

Efficient and Formally Proven Reduction of Large Integers by Small Moduli

Technical report: ICIS-R07033, February, Radboud University Nijmegen, 2008

Luc Rutten	Marko van Eekelen
Delft Development Lab	Institute for Computing and Information Sciences
Tivoli, Software Group, IBM	Radboud University Nijmegen
luc.rutten@nl.ibm.com	marko@cs.ru.nl

February 6, 2008

Abstract

On w -bit processors which are much faster at multiplying two w -bit integers than at dividing $2w$ -bit integers by w -bit integers, reductions of large integers by moduli M smaller than 2^{w-1} are often implemented sub-optimally, leading applications to take excessive processing time.

We present a modular reduction algorithm implementing division by a modulus through multiplication by a reciprocal of that modulus, a well-known method for moduli larger than 2^{w-1} . We show that application of this method to smaller moduli makes it possible to express certain modular sums and differences without having to compensate for word overflows.

By embedding the algorithm in a loop and applying a few transformations to the loop, we obtain an algorithm for reduction of large integers by moduli up to 2^{w-1} . Implementations of this algorithm can run considerably faster than implementations of similar algorithms that allow for moduli up to 2^w . This is substantiated by measurements on processors with relatively fast multiplication instructions.

It is notoriously hard to specify efficient mathematical algorithms on the level of abstract machine instructions in an error-free manner. In order to eliminate the chance of errors as much as possible, we have created formal correctness proofs of our algorithms, checked by a mechanized proof assistant.

1 Introduction

There are many applications of *modular reductions*, which are computations of *residues* $x \bmod M$ (“ x modulo M ”) where $x \in \mathbb{Z}$ and $M \in \mathbb{Z}^+$. The positive integer M is called the *modulus*. The computation of $x \bmod M$ is called the *reduction of x by M* . The residue $x \bmod M$ is sometimes called the *remainder (after division of x by M)*.

Often, applications of modular reductions employ moduli which may be (much) larger than 2^w , where $w \in \mathbb{Z}^+$ is the *word size* of a computer processor. Some applications though only employ moduli $M \leq 2^{w-1}$,

which we call *small moduli*. Such applications include small modulus specializations of large integer reduction [GMP], modular exponentiation, and modular multiplication, for example employed in residue arithmetic and in Garner's algorithm [Knu98]. Implementations of such small-modulus specializations benefit from efficient reduction by small moduli.

Unsigned integer arithmetic sets and operations are presented in section 2. In section 3 we present an efficient small modulus reduction algorithm expressed in terms of these sets and operations. The algorithm – ModRed – can be employed to compute residues $x \bmod M$ for $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$ and $x \in \mathbb{Z}_{2^{\lceil \lg M \rceil + w}}$.

In section 4 we show that ModRed can be embedded in a loop in order to reduce large integers. Some transformations are applied to the loop to obtain a more efficient algorithm – MultiRed – which can be employed to compute residues $x \bmod M$ for $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$ and $x \in \mathbb{Z} \mid 0 \leq x$.

Algorithms ModRed and MultiRed contain low-word multiplications and high-word multiplications but no divisions except by powers of 2. This indicates that they can be relatively efficiently implemented on processors where division of $2w$ -bit words by w -bit words is expensive relative to low-word and high-word multiplications. While the number and kind of operations of an algorithm may give an impression of its efficiency, many other factors affect real life performance. In section 5, we therefore present performance measurements of ModRed and MultiRed implementations. The performance of these implementations is compared with the performance of a similar implementation from [GMP], whose method for reducing an unsigned double word by an unsigned word is based on an algorithm proposed in [GM94].

It is notoriously hard to specify mathematical algorithms like ModRed and MultiRed in an error-free manner. In order to obtain a high degree of certainty that ModRed and MultiRed are free of errors, we have created formal correctness proofs of the algorithms, checked by a mechanized proof assistant. This is discussed in section 6.

Related and future work are discussed in sections 7 and 8 while conclusions are presented in section 9.

2 Preliminaries

The algorithms proposed in section 3 and section 4 are expressed in terms of the standard unsigned integer arithmetic sets and operations (on the level of abstract machine instructions) shown in the current section, which may be skimmed by readers familiar with these concepts.

In this text, \mathbb{Z}^+ is defined as $\{i \in \mathbb{Z} \mid 0 < i\}$. For each $n \in \mathbb{Z}$, \mathbb{Z}_n is defined as $\{i \in \mathbb{Z} \mid 0 \leq i < n\}$. For each $x \in \mathbb{Z}$ and $M \in \mathbb{Z}^+$, $\lfloor \frac{x}{M} \rfloor$ is defined as the unique integer $q \in \mathbb{Z}$ and $x \bmod M$ is defined as the unique integer $r \in \mathbb{Z}_M$ for which it holds: $x = qM + r$. Expressions $\lfloor \frac{x}{M} \rfloor$ are sometimes written like $\lfloor x/M \rfloor$. We will write mod as a left associative operator with the same precedence as the multiplication operator. For each $x \in \mathbb{Z}^+$, the *binary logarithm of x rounded up to the nearest integer*, $\lceil \lg x \rceil$, is defined as the (unique) nonnegative integer i such that $\lfloor \frac{2^i}{2} \rfloor < x \leq 2^i$. For each $x \in \mathbb{Z}^+$, $\lfloor \lg x \rfloor$ is similarly defined as the (unique) nonnegative integer i such that $2^i \leq x < 2^{i+1}$.

For each word size $w \in \mathbb{Z}^+$, each $x \in \mathbb{Z}_{2^w}$ is called a *w-bit word* in *unsigned w-bit integer arithmetic*. Such a word x may be identified with a tuple of w bits: $(\lfloor \frac{x}{2^i} \rfloor \bmod 2)_{i=0}^{w-1} \in \mathbb{Z}_2^w$. A *w-bit processor* provides

instructions for direct manipulation of representations of such tuples. For $x, y \in \mathbb{Z}_{2^w}$ and $i \in \mathbb{Z}_w$, instructions for computation of $(x + y) \bmod 2^w$ (addition), $(x - y) \bmod 2^w$ (subtraction), $xy \bmod 2^w$ (low-word multiplication), $\lfloor \frac{xy}{2^w} \rfloor$ (high-word multiplication), $2^i x \bmod 2^w$ (left shift), and $\lfloor \frac{x}{2^i} \rfloor$ (right shift) are commonly provided. Each of these expressions usually takes just a single instruction, while some processors are able to compute $(xy \bmod 2^w, \lfloor \frac{xy}{2^w} \rfloor)$ with a single instruction. Conditional operations are also commonly provided, usually taking a comparison instruction and a conditional move or branch instruction. An additional unconditional branch instruction may also be involved.

The greatest integer that can be represented in unsigned w -bit integer arithmetic is $2^w - 1$. Greater integers are said to *overflow* a word, and $\lfloor \frac{x}{2^w} \rfloor$ is called the *nonzero word overflow* of such an integer x .

