

Sparse multivariate polynomial factorization: a high-performance design and implementation.

Michael Monagan¹ and Baris Tuncer¹

Department of Mathematics, Simon Fraser University, Vancouver, Canada
mmonagan@sfu.ca and tuncer@sfu.ca,

Abstract. Our goal is to develop a high-performance code for factoring a multivariate polynomial in n variables with integer coefficients which is polynomial time in the sparse case and efficient in the dense case. Maple, Magma, Macsyma, Singular and Mathematica all implement Wang's multivariate Hensel lifting, which, for sparse polynomials, can be exponential in n . Wang's algorithm is also highly sequential.

In this work we reorganize multivariate Hensel lifting to facilitate a high-performance parallel implementation. We identify multivariate polynomial evaluation and bivariate Hensel lifting as two core components. We have also developed a library of algorithms for polynomial arithmetic which allow us to assign each core an independent task with all the memory it needs in advance so that memory management is eliminated and all important operations operate on dense arrays of 64 bit integers. We have implemented our algorithm and library using Cilk C for the case of two monic factors. We discuss details of the implementation and present experimental timing results.

Keywords: Hensel Lifting, Polynomial Factorization, Cilk C

1 Introduction

Let $a = fg$ where f and g are two irreducible polynomials in $\mathbb{Z}[x_1, x_2, \dots, x_n]$. Let $\alpha := (\alpha_2, \alpha_3, \dots, \alpha_n) \in \mathbb{Z}^{n-1}$ be an evaluation point. For a given polynomial $h \in \mathbb{Z}[x_1, x_2, \dots, x_n]$ let us use the notation $h_j = h(x_1, \dots, x_j, \alpha_{j+1}, \dots, \alpha_n)$ so that $a_1 = a(x_1, \alpha_2, \dots, \alpha_n)$. To factor a we first factor the image a_1 over \mathbb{Z} . With high probability $f(x_1, \alpha)$ and $g(x_1, \alpha)$ will be irreducible so we obtain f_1 and g_1 . Next we use a process known as Multivariate Hensel Lifting (MHL) to recover f and g from a, f_1, g_1 . Maple, Magma, Macsyma, Singular and Mathematica all implement Wang's MHL from [7, 8]. A complete description of Wang's MHL may be found in Ch 6 of Geddes et. al. [2].

The input to Wang's MHL is a, α, f_1, g_1 and a lifting prime p . The evaluation point α and prime p must satisfy $\gcd(f_1, g_1) = 1$ in $\mathbb{Z}_p[x_1]$. The algorithm lifts the factors f_1, g_1 to f_2, g_2 then f_2, g_2 to f_3, g_3 etc. until we obtain f_n, g_n . At the j 'th step we have $a_j - f_j g_j \pmod{p} = 0$. At the end of this iteration we have $a - f_n g_n \pmod{p} = 0$. Thus for sufficiently large p we obtain the factorization $a = fg$ over \mathbb{Z} . The reason Hensel lifting is done modulo a prime p is to avoid an expression swell that would otherwise occur over \mathbb{Q} .

Algorithm 1 below shows the j 'th step of MHL. Throughout the paper we restrict our presentation to f and g monic in x_1 . We refer the reader to [2] for how to modify MHL for the non-monic case.

Algorithm 1 j^{th} step of Multivariate Hensel Lifting for $j > 1$: Monic Case.

Input : $p, \alpha_j \in \mathbb{Z}_p, a_j \in \mathbb{Z}_p[x_1, \dots, x_j], f_{j-1}, g_{j-1} \in \mathbb{Z}_p[x_1, \dots, x_{j-1}]$ where a_j, f_{j-1}, g_{j-1} are monic in x_1 and $a_j(x_j = \alpha_j) = f_{j-1}g_{j-1}$.

Output : $f_j, g_j \in \mathbb{Z}_p[x_1, \dots, x_j]$ such that $a_j = f_j g_j$ or FAIL.

```

1:  $f_j \leftarrow f_{j-1}; g_j \leftarrow g_{j-1}$ .
2:  $error \leftarrow a_j - f_j g_j$ .
3: for  $i$  from 1 while  $error \neq 0$  and  $\deg(f_j, x_j) + \deg(g_j, x_j) < \deg(a_j, x_j)$  do
4:    $c_i \leftarrow$  Taylor coefficient of  $(x_j - \alpha_j)^i$  of  $error$ 
5:   if  $c_i \neq 0$  then
6:     Solve the MDP  $\sigma_i g_{j-1} + \tau_i f_{j-1} = c_i$  in  $\mathbb{Z}_p[x_1, \dots, x_{j-1}]$  for  $\sigma_i$  and  $\tau_i$ .
7:      $(f_j, g_j) \leftarrow (f_j + \sigma_i \times (x_j - \alpha_j)^i, g_j + \tau_i \times (x_j - \alpha_j)^i)$ 
8:      $error \leftarrow a_j - f_j g_j$ 
9:   end if
10: end for
11: if  $error = 0$  then return  $f_j, g_j$  else return FAIL end if

