmecrt>.
- [8] Daniel J. Bernstein, *Fast multiplication and its applications*, to appear. URL: <http://cr.yp.to/papers.html#multapps>. ID 8758803e61822d485d54251b27b1a20d.
- [9] Daniel J. Bernstein, *Pippenger's exponentiation algorithm*, to be incorporated into author's *High-speed cryptography* book. URL: <http://cr.yp.to/papers.html#pippenger>.
- [10] Allan Borodin, Robert T. Moenck, *Fast modular transforms*, *Journal of Computer and System Sciences* **8** (1974), 366–386; older version, not a subset, in [21]. ISSN 0022-0000. URL: <http://cr.yp.to/bib/entries.html#1974/borodin>. MR0371144 (51:7365)
- [11] Henri Cohen, *A course in computational algebraic number theory*, *Graduate Texts in Mathematics*, 138, Springer-Verlag, Berlin, 1993. ISBN 3-5440-55640-0. MR1228206 (94i:11105)
- [12] Richard Cole, Uzi Vishkin, *Faster optimal parallel prefix sums and list ranking*, *Information and Computation* **81** (1989), 334–352. ISSN 0890-5401. MR1000070 (90k:68056)
- [13] Jean-Marc Couveignes, *Computing a square root for the number field sieve*, in [19] (1993), 95–102. MR1321222
- [14] Shawna Meyer Eikenberry, Jonathan P. Sorenson, *Efficient algorithms for computing the Jacobi symbol*, *Journal of Symbolic Computation* **26** (1998), 509–523. ISSN 0747-7171. MR1646683 (99h:11146)
- [15] Faith E. Fich, *The complexity of computation on the parallel random access machine*, in [26] (1993), 843–899. MR1212591 (94c:68086)
- [16] Daniel M. Gordon, *A survey of fast exponentiation methods*, *Journal of Algorithms* **27** (1998), 129–146. ISSN 0196-6774. URL: <http://www.ccrwest.org/gordon/dan.html>. MR1613189 (99g:94014)
- [17] Richard M. Karp (chairman), *13th annual symposium on switching and automata theory*, IEEE Computer Society, Northridge, 1972.
- [18] Donald E. Knuth, *The art of computer programming, volume 2: seminumerical algorithms*, 3rd edition, Addison-Wesley, Reading, 1997. ISBN 0-201-89684-2.
- [19] Arjen K. Lenstra, Hendrik W. Lenstra, Jr. (editors), *The development of the number field sieve*, *Lecture Notes in Mathematics*, 1554, Springer-Verlag, Berlin, 1993. ISBN 3-540-57013-6. MR1321216 (96m:11116)
- [20] Alfred J. Menezes, Paul C. van Oorschot, Scott A. Vanstone, *Handbook of applied cryptography*, CRC Press, Boca Raton, Florida, 1996. ISBN 0-8493-8523-7. URL: <http://cacr.math.uwaterloo.ca/hac>. MR1412797 (99g:94015)
- [21] Robert T. Moenck, Allan Borodin, *Fast modular transforms via division*, in [17] (1972), 90–96; newer version, not a superset, in [10]. URL: <http://cr.yp.to/bib/entries.html#1972/moenck>. MR0371144 (51:7365)
- [22] Peter L. Montgomery, *An FFT extension of the elliptic curve method of factorization*, Ph.D. thesis, University of California at Los Angeles, 1992. URL: <http://cr.yp.to/bib/entries.html#1992/montgomery>.
- [23] Peter L. Montgomery, Robert D. Silverman, *An FFT extension to the $P - 1$ factoring algorithm*, *Mathematics of Computation* **54** (1990), 839–854. ISSN 0025-5718. MR1011444 (90j:11142)
- [24] Andrew M. Odlyzko, Gary Walsh, Hugh Williams (editors), *Conference on the mathematics of public key cryptography: the Fields Institute for Research in the Mathematical Sciences, Toronto, Ontario, June 12–17, 1999*, book of preprints distributed at the conference, 1999.
- [25] Christos M. Papadimitriou, *Computational complexity*, Addison-Wesley, Reading, Massachusetts, 1994. ISBN 0-201-53082-1. MR1251285 (95f:68082)
- [26] John H. Reif (editor), *Synthesis of parallel algorithms*, Morgan Kaufman, San Mateo, California, 1993. ISBN 1-55860-135-X. MR1212591 (94c:68086)

- [27] J. Barkley Rosser, Lowell Schoenfeld, *Approximate formulas for some functions of prime numbers*, Illinois Journal of Mathematics **6** (1962), 64–94. ISSN 0019-2082. URL: <http://cr.yp.to/bib/entries.html#1962/rosser>. MR0137689 (25:1139)
- [28] Arnold Schönhage, Volker Strassen, *Schnelle Multiplikation großer Zahlen*, Computing **7** (1971), 281–292. ISSN 0010-485X. URL: <http://cr.yp.to/bib/entries.html#1971/schoenrage-mult>. MR0292344 (45:1431)
- [29] Jonathan P. Sorenson, *A sublinear-time parallel algorithm for integer modular exponentiation*, in [24] (1999).
- [30] Jonathan P. Sorenson, Ian Parberry, *Two fast parallel prime number sieves*, Information and Computation **144** (1994), 115–130. ISSN 0890-5401. MR1294306 (95h:11097)
- [31] Larry Stockmeyer, Uzi Vishkin, *Simulation of parallel random access machines by circuits*, SIAM Journal on Computing **13** (1984), 409–422. ISSN 0097-5397. MR0739997 (85g:68018)
- [32] Stephen R. Tate, *Newton iteration and integer division*, in [26] (1993), 539–572. MR1212591 (94c:68086)
- [33] Alf van der Poorten, Andreas Stein (editors), *High primes and misdemeanours: lectures in honour of the 60th birthday of Hugh Cowie Williams*, American Mathematical Society, Providence, 2004. ISBN 0-8218-3353-7. MR2005b:11003
- [34] Uzi Vishkin, *Advanced parallel prefix-sums, list ranking and connectivity*, in [26] (1993), 215–257. MR1212591 (94c:68086)
- [35] Joachim von zur Gathen, *Computing powers in parallel*, SIAM Journal on Computing **16** (1987), 930–945. MR0908877 (89j:68060)

DEPARTMENT OF MATHEMATICS, STATISTICS, AND COMPUTER SCIENCE (M/C 249), THE UNIVERSITY OF ILLINOIS AT CHICAGO, CHICAGO, IL 60607–7045

E-mail address: djb@cr.yp.to

DEPARTMENT OF COMPUTER SCIENCE AND SOFTWARE ENGINEERING, BUTLER UNIVERSITY, 4600 SUNSET AVENUE, INDIANAPOLIS, INDIANA 46208

E-mail address: sorenson@butler.edu