Tuples of words $(\lfloor \frac{x}{2^{wi}} \rfloor \bmod 2^w)_{i=0}^{k-1} \in \mathbb{Z}_{2^w}^k$ can represent *multi-word integers* $x \in 2^{wk}$, also called *large integers* because they may be much larger than $2^w - 1$. For example, each *double word* $x \in \mathbb{Z}_{2^{2w}}$ can be represented with a *low word* $x \bmod 2^w$ and a *high word* $\lfloor \frac{x}{2^w} \rfloor$. Accessing the words of a multi-word integer and composing multi-word integers from words may require the use of *load* and *store* instructions.

3 A new algorithm for small modulus reduction

In this section, we present a new algorithm which can be employed to compute residues $x \bmod M$ for $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$ and $x \in \mathbb{Z}_{2^{\lceil \lg M \rceil + w}}$. The algorithm – ModRed – is defined on the level of an abstract machine using instructions introduced in section 2. Some applications of ModRed are mentioned in section 1.

The main idea behind the ModRed algorithm is to compute $x \bmod M$ as $(\lfloor \frac{x}{a} \rfloor a \bmod M + x \bmod a \bmod M) \bmod M$ where $a = 2^{\lceil \lg M \rceil}$, and to approximate $\lfloor \frac{x}{a} \rfloor a \bmod M$ by $d = \lfloor \frac{x}{a} \rfloor a - qM$ where $q = \lfloor \lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{M} \rfloor / b \rfloor$ and $b = 2^w$.

The value of q is only 0 or 1 smaller than $\lfloor \lfloor \frac{x}{a} \rfloor a / M \rfloor$, as

$$\forall y, a, b \in \mathbb{Z}^+ \quad \forall x \in \mathbb{Z}_{ab+a} : \lfloor \frac{\lfloor \frac{x}{a} \rfloor a}{y} \rfloor - 1 \leq \lfloor \frac{\lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{y} \rfloor}{b} \rfloor \leq \lfloor \frac{\lfloor \frac{x}{a} \rfloor a}{y} \rfloor \quad (1)$$

Therefore, d can be equal to $\lfloor \frac{x}{a} \rfloor a \bmod M$ or $\lfloor \frac{x}{a} \rfloor a \bmod M + M$. When $M \leq 2^{w-1}$, both values lie in \mathbb{Z}_{2^w} so word overflows do not occur: $\lfloor \frac{d}{2^w} \rfloor = 0$. The value of $\lfloor \frac{x}{a} \rfloor a \bmod M$ is therefore easy to derive from d with a conditional subtraction. To the resulting value $\lfloor \frac{x}{a} \rfloor a \bmod M$, $x \bmod a \bmod M$ can be added modulo M to obtain $x \bmod M$, again using only w -bit integer arithmetic and without word overflows. If M would have been greater than 2^{w-1} , word overflows would have had to be handled with additional operations. Additional operations are also needed if $x \bmod M$ is approximated more directly, using equation 65 of section 7. That can be seen in the `udiv_qrnd_preinv1` macro of [GMP]

By choosing $b = 2^w$, the low word of $\lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{M} \rfloor$ is not needed for computation of q . When $\lceil \lg M \rceil < w - 1$, then on processors with separate low-word and high-word multiplication instructions, this saves a low-word multiplication w.r.t. Barrett's algorithm [Bar87], which is reviewed in section 7. Furthermore, when $b = 2^w$, the high word of $\lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{M} \rfloor$ can on most

processors be obtained without carrying out any operation after the multiplication because the high word is already present in a machine register after a high-word (or double-word) multiplication. This contrasts with the left shift, right shift, and addition operations which are usually employed to divide a double word by $2^{\lceil \lg M \rceil + 1}$ on a w -bit processor when $\lceil \lg M \rceil < w - 1$.

These advantages of choosing $b = 2^w$ are shared with the algorithm proposed in section 8 of [GM94]. The low values of M that ModRed applies to $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$ give rise to an additional advantage relative to this algorithm, which applies to $M \in \mathbb{Z} \mid 1 \leq M < 2^w$: because ModRed does not apply to moduli $M \in \mathbb{Z} \mid 2^{w-1} < M < 2^w$, it needs fewer operations for taking care of word overflows.

Algorithm 1 ModRed

Inputs: $M \in \mathbb{Z}_{2^w}$, $p, t \in \mathbb{Z}_w$, $M', v, u \in \mathbb{Z}_{2^w}$

Output: $\text{ModRed}(M, p, t, M', v, u) = r'''$

where $s, s', h, h', q, q', y, d, r, r', r'', r''' \in \mathbb{Z}_{2^w}$ are defined with

$$s = \lfloor \frac{u}{2^p} \rfloor \quad (2)$$

$$s' = 2^p s \bmod 2^w \quad (3)$$

$$h = 2^t v \bmod 2^w \quad (4)$$

$$h' = (h + s) \bmod 2^w \quad (5)$$

$$q = \lfloor \frac{h' M'}{2^w} \rfloor \quad (6)$$

$$q' = (q + h') \bmod 2^w \quad (7)$$

$$y = q' M \bmod 2^w \quad (8)$$

$$d = (s' - y) \bmod 2^w \quad (9)$$

$$r = (u - y) \bmod 2^w \quad (10)$$

$$r' = \begin{cases} r & \text{if } d < M \\ (r - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (11)$$

$$r'' = \begin{cases} r' & \text{if } r' < M \\ (r' - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (12)$$

$$r''' = \begin{cases} r'' & \text{if } r'' < M \\ (r'' - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (13)$$

Let $C \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \rightarrow \mathbb{Z}_w \times \mathbb{Z}_w \times \mathbb{Z}_{2^w}$ be defined with

$$C(M) = (\lceil \lg M \rceil, w - \lceil \lg M \rceil - \lfloor \frac{1}{M} \rfloor, \lfloor \frac{2^{\lceil \lg M \rceil + w}}{M} \rfloor - 2^w) \quad (14)$$

where $M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\}$. Then, splitting the integer to be reduced (x) in its high and low word as explained in section 2, the correctness of algorithm ModRed is expressed by the following theorem:

Theorem 3.1 (Correctness of ModRed)

$\forall M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \forall x \in \mathbb{Z}_{2^{\lceil \lg M \rceil + w}} :$

$$\text{ModRed}(M, C(M), \lfloor \frac{x}{2^w} \rfloor, x \bmod 2^w) = x \bmod M \quad (15)$$

Determination of the components of $C(M)$ may take place according to methods described in e.g. [Knu98] or implemented in e.g. [GMP]. Computation of such values can be time consuming in comparison with evaluating an application of ModRed. It is therefore recommended to use ModRed only if ModRed($M, C(M), \lfloor \frac{x}{2^w} \rfloor, x \bmod 2^w$) for a single value of M is to be computed for many values of x . A suitable minimum number of such values can be determined with performance measurements. When a compiler uses ModRed to generate better code for $x \bmod M$ where M is known at compile-time, the minimum number may be equal to 1 as the compiler can precompute $C(M)$.

The outcome of one ModRed application may be passed to the next application. This can be seen for example in the algorithm proposed in section 4.

4 Reduction of large integers by a small modulus

In this section, it will be seen that algorithm ModRed can be placed in a loop to obtain a multi-word reduction algorithm. We will demonstrate how that algorithm can be transformed to a more efficient multi-word reduction algorithm, called MultiRed.