```

There are two main computations in Algorithm 1, namely, the multivariate polynomial diophantine equation (MDP) in Step 6, which typically dominates the cost, and the multivariate multiplication of $f_j \times g_j$ in Steps 2 and 8. Wang's method for solving an MDP resembles his Hensel lifting. He first solves the univariate polynomial diophantine equation

$$\sigma_{i1} g_{j-1}(x_1, \alpha_2, \dots, \alpha_{j-1}) + \tau_{i1} f_{j-1}(x_1, \alpha_2, \dots, \alpha_{j-1}) = c_i(x_1, \alpha_2, \dots, \alpha_{j-1})$$

using the Euclidean algorithm then recovers x_2 then x_3 etc. in σ_i and τ_i . For each x_j there is an iteration on the degree of x_j similar to Algorithm 1. This results in a highly serial algorithm which precludes a parallel implementation.

Wang's solution to the MDP is exponential in $j - 1$ when the evaluation points α_i are non-zero. This makes the whole Hensel lifting process exponential for sparse f and g . Polynomial time algorithms were developed by Zippel in 1981 [9], Kaltofen in 1985 [5], and Monagan and Tuncer in 2016 [6].

Let us use the notation $\text{supp}(h)$ to denote the set of monomials appearing in the polynomial h . Monagan and Tuncer [6] solved this exponential problem by observing that if α_j in Algorithm 1 is chosen at random from a sufficiently large set then with high probability the monomials in σ_i for $i \geq 1$ will be contained in the monomials in f_{j-1} , that is $\text{supp}(\sigma_i) \subseteq \text{supp}(f_{j-1})$. Similarly, $\text{supp}(\tau_i) \subseteq \text{supp}(g_{j-1})$ with high probability. They interpolate σ_i and τ_i by picking $\beta_2, \dots, \beta_{j-1}$ at random from \mathbb{Z}_p , computing sufficiently many images of $\sigma_{ik} = \sigma_i(x_1, \beta_2^k, \beta_3^k, \dots, \beta_{j-1}^k)$ and $\tau_{ik} = \tau_i(x_1, \beta_2^k, \dots, \beta_{j-1}^k)$ for $1 \leq k$ by solving univariate diophantine equations

$$\sigma_{ik} g_{j-1}(x_1, \beta_2^k, \dots, \beta_{j-1}^k) + \tau_{ik} f_{j-1}(x_1, \beta_2^k, \dots, \beta_{j-1}^k) = c_i(x_1, \beta_2^k, \dots, \beta_{j-1}^k)$$

for σ_{ik} and τ_{ik} in $\mathbb{Z}_p[x_1]$. Equating coefficients we obtain linear systems. The linear systems are Vandermonde systems which can be solved efficiently in quadratic time and linear space – see Zippel [10]. This improves on Kaltofen’s solution to the MDP which results in large unstructured linear systems. The second author has installed this new approach in Maple. It will be available in Maple 2019.

2 High performance considerations

As a first step we reorganize the computation of c_i in Algorithm 1 to avoid recomputing the entire product $f_j \times g_j$. At the i ’th iteration of the loop we have $f_j = f_{j-1} + \sum_{k=1}^{i-1} \sigma_k (x_j - \alpha_j)^k$ and $g_j = g_{j-1} + \sum_{k=1}^{i-1} \tau_k (x_j - \alpha_j)^k$ and

$$c_i = \text{coeff}(a_j - f_j g_j, (x_j - \alpha_j)^i) = \frac{a_j^{(i)}(\alpha_j)}{i!} - \sum_{k=1}^{i-1} \sigma_k \tau_{i-k}$$

where $a^{(i)}$ is the i ’th derivative of a_j wrt x_j . So we may write the loop in Algorithm 1 as follows.