Rationale

Algorithm ModRed can be employed in a loop of k iterations in order to reduce an integer in $\mathbb{Z}_{2^{wk}}$ by a modulus $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$. This is illustrated by the reduction of a 3-word integer $x = x_2 2^{w \cdot 2} + x_1 2^w + x_0$ where $x_2, x_1, x_0 \in \mathbb{Z}_{2^w}$: $x \bmod M = ((0 \cdot 2^w + x_2) 2^w + x_1) 2^w + x_0 \bmod M = (((0 \cdot 2^w + x_2) \bmod M) 2^w + x_1) \bmod M) 2^w + x_0 \bmod M = \text{ModRed}(M, C(M), \text{ModRed}(M, C(M), \text{ModRed}(M, C(M), 0, x_2), x_1), x_0)$. In the first loop iteration, the innermost ModRed application is evaluated, while in the the last loop iteration, the outermost ModRed application is evaluated. Performance measurements of ModRed in a multi-word reduction program (see section 5) show it to be slower than or about as fast as the `udiv_qrmod_preinv1` macro of [GMP]. By applying some transformations to the ModRed loop, it is possible obtain a more efficient multi-word reduction program.

Transformations

In the descriptions of the transformations, the loop iterations employed to reduce $x \in \mathbb{Z}_{2^{wn}}$ will be numbered $n-1, n-2, \dots, 0$. For each $i \in \mathbb{Z}_n$, x_i will be defined as $\lfloor \frac{x}{2^{wi}} \rfloor \bmod 2^w$, so $x = \sum_{i=0}^n 2^{wi} x_i$. The ModRed variables of each iteration i will get suffix i . We define r_n''' as 0. The inputs of iteration i are defined as $v_i = r_{i+1}'''$ and $u_i = x_i$. The value of r_0''' is computed in iteration 0, i.e. the last iteration. The value of r_0''' equals $x \bmod M$.

At the end of each iteration i , r_i''' is derived from r_i'' using a conditional subtraction. The r_i''' value then enters the next loop iteration as v_{i-1} (if $i > 0$). Because $r_i'' = r_i''' \vee r_i'' = r_i'' + M$ and because $r_i'' < 2^{\lceil \lg M \rceil}$, we can let v_{i-1} be defined as r_i'' instead of r_i''' . That is because for all $u \in \mathbb{Z}_{2^w}$ and for all $v \in \mathbb{Z}_{2^{\lceil \lg M \rceil}}$ such that $v + M < 2^{\lceil \lg M \rceil}$, it holds $\text{ModRed}(M, C(M), v + M, u) = (2^w(v + M) + u) \bmod M = (2^w v + u) \bmod M$.

$M = \text{ModRed}(M, C(M), v, u)$. By defining v_{i-1} as r_i'' instead of r_i''' , one conditional subtraction is saved per loop iteration (except for the last iteration) as was already suggested at the end of section 3. Therefore, let's define r_k'' as 0 and let's define v_i as r_{i+1}'' .

After transforming the loop by replacing $v_i = r_{i+1}'''$ with $v_i = r_{i+1}''$, each loop iteration ends with two conditional subtractions. Also, each loop starts with a right shift of x_i , which is preceded by loading x_i from memory. The instruction level parallelism is increased by moving the two conditional subtractions at the end of iteration i to the beginning of the next iteration, $i - 1$. After this transformation, the conditional subtractions can in principle be computed in parallel with the right shift (and subsequent left shift).

Now h_i' depends on h_i and s_i , where h_i depends on r_{i+1}'' , which depends on r_{i+1}' , which depends on r_{i+1} . On the other hand, s_i just depends on x_i . Therefore, it is likely that s_i will have been computed well before h_i . The conditional subtraction from r_{i+1}' is then replaced by a conditional subtraction from s_i . After this transformation, the two dependency chains of the (new) components of h_i' have equal lengths, increasing instruction level parallelism.

In the second conditional subtraction of the resulting algorithm there is a comparison between r' and M . This can be optimized for some processors which take less time to evaluate r'' if $r' < c$ than if $r' \geq c$, especially for values of M close to $2^{\lceil \lg M \rceil - 1} + 1$. The transformation is to replace M by $c = 2^{\lceil \lg M \rceil}$. This does not cause overflows in the rest of the computation. So, no extra corrections are needed.

After performing the four transformations some extra equations have to be added after the loop to make up for the moved conditional subtractions. This results in the algorithm MultiRed given below.

For expressing the correctness of the algorithm MultiRed we need another auxiliary definition. $C' \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \rightarrow \mathbb{Z}_w \times \mathbb{Z}_w \times \mathbb{Z}_{2^w} \times \mathbb{Z}_{2^w} \times \mathbb{Z}_{2^w}$ is defined with

$$C'(M) = (C(M), 2^{\lceil \lg M \rceil}, 2^{w - \lceil \lg M \rceil - \lfloor \frac{1}{M} \rfloor} M \bmod 2^w) \quad (32)$$

where $M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\}$. Using this definition we formulate the MultiRed correctness theorem:

Theorem 4.1 (Correctness of MultiRed)

$$\forall M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \forall n \in \mathbb{N} \forall x \in \mathbb{Z}_{2^{nw}} : \\ \text{MultiRed}(M, C'(M), n, x) = x \bmod M \quad (33)$$

4.1 Parallellization

It is possible to compute $x \in \mathbb{Z}_{2^{2kw}} \bmod M$ using the formula $x \bmod M = ((\lfloor \frac{x}{2^{kw}} \rfloor \bmod M) \cdot (2^{kw} \bmod M) + (x \bmod 2^{kw}) \bmod M) \bmod M$. The components $\lfloor \frac{x}{2^{kw}} \rfloor$ and $x \bmod 2^{kw}$ of x can be reduced (e.g. with MultiRed) in parallel, independently of each other. The value of $2^{kw} \bmod M$ can be determined with fewer than $2 \lceil \lg k \rceil$ low-word multiplications, $2 \lceil \lg k \rceil$ high-word multiplications, and $2 \lceil \lg k \rceil$ reductions of integers in $\mathbb{Z}_{2^{\lceil \lg M \rceil + w}}$ by M e.g. employing ModRed. When k is large, computation of $2^{kw} \bmod M$ therefore only takes a small fraction of the time to reduce the components of x . Reducing the components in parallel and combining the residues may then take only a bit more than half the (wall-clock) time taken to reduce the components after one another.

Algorithm 2 MultiRed

Inputs: $M \in \mathbb{Z}_{2^w}$, $p, t \in \mathbb{Z}_w$, $M', c, M'' \in \mathbb{Z}_{2^w}$, $n \in \mathbb{N}$, $x \in \mathbb{Z}_{2^{nw}}$

Output: $\text{MultiRed}(M, p, t, M', c, M'', n, x) = r_0'''$

where $d_n, r_n \in \mathbb{Z}_{2^w}$ are defined with

$$d_n = 0 \tag{16}$$

$$r_n = 0 \tag{17}$$

and where for $i = n-1, n-2, \dots, 0$, x_i is defined as $\lfloor \frac{x}{2^{wi}} \rfloor \bmod 2^w$ and $r'_i, s_i, s'_i, s''_i, h_i, h'_i, q_i, q'_i, y_i, d_i, r_i \in \mathbb{Z}_{2^w}$ are defined with

$$r'_i = \begin{cases} r_{i+1} & \text{if } d_{i+1} < M \\ (r_{i+1} - M) \bmod 2^w & \text{otherwise} \end{cases} \tag{18}$$

$$s_i = \lfloor \frac{x_i}{2^p} \rfloor \tag{19}$$

$$s'_i = 2^p s_i \bmod 2^w \tag{20}$$

$$s''_i = \begin{cases} s_i & \text{if } r'_i < c \\ (s_i - M'') \bmod 2^w & \text{otherwise} \end{cases} \tag{21}$$

$$h_i = 2^t r'_i \bmod 2^w \tag{22}$$

$$h'_i = (h_i + s''_i) \bmod 2^w \tag{23}$$

$$q_i = \lfloor \frac{h'_i M'}{2^w} \rfloor \tag{24}$$

$$q'_i = (q_i + h'_i) \bmod 2^w \tag{25}$$

$$y_i = q'_i M \bmod 2^w \tag{26}$$

$$d_i = (s'_i - y_i) \bmod 2^w \tag{27}$$

$$r_i = (x_i - y_i) \bmod 2^w \tag{28}$$

and where $r'_0, r''_0, r'''_0 \in \mathbb{Z}_{2^w}$ are defined with

$$r'_0 = \begin{cases} r_0 & \text{if } d_0 < M \\ (r_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \tag{29}$$

$$r''_0 = \begin{cases} r'_0 & \text{if } r'_0 < c \\ (r'_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \tag{30}$$

$$r'''_0 = \begin{cases} r''_0 & \text{if } r''_0 < M \\ (r''_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \tag{31}$$

5 Performance

The ModRed and MultiRed algorithms introduced in section 3 and section 4 have been designed with efficiency in mind. To gain some insight in the algorithms' efficiency, we will compare run-times of implementations of the algorithms, each allowing for moduli up to 2^{w-1} , with run-times of the `mpn_mod_1` function of [GMP], which uses the `udiv_qrnd_preinv1` macro allowing for moduli between 2^{w-1} and 2^w . A variation of `mpn_mod_1` is also measured.

We benchmarked the reduction of a multi-word integer (in particular we took $\sum_{i=0}^{\lfloor\sqrt{s}\rfloor-1} (16807^i \bmod (2^{31}-1)) \cdot 2^{wi}$ where s denotes the processor speed in Hz) by several moduli smaller than 2^{w-1} (in particular $2^{w-1} - i \lfloor \frac{2^{w-1}}{\lfloor\sqrt{s}\rfloor} \rfloor$ for each $i = 1.. \lfloor\sqrt{s}\rfloor - 1$). Implementations of ModRed and MultiRed were compared with an implementation (denoted as `mm1`) directly based on the `mpn_mod_1` function from the GNU Multiple Precision Arithmetic Library [GMP]

The `mpn_mod_1` function uses the `udiv_qrnd_preinv1` macro which is based on an algorithm of [GM94]. It can be employed to reduce integers in $\mathbb{Z}_{2^{2w}}$ by moduli $M \in \mathbb{Z} \mid 2^{w-1} \leq M < 2^w$. The following equations reflect that part of the macro which computes $r' = (2^w v + u) \bmod M$. Equations to determine $\lfloor \frac{2^w v + u}{M} \rfloor$ are not included because this quotient is not needed for computing residues. In the pseudocode, M' denotes $\lfloor \frac{2^{2w}-1}{M} \rfloor - 2^w$. The integers Y , Z , and Z' are double words. When $\lfloor \frac{Z}{2^w} \rfloor = 0$, $r = Z \bmod 2^w$ so $\lfloor \frac{Z'}{2^w} \rfloor$ does not really have to be computed.

$$q = \lfloor \frac{vM'}{2^w} \rfloor \tag{a}$$

$$q' = (q + v) \bmod 2^w \tag{b}$$

$$Y = q' M \tag{c}$$

$$Z = (2^w v + u) - Y \tag{d}$$

$$Z' = \begin{cases} Z & \text{if } \lfloor \frac{Z}{2^w} \rfloor = 0 \\ Z - M & \text{otherwise} \end{cases} \tag{e}$$

$$r = \begin{cases} Z' \bmod 2^w & \text{if } \lfloor \frac{Z'}{2^w} \rfloor = 0 \\ Z' \bmod 2^w - M & \text{otherwise} \end{cases} \tag{f}$$

$$r' = \begin{cases} r & \text{if } r < M \\ r - M & \text{otherwise} \end{cases} \tag{g}$$

We also compared MultiRed with with a variation `mm1'` of the `mm1` implementation. We made this variation in order to view the effects of some transformations on `mm1` that are analogous to the transformations on the ModRed loop to obtain MultiRed. Preliminary performance measurements indicated that MultiRed can be considerably more efficient than ModRed, leading to the expectation that `mm1'` is more efficient than `mm1`. The `mm1` implementation is transformed to the `mm1'` implementation by moving the last three equations ((e), (f), and (g)) of the `mm1` loop body to the front of the loop, and by putting some operations in front of the loop and after the loop to compensate for this.

Note that `mm1` and `mm1'` implement reduction of multi-word integers by moduli smaller than 2^{w-1} , rather than by moduli between 2^{w-1} and 2^w . To be able to use the `udiv_qrnd_preinv1` macro, a double word is reduced by $2^{w-\lceil \lg M \rceil} M$ rather than by M in each iteration of the `mm1` an `mm1'` inner loops. The resulting residue in \mathbb{Z}_{2^w} is reduced by M to

obtain the final result. While this reduction may take a (costly) division operation, the number of loop iterations is rather large so the time spent by the division operation is rather small in comparison with the time spent by the loop iteration.

The benchmarks were carried out on two 64-bit processors: a 1.6 GHz AMD Turion* 64 X2 ML-50 [AMD], and a 2 GHz IBM PowerPC† 970FX [IBM]. 64-Bit processors are often employed in workstations and they are beginning to emerge in personal computers. Even to reduce a multi-word integer by a 32-bit word, it is profitable to use a 64-bit ModRed implementation on a 64-bit processor, as the number of 64-bit ModRed iterations needed to reduce a $64n$ -bit multi-word integer is equal to n , while the number of 32-bit ModRed iterations needed to reduce the same $64n$ -bit multi-word integer is equal to $2n$.

Each benchmark was compiled with a version of the GNU C compiler [GCC], using option `-O3` for high optimization: gcc 4.1.2 with the Turion 64 X2 ML-50, and gcc 4.3.0 with the PowerPC 970FX. The benchmarks were carried out on quiet systems, with no other processor, memory, or disk intensive tasks at hand. The benchmarks were repeated in the course of several days and at different times of a day.

The following table lists (average) numbers of clock cycles spent by a single iteration of four implementations for reducing multi-word integers by moduli smaller than 2^{w-1} .