- 1: $f_j \leftarrow f_{j-1}$; $g_j \leftarrow g_{j-1}$; $ap \leftarrow a_j$; $df \leftarrow 0$; $dg \leftarrow 0$.
- 2: **for** i from 1 while $df + dg < \text{deg}(a_j, x_j)$ **do**
- 3: $ap \leftarrow i^{-1} \partial ap / \partial x_j$
- 4: $c_i \leftarrow ap(\alpha_j) - \sum_{k=1}^{i-1} \sigma_k \tau_{i-k}$
- 5: Solve the MDP $\sigma_i g_{j-1} + \tau_i f_{j-1} = c_i$ in $\mathbb{Z}_p[x_1, \dots, x_{j-1}]$ for σ_i and τ_i .
- 6: **if** $\sigma_i \neq 0$ set $df \leftarrow i$. **if** $\tau_i \neq 0$ set $dg \leftarrow i$.
- 7: **end for**
- 8: $f_j \leftarrow f_{j-1} + \sum_{k=1}^{i-1} \sigma_k \times (x_j - \alpha_j)^k$
- 9: $g_j \leftarrow g_{j-1} + \sum_{k=1}^{i-1} \tau_k \times (x_j - \alpha_j)^k$

How can we parallelize this for a multi-core computer? We are using Cilk C, a parallel extension of C available with the gcc compiler. See [1, 3]. Because of the time needed to start a Cilk process, the units of work should be of size at least 10^4 clock cycles, equivalently, at least 10^3 multiplications in \mathbb{Z}_p . Also, small units of work must require no memory allocations, otherwise memory management will become a parallel bottleneck. We propose to reduce the multivariate Hensel lifting in $\mathbb{Z}_p[x_1, \dots, x_j]$ to Hensel lift bivariate images in $\mathbb{Z}_p[x_1, x_j]$. That is we will Hensel lift x_j in

$$a_j(x_1, \beta_2^k, \dots, \beta_{j-1}^k, x_j), f_{j-1}(x_1, \beta_2^k, \dots, \beta_{j-1}^k), \text{ and } g_{j-1}(x_1, \beta_2^k, \dots, \beta_{j-1}^k).$$

Algorithm 2 is our main unit of work. The complexity estimates on the right count arithmetic operations in \mathbb{Z}_p . Here $d_1 = \text{deg}(a, x_1)$, $d_j = \text{deg}(a, x_j)$.

In Algorithm 2 the loop runs to either $df = \text{deg}(f_j, x_j)$ or $dg = \text{deg}(g_j, x_j)$, whichever is greater. Now since $d_f + d_g = d_j$ the most expensive step is the sum of products $\Sigma = \sum_{k=1}^{i-1} \sigma_i(x_1) \tau_{k-i}(x_1)$ in Step 5 which costs $\sum_{i=1}^{d_j} O(i d_1^2) \in O(d_j^2 d_1^2)$ in total. To reduce this we evaluate and interpolate x_1 in Σ . We first

Algorithm 2 HenselLift1: Bivariate Hensel Lift of x_j for $j > 1$.

Input: $p, \alpha_j \in \mathbb{Z}_p, a \in \mathbb{Z}_p[x_1, x_j], f_0, g_0 \in \mathbb{Z}_p[x_1]$ where a, f_0, g_0 are monic in x_1 , $a(x_1, \alpha_j) = f_0 g_0$ and $\gcd(f_0, g_0) = 1$.

Output : $f_j, g_j \in \mathbb{Z}_p[x_1, x_j]$ such that $a_j = f_j g_j$.

```

1:  $df \leftarrow 0; dg \leftarrow 0; da \leftarrow \deg(a, x_j)$ 
2: Solve  $sg_0 + tf_0 = 1$  for  $s, t \in \mathbb{Z}_p[x_1]$  using the EEA .....  $O(d_1^2)$ 
3: for  $i$  from 1 while  $df + dg < da$  do
4:    $a \leftarrow i^{-1} \partial a / \partial x_j$  .....  $O(d_1 d_j)$ 
5:    $c_i \leftarrow a(x_1, x_j = \alpha_j) - \sum_{k=1}^{i-1} \sigma_k(x_1) \tau_{i-k}(x_1)$  .....  $O(d_1 d_j) + O(id_1^2)$ 
6:   Solve  $\sigma_i g_0 + \tau_i f_0 = c_i$  for  $\sigma_i, \tau_i \in \mathbb{Z}_p[x_1]$  as follows .....  $O(d_1^2)$ 
7:      $\sigma_i \leftarrow (c_i s) \text{ rem } f_0; \tau_i \leftarrow (c_i - \sigma_i g_0) \text{ quo } f_0$ 
8:   if  $\sigma_i \neq 0$  then  $df \leftarrow i$  end if
9:   if  $\tau_i \neq 0$  then  $dg \leftarrow i$  end if
10: end for
11: We have  $f_j = f_0 + \sum_{i=1}^{df} \sigma_i(x_1)(x_j - \alpha_j)^i$  and  $g_j = g_0 + \sum_{i=1}^{dg} \tau_i(x_1)(x_j - \alpha_j)^i$ .
12: return  $[f_0, \sigma_1, \dots, \sigma_{df}]$  and  $[g_0, \tau_1, \dots, \tau_{dg}]$ 