	Turion 64 X2	PowerPC 970FX
mm1	<i>24.5</i>	40.6
mm1'	25.0	<i>38.4</i>
ModRed	24.0	47.6
MultiRed	<i>18.0</i>	<i>36.2</i>
$\lfloor \frac{\text{bestmm1} \cdot 100}{\text{bestModRedfamily}} \rfloor / 100$	1.36	1.06

The performance measurement of MultiRed shows only a small improvement (1.06) for the PowerPC 970FX processor. In our opinion this is caused by internal scheduling and instruction level parallelism differences between the considered processors. It suggests that a different order of the calculations might give better performance on the PowerPC. After several experiments we found a calculation sequence showing an improvement which is more alike the improvement found with the Turion. The calculation scheme does not only concern changing the order of the instructions. It also encompasses changes in the actual instructions employed and even in the number of instructions. For this reason it is justifiable to give the resulting variant of MultiRed a separate name – MultiRed' – and to specify and verify the algorithm separately. This variation of MultiRed is described in the next section.

5.1 A variation of MultiRed

The lower performance of MultiRed on the PowerPC 970FX suggests that the first conditional subtraction of its loop may be a bottleneck on some processors. It is possible to apply a few transformations to the MultiRed loop in order to obtain an algorithm which runs more efficiently on some

*AMD Turion is a trademark of Advanced Micro Devices, Inc.

†PowerPC is a trademark of International Business Machines Corporation

of those processors. Surprisingly, this involves the introduction of yet another conditional subtraction.

In each loop iteration, we determine $f_i = x_i \bmod 2^{\lceil \lg M \rceil}$ (e.g. with $(x_i - s'_i) \bmod 2^w$).

At the end of each loop iteration, we determine g_i , which equals r_i if $f_i < M$, and $(r_i - M) \bmod 2^w$ otherwise. Instead of r_i , we pass g_i to the next loop iteration (if any).

Finally, we replace the comparison $d_{i+1} < M$ with $d_{i+1} < c$, which may lead to a subtraction being carried out less often on average when $c = 2^{\lceil \lg M \rceil}$, like described at the end of section 3.

The correctness theorem of MultiRed' much resembles the MultiRed correctness theorem, once more using the function C' defined with equation 32.

Theorem 5.1 (Correctness of MultiRed')

$$\forall M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \forall n \in \mathbb{N} \forall x \in \mathbb{Z}_{2^{nw}} : \\ \text{MultiRed}'(M, C'(M), n, x) = x \bmod M \quad (52)$$

Below, we give the performance measurements including the new algorithm MultiRed'.

	Turion 64 X2	PowerPC 970FX
mm1	24.5	40.6
mm1'	25.0	38.4
ModRed	24.0	47.6
MultiRed	18.0	36.2
MultiRed'	21.0	30.8
$\lfloor \frac{\text{bestmm1} \cdot 100}{\text{bestModRedfamily}} \rfloor / 100$	1.36	1.24

The table shows that for the PowerPC MultiRed' indeed improves over MultiRed as expected. For the Turion however, its performance is slightly worse than that of MultiRed.

5.2 Evaluation of performance measurements

The benchmark results indicate that implementation of MultiRed and MultiRed' can be profitable on processors of different kinds. It can also be seen that mm1' is not much faster than mm1. The PowerPC 970FX has no instructions to divide 128-bit integers by 64-bit integers. The Turion 64 X2 has a 128-bit by 64-bit unsigned integer division instruction but its instruction for multiplying two unsigned 64-bit integers to obtain a 128-bit result is much faster: the MultiRed implementation ran more than four times as fast as an implementation using the 128-bit integer by 64-bit integer division instruction. For these reasons, benchmarks using double-word by single-word divisions have not been listed in the table.

Besides on the number and kind of operations of an algorithm, the execution time of an implementation of the algorithm depends on factors like instruction arities, number of registers, handling of constants, pipelining, the programming language, compiler, and compiler options. Because the efficiency aspects of these factors can hardly be formalized, benchmarks can be used to measure an implementation's execution time. Beware though that benchmark results depend on uncalculated factors. In order to determine the relative efficiency of ModRed, MultiRed, and MultiRed'

Algorithm 3 MultiRed'

Inputs: $M \in \mathbb{Z}_{2^w}$, $p, t \in \mathbb{Z}_w$, $M', c, M'' \in \mathbb{Z}_{2^w}$, $n \in \mathbb{N}$, $x \in \mathbb{Z}_{2^{nw}}$

Output: MultiRed'(M, p, t, M', c, M'', n, x) = r_0'''

where $d_n, g_n \in \mathbb{Z}_{2^w}$ are defined with

$$d_n = 0 \quad (34)$$

$$g_n = 0 \quad (35)$$

and where for $i = n - 1, n - 2, \dots, 0$, x_i is defined as $\lfloor \frac{x}{2^{wi}} \rfloor \bmod 2^w$ and $r'_i, s_i, s'_i, f_i, s''_i, h_i, h'_i, q_i, q'_i, y_i, d_i, r_i, g_i \in \mathbb{Z}_{2^w}$ are defined with

$$r'_i = \begin{cases} g_{i+1} & \text{if } d_{i+1} < c \\ (g_{i+1} - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (36)$$

$$s_i = \lfloor \frac{x_i}{2^p} \rfloor \quad (37)$$

$$s'_i = 2^p s_i \bmod 2^w \quad (38)$$

$$f_i = (x_i - s'_i) \bmod 2^w \quad (39)$$

$$s''_i = \begin{cases} s_i & \text{if } r'_i < c \\ (s_i - M'') \bmod 2^w & \text{otherwise} \end{cases} \quad (40)$$

$$h_i = 2^t r'_i \bmod 2^w \quad (41)$$

$$h'_i = (h_i + s''_i) \bmod 2^w \quad (42)$$

$$q_i = \lfloor \frac{h'_i M'}{2^w} \rfloor \quad (43)$$

$$q'_i = (q_i + h'_i) \bmod 2^w \quad (44)$$

$$y_i = q'_i M \bmod 2^w \quad (45)$$

$$d_i = (s'_i - y_i) \bmod 2^w \quad (46)$$

$$r_i = (x_i - y_i) \bmod 2^w \quad (47)$$

$$g_i = \begin{cases} r_i & \text{if } f_i < M \\ (r_i - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (48)$$

and where $r'_0, r''_0, r'''_0 \in \mathbb{Z}_{2^w}$ are defined with

$$r'_0 = \begin{cases} g_0 & \text{if } d_0 < M \\ (g_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (49)$$

$$r''_0 = \begin{cases} r'_0 & \text{if } r'_0 < c \\ (r'_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (50)$$

$$r'''_0 = \begin{cases} r''_0 & \text{if } r''_0 < M \\ (r''_0 - M) \bmod 2^w & \text{otherwise} \end{cases} \quad (51)$$

on processors of other kinds than the ones mentioned in this section, or using compilers other than the ones mentioned, the benchmarks should be carried out with those processors and compilers.

6 Correctness proofs

Writing an algorithm conform to a given specification is an error-prone task. Some errors may be found by testing an implementation of the algorithm but typically some errors go unnoticed because they have only a very small chance of showing up in tests. Algorithms involving integer arithmetic are no exception; to the contrary, they may contain errors which show up in an extremely small fraction of all possible tests. Such errors can be avoided by *formally proving* the correctness of algorithms. It is possible to obtain a high degree of assurance that a “correctness proof” is not in error itself by constructing or checking it with a computer.

6.1 The choice of proof assistant

Several computerized proof assistants are available, e.g. Coq, PVS, and Isabelle. For a comparison see [Wie03].

We have chosen to use the Coq proof assistant [Coq] here because it produces explicit proof terms which can be checked independently with a relatively simple proof checker. This satisfies the *de Bruijn criterion*, named after the Dutch mathematician N.G. de Bruijn, who is considered to be the principal founder of machine verification of formalized proofs. He emphasized the following criterion [dB70] for reliable automated proof-checkers: *their programs must be small, so small that a human can (easily) verify the code by hand*. We feel that the trust which is required for mathematical algorithms like ModRed is best obtained by using a proof checker that satisfies this criterion. Of course this does not give 100% certainty of correctness since to a certain degree, the correctness of the theorems proved with Coq depends on the correctness of the Coq implementation and of the correct operation of the computer which runs Coq.

6.2 The overall proof methodology

The correctness of the ModRed, MultiRed, and MultiRed' algorithms, defined in this text as theorem 3.1, 4.1 and 5.1, has been formally verified with a computer [Rut], using the Coq proof assistant [Coq].

Before giving the essential parts of the proofs in section 6.3, we will first explain the methodology with which we maintained the correspondence between the proofs in this paper and the Coq proofs. After that we will shortly discuss the overall structure of the proofs.

Guaranteeing the connection between computerized proof and paper proof

A computerized proof has a very high degree of reliability. It may however occur that the properties that are proven do not fully correspond to the properties that have to be proven. When this happens, the proof is not wrong but it is the wrong proof. There are many syntactical differences between the Coq level and a general mathematical description. An error is easily made. Therefore we pay extra attention to guaranteeing the

correspondence between the two levels. Below we explain an approach to avoid discrepancies between the Coq proof and the mathematical proof.

As an attempt to reduce the number of errors in the numbered definitions and theorems in this text, we automated the translation from Coq definitions and theorems to numbered definitions and theorems in this text. The translation is performed by a straightforward Python [Mar06] script.

All numbered (and some unnumbered) definitions and theorems in this text correspond with definitions and theorems written in Gallina, the specification language of Coq. The definition of the ModRed algorithm is entirely included with equation 2 to equation 13. Similarly, the entire definitions of MultiRed and MultiRed' are included through equations.

In the correspondences, the set \mathbb{Z} is identified with the Coq set `Z`, while sets \mathbb{Z}_n are identified with Coq sets `Z_n`. To illustrate the correspondences, let us recall the ModRed correctness theorem 3.1, expressed by equation 15:

$$\forall M \in \{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\} \forall x \in \mathbb{Z}_{2^{\lceil \lg M \rceil + w}} : \\ \text{ModRed}(M, C(M), \lfloor \frac{x}{2^w} \rfloor, x \bmod 2^w) = x \bmod M$$

In Coq, this theorem looks like

```
Theorem ModRed_eq :
  forall (M : Mset)(x : Z_ (2 ^ (Zlog_sup (Z_from_Mset M) + w))),
    _Z (ModRed
      (Z_from_Mset M)
      (C M)
      (ex2w (_Z x / 2 ^ w) (xhex M x))
      (ex2w (_Z x mod 2 ^ w) (xlex (_Z x)))) =
      _Z x mod Z_from_Mset M.
```

Each expression `ex2w z p` is shorthand for `exist (in_Z_ (2 ^ w)) z p`. Such an expression represents a value $z \in \mathbb{Z}$ and a proof p of the fact that z is an element of \mathbb{Z}_{2^w} . The Coq expression `_Z` represents a “type conversion” injection which remains implicit in equation 15 because most readers need not be reminded that \mathbb{Z}_{2^w} is a subset of \mathbb{Z} . Similarly, `Z_from_Mset` represents an injection from $\{M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}\}$ to \mathbb{Z} .

Structure of the Proofs

The structure of the proofs closely corresponds to the structure of the algorithms. The goal of the first proof is to show that after the last step of the ModRed algorithm, the value of the result variable r''' is equal to $x \bmod M$, where $x = 2^w v + u$. For each step the proof introduces a lemma that captures an essential property of the variable defined at that step. For instance, for equation 2 in the algorithm:

$$s = \lfloor \frac{u}{2^p} \rfloor$$

we create the following lemma (which will be equation 53 in the proof):

$$s = \lfloor \frac{x \bmod 2^w}{2^p} \rfloor$$

Using such lemmas the ModRed correctness theorem is proved below in a bottom-up fashion.

6.3 Proof of the algorithm ModRed, theorem 3.1

In the rest of this section, the following assumptions hold for $M, x_1, x_0 \in \mathbb{Z}_{2^w}$ and $p \in \mathbb{Z}_w$: $1 \leq M$, $M \leq 2^{w-1}$, $p = \lceil \lg M \rceil$, $x_1 < 2^p$, and $x = 2^w x_1 + x_0$. The variables s to r''' are defined with equation 2 to equation 13.

At the beginning of the proof, the following identities are established.

$$s = \lfloor \frac{x \bmod 2^w}{2^p} \rfloor \quad (53)$$

$$s' = x \bmod 2^w - x \bmod 2^p \quad (54)$$

$$h = \lfloor \frac{x}{2^w} \rfloor 2^{w-p} \quad (55)$$

$$h' = \lfloor \frac{x}{2^p} \rfloor \quad (56)$$

$$q = \lfloor \frac{\lfloor \frac{x}{2^p} \rfloor (\lfloor \frac{2^{p+w}}{M} \rfloor \bmod 2^w)}{2^w} \rfloor \quad (57)$$

Central to the proof is to show that q' has only two possible values:

$$q' = \lfloor \frac{\lfloor \frac{x}{2^p} \rfloor 2^p}{M} \rfloor \bmod 2^w \vee q' = (\lfloor \frac{\lfloor \frac{x}{2^p} \rfloor 2^p}{M} \rfloor - 1) \bmod 2^w \quad (58)$$

Equation 58 can be proved easily using equation 1. Equation 53 to equation 58 lead to the following equations in turn:

$$y = \lfloor \frac{\lfloor \frac{x}{2^p} \rfloor 2^p}{M} \rfloor M \bmod 2^w \vee y = (\lfloor \frac{\lfloor \frac{x}{2^p} \rfloor 2^p}{M} \rfloor M - M) \bmod 2^w \quad (59)$$

$$d = \lfloor \frac{x}{2^p} \rfloor 2^p \bmod M \vee d = \lfloor \frac{x}{2^p} \rfloor 2^p \bmod M + M \quad (60)$$

$$r = (d + x \bmod 2^p) \bmod 2^w \quad (61)$$

$$r' = \lfloor \frac{x}{2^p} \rfloor 2^p \bmod M + x \bmod 2^p \quad (62)$$

$$r'' = x \bmod M \vee r'' = x \bmod M + M \quad (63)$$

$$r''' = x \bmod M \quad (64)$$

Equation 64 immediately leads to the ModRed correctness theorem, expressed by equation 15.