```

tried Horner evaluation at $x_1 = 0, 1, 2, \dots, d_j$ and Newton interpolation which are both $O(d_j^2)$ but this didn't really help until $d_j > 100$ so we instead tried an FFT which worked better. We give details here for clarity. Let ω be a primitive n 'th root of unity in \mathbb{Z}_p with $2^{k-1} \leq d_j < n \leq 2^k$. Replace Step 5 with

```

5 Evaluate  $\sigma_{il} \leftarrow \sigma_{i-1}(\omega^l)$  and  $\tau_{il} \leftarrow \tau_{i-1}(\omega^l)$  for  $0 \leq l < n$  using an FFT.
for  $l = 0$  to  $n - 1$  do  $c_{il} \leftarrow \sum_{k=1}^{i-1} \sigma_{kl} \times \tau_{(i-k)l}$  end for .....  $O(i)$ 
Interpolate  $\Sigma(x_1)$  from values  $\{(\omega^l, c_{il}) : 0 \leq l < n\}$  using an inverse FFT.
 $c_i(x_1) \leftarrow a(x_1, x_j = \alpha_j) - \Sigma(x_1)$ .

```

This reduces the cost of Step 5 to $O(d_1 \log d_1) + O(id_1) + O(d_1 \log d_1) + O(d_j d_1)$. The total cost of Algorithm 2 is reduced to $O(d_1^2 d_j + d_1 d_j^2)$.

2.1 Implementation of HenselLift1

In Algorithm 2 there are univariate operations in $\mathbb{Z}_p[x_1]$ in steps 2, 5 and 7, bivariate operations in $\mathbb{Z}_p[x_1, x_j]$ in Steps 4 and 5 and FFTs in Step 5. For a high performance implementation we have designed a library of polynomial arithmetic for $\mathbb{Z}_p[x_1]$ and $\mathbb{Z}_p[x_1, \dots, x_n]$. The data structure for $\mathbb{Z}_p[x_1]$ is just a dense array of coefficients. For $\mathbb{Z}_p[x_1, \dots, x_n]$ we use a sparse representation. We encode, e.g., the trivariate polynomial $\sum_{i=1}^t a_i M_i(x_1, x_2, x_3)$ as a two arrays of integers $A = [a_1 | a_2 | \dots | a_t]$ and the monomials $X = [M_1 | M_2 | \dots | M_t]$ also stored as an array of integers, that is, each monomial $x_1^i x_2^j x_3^k$ in X is stored as the 64 bit integer $2^{42}i + 2^{21}j + k$. Each subroutine in our library has inputs which are either integers or arrays of integers or arrays of arrays of integers. The arrays may be for inputs, outputs, and, if needed, temporary storage. For example, in Step 7 the multiplications $c_i s$ and $\sigma_i g_0$ are done using the C routine

```

# define LONG long long int // 64 bit signed C integer
int polmul64s( LONG *A, LONG *B, LONG *C, int da, int db, LONG p );

```

Here $da = \deg(a, x)$, $db = \deg(b, x)$, the coefficients of $a(x)$ and $b(x)$ are stored in the arrays A and B . The product $c(x) = a(x)b(x) \bmod p$ is computed in the array C which must be an array of size at least $da + db + 1$.

As a second example, in Step 4 we differentiate $a(x_1, x_j)$ with respect to x_j by calling the routine `poldiff64s(A, X, t, 2, 2, p)` below. Here $a(x_1, x_j)$ is input in the arrays (A, X) and the routine overwrites (A, X) with the derivative $\partial a / \partial x_j$.

```
int poldiff64s( LONG *A, LONG *X, int t, int n, int j, LONG p ) {
// diff(a,x[j]): a is stored as pair (