The correctness of the core of the MultiRed algorithm (eventually leading to theorem 4.1) is largely proved in the same way as the correctness of the ModRed algorithm. The MultiRed algorithm applies this core within a loop. The correctness of the loop is proved with natural induction using an induction step expressed as follows, where $\text{fst}(a, b) = a$, $\text{snd}(a, b) = b$, $x_i = \lfloor \frac{x}{2^{wi}} \rfloor \bmod 2^w$, and $f'((d_{i+1}, r_{i+1}), x_i) = (d_i, r_i)$, the latter definition assuming that M and the components of $C'(M)$ are given as inputs.

$$\begin{aligned} \forall d, r \in \mathbb{Z}_{2^w} \forall n \in \mathbb{N} : P(M, p, d, r, \lfloor \frac{x}{2^{(n+2) \cdot w}} \rfloor, x_{n+1}) \implies \\ P(M, p, \text{fst}(f'((d, r), x_n)), \text{snd}(f'((d, r), x_n)), \lfloor \frac{x}{2^{(n+1) \cdot w}} \rfloor, x_n) \end{aligned}$$

The definition of this induction step is similar to a loop invariant. It uses a predicate P expressing the required relations between its arguments such that at the end of the loop the required property holds. At the end of the loop, values of d and r have been obtained such that $P(M, p, d, r,$

$\lfloor \frac{x}{2^w} \rfloor, x_0$) holds. Because of this property of d and r , $x \bmod M$ is easy to derive from d and r . The predicate P is defined as follows.

$$\begin{aligned}
P(M, p, d, r, h, l) = & \left(\begin{aligned}
d &= \lfloor \frac{(h \bmod M)2^w + l}{2^p} \rfloor 2^p \bmod M \\
\vee \quad d &= \lfloor \frac{(h \bmod M)2^w + l}{2^p} \rfloor 2^p \bmod M + M \\
\vee \quad d &= \lfloor \frac{(h \bmod M + M)2^w + l}{2^p} \rfloor 2^p \bmod M \\
\vee \quad d &= \lfloor \frac{(h \bmod M + M)2^w + l}{2^p} \rfloor 2^p \bmod M + M \\
& \left. \begin{aligned}
& \right) \\
\wedge \quad r &= (d + l \bmod 2^p) \bmod 2^w
\end{aligned} \right)
\end{aligned}$$

The induction step infers a property of n from a property of $n + 1$. Proofs by induction usually employ a step where a property of $n + 1$ is inferred from a property of n , but here the order of n and $n + 1$ is reversed, starting with the highest nonzero word of x and iterating through the words of x until its lowest word is reached.

With the MultiRed' correctness proof (theorem 5.1), the slightly different predicate P' defined as follows is employed in an induction step.

$$\begin{aligned}
P'(M, p, d, r, h, l) = & \left(\begin{aligned}
d &= \lfloor \frac{(h \bmod M)2^w + l}{2^p} \rfloor 2^p \bmod M \\
\vee \quad d &= \lfloor \frac{(h \bmod M)2^w + l}{2^p} \rfloor 2^p \bmod M + M \\
\vee \quad d &= \lfloor \frac{(h \bmod M + M)2^w + l}{2^p} \rfloor 2^p \bmod M \\
\vee \quad d &= \lfloor \frac{(h \bmod M + M)2^w + l}{2^p} \rfloor 2^p \bmod M + M \\
& \left. \begin{aligned}
& \right) \\
\wedge \quad r &= (d + l \bmod 2^p \bmod M) \bmod 2^w
\end{aligned} \right)
\end{aligned}$$

6.4 Evaluation of Proof Methodology

We estimate that the effort of producing the proofs was about 6 weeks full-time work. In our proof scripts, 272 theorems and lemmas are found at the time of this writing. The proof scripts are comprised of 2801 lines, amounting to 142264 characters. The proofs employ large distributed Coq libraries. We made numerous revisions and variations of our proofs, and in the course of time, due to an increased experience in proving, the size of the proofs, and the time to construct them, shrank considerably. The lemmas and theorems were first proved in mind, then on paper, and finally, using Coq, on a computer. The proofs in mind made implicit use of many properties of integers. The proofs on paper required these properties to be explicated.

During the proof process no errors in the algorithms were found. The Coq proofs provide the best possible trust one can have in the correctness of these algorithms.

Many definitions and small lemmas that were needed for the proofs can also be used in other proofs concerning computer arithmetic. To facilitate that, we placed the proofs and the more general definitions and lemmas they rely on in separate files on the Coq site [Rut]. This may serve as the first step in creating a large computer arithmetic library for formal proofs.

7 Related work

With Barrett's division-free modular reduction algorithm [Bar87], $x \bmod M$ is computed by first approximating $\lfloor \frac{x}{M} \rfloor$ with $q = \lfloor \lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{M} \rfloor / b \rfloor$, for $a = d^{\lceil \log_d M \rceil}$ and $b = d^{\lceil \log_d M \rceil + 2}$, where $d \in \mathbb{Z} \mid d \geq 4$. As can be seen

in e.g. [BGV94], and in equations 1 and 65 of this text, it is also possible to choose $d = 2$. The value of $q = \lfloor \lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{M} \rfloor / b \rfloor$ can be $\lfloor \frac{x}{M} \rfloor$, $\lfloor \frac{x}{M} \rfloor - 1$, or $\lfloor \frac{x}{M} \rfloor - 2$, as

$$\forall y, b \in \mathbb{Z}^+ \quad \forall a \in \mathbb{Z}_y \quad \forall x \in \mathbb{Z}_{ab+a} : \lfloor \frac{x}{y} \rfloor - 2 \leq \lfloor \frac{\lfloor \frac{x}{a} \rfloor \lfloor \frac{ab}{y} \rfloor}{b} \rfloor \leq \lfloor \frac{x}{y} \rfloor \quad (65)$$

From the 3 possible values of q , it immediately follows that $x - qM$ can be equal to $x \bmod M$, $x \bmod M + M$, or $x \bmod M + 2M$. If $M > \lfloor \frac{2^w}{3} \rfloor$, the value of $x - qM$ may be greater than or equal to 2^w , which means there may be a nonzero word overflow $\lfloor \frac{x - qM}{2^w} \rfloor$. In that case, $x \bmod M$ cannot be derived from $(x - qM) \bmod 2^w$ just using at most two comparisons and two subtractions. In contrast, with equation 1 a similar overflow can only occur when $M > \lfloor \frac{2^w}{2} \rfloor$.

While Barrett's algorithm with $d = 2$ is based on equation 65, just like ModRed, it uses $a = 2^{\lfloor \lg M \rfloor}$ and $b = 2^{\lfloor \lg M \rfloor + 2}$ instead of $a = 2^{\lfloor \lg M \rfloor}$ and $b = 2^w$. Barrett's algorithm is usually employed in multi-word integer arithmetic, which is why it does not need to take special measurements to avoid word overflows. Only the upper half of the product $\lfloor \frac{x}{2^{\lfloor \lg M \rfloor}} \rfloor \lfloor \frac{2^{2(\lfloor \lg M \rfloor + 1)}}{M} \rfloor$ needs to be calculated in principle though. In the context of small moduli, the `udiv_qrnd_preinv1` macro of (version 4.1.4 of) the GNU Multiple Precision Arithmetic Library [GMP] shows that the low word of $\lfloor \frac{x}{2^{\lfloor \lg M \rfloor}} \rfloor \lfloor \frac{2^{\lfloor \lg M \rfloor + w}}{M} \rfloor$ need not be computed at all for $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$. A similar thing can also be seen in our ModRed algorithm.

One of the algorithms from [GM94] is for division and reduction of integers in $\mathbb{Z}_{2^{2w}}$ by positive integers in \mathbb{Z}_{2^w} . When the divisor is restricted to the smaller domain $\{M \in \mathbb{Z} \mid 2^{w-1} \leq M < 2^w\}$, some operations of this algorithm need not be carried out. This leads to the `udiv_qrnd_preinv1` macro of [GMP]. The macro yields $\lfloor \frac{x}{M} \rfloor$ and $x \bmod M$ for $x \in \mathbb{Z}_{2^{2w}}$ and $M \in \mathbb{Z} \mid 2^{w-1} \leq M < 2^w$. The algorithm and macro employ $\lfloor \frac{2^{\lfloor \lg M \rfloor + 1 + w}}{M} \rfloor \bmod 2^w$ as a predefined value, in contrast with ModRed, which uses $\lfloor \frac{2^{\lfloor \lg M \rfloor + w}}{M} \rfloor - 2^w$, and in contrast with Barrett's algorithm, which employs $\lfloor \frac{2^{2(\lfloor \lg M \rfloor + 1)}}{M} \rfloor$.

The algorithms passed in review above are for reduction of certain double words, including products of two integers $m, n \in \mathbb{Z}_M$. That is because $\forall x, y \in \mathbb{Z}_M : 0 \leq xy < M^2 \leq 2^{\lfloor \lg M \rfloor} 2^{w-1} < 2^{\lfloor \lg M \rfloor + w}$. Schrage's algorithm [BFS87] is designed to perform modular multiplications $mn \bmod M$, having inputs M and $m, n \in \mathbb{Z}_M$. It contains divisions rather than high-word multiplications. When high-word and double-word multiplication instructions are not available (and thus have to be emulated in software), Schrage's algorithm may evaluate modular multiplications more efficiently than with the previously mentioned modular reduction algorithms.

Some algorithms perform modular multiplications with very special moduli, for example with $M = 2^{31} - 1$ [PRB69]. Such algorithms are much less generally applicable than general modular reduction algorithms but their implementations can be much more efficient.

8 Future Work

Further transformations may be of interest for investigation in the future. It may be possible to replace comparisons with M by comparisons with

$2^{\lceil \lg M \rceil}$ in more places. This may require additional corrections at other places in order to maintain correctness. These corrections might influence performance negatively. Different variants will have to be considered, correctness will have to be proven, and performance measurements will have to be done.

Recently, ideas for a new reduction method have been communicated by Peter Montgomery on the GMP site ([GMP]). For many cases these ideas seem to incorporate an improvement. It is less clear yet whether they will also be an improvement for reducing relatively small large integers with moduli that are different for each call. As future work it seems worthwhile to define and implement an algorithm using these ideas as a starting point, to prove its correctness with a proof assistant, and to compare its performance with the performance of the algorithms proposed in this paper.

9 Conclusions

We have proposed an algorithm – ModRed – for reduction of integers in $\mathbb{Z}_{2^{\lceil \lg M \rceil + w}}$ by a modulus M on w -bit processors, where $w \in \mathbb{Z}^+$ and $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$. We have also proposed algorithms MultiRed and MultiRed' – based on ModRed – for reduction of multi-word integers by moduli $M \in \mathbb{Z} \mid 1 \leq M \leq 2^{w-1}$.

With measurements on processors which provide relatively slow divisions of double words by single words, we have shown that implementations of MultiRed can sometimes be over 30% more efficient than comparable implementations based on the algorithm proposed by Granlund and Montgomery for division and reduction of an unsigned double word by an unsigned word. That algorithm applies to larger moduli $M \in \mathbb{Z} \mid 2^{w-1} \leq M < 2^w$ as well, so with respect to the algorithm, ModRed trades some generality for some efficiency.

The formal correctness of the algorithms has been proved with the aid of the Coq proof assistant. This gives a very high degree of trust in the correctness of these algorithms that are expressed on the level of abstract machine instructions.

Acknowledgment

The authors would like to thank John van Groningen for insight and support relevant to this research.

References

- [AMD] Advanced Micro Devices, AMD - Processor Homepage. www.amd.com.
- [Bar87] P. Barrett. Implementing the Rivest Shamir and Adleman public key encryption algorithm on a standard digital signal processor. In A. Odlyzko, editor, *Advances in Cryptology, Proc. Crypto '86*, volume 263 of *Lecture Notes in Computer Science*, pages 311–323. Springer-Verlag, 1987.
- [BFS87] P. Bratley, B.L. Fox, and L.E. Schrage. *A Guide to Simulation*. Springer-Verlag New York Inc., second edition, 1987.

- [BGV94] A. Bosselaers, R. Govaerts, and J. Vandewalle. Comparison of three modular reduction functions. In *Advances in Cryptology, Proc. Crypto '93*, volume 773 of *Lecture Notes in Computer Science*, pages 175–186. Springer-Verlag, 1994.
- [Coq] The Coq proof assistant. pauillac.inria.fr/coq.
- [dB70] N. G. de Bruijn. The mathematical language automath, its usage, and some of its extensions. In *Symposium on Automatic Demonstration*, volume 125 of *Lecture Notes in Mathematics*, pages 29–61. INRIA, Versailles, 1968, Springer, Berlin, 1970.
- [GCC] GCC Home Page - GNU Project - Free Software Foundation (FSF). gcc.gnu.org.
- [GM94] T. Granlund and P.L. Montgomery. Division by Invariant Integers using Multiplication. *SIGPLAN Notices*, 29(6):61–72, 1994.
- [GMP] The GNU MP Bignum Library. www.swox.com/gmp.
- [IBM] IBM Power Architecture. www.ibm.com/chips/power.
- [Knu98] D.E. Knuth. *The Art of Computer Programming*, volume 2 – Seminumerical Algorithms. Addison Wesley Longman, third edition, 1998.
- [Mar06] Alex Martelli. *Python in a Nutshell*. O'Reilly, second edition, 2006.
- [PRB69] W.H. Payne, J.R. Rabung, and T.P. Bogyo. Coding the Lehmer Pseudo-random Number Generator. *Communications of the ACM*, 12(2):85–86, 1969.
- [Rut] L. Rutten. Fast Modular Reduction Proof Scripts in Coq. coq.inria.fr/contribs/modred.html.
- [Wie03] F. Wiedijk. Comparing Mathematical Provers. In G. Goos, J. Hartmanis, and J. van Leeuwen, editors, *Mathematical Knowledge Management: Second International Conference, Proceedings*, volume 2594 of *Lecture Notes in Computer Science*, pages 188–202. Springer, 2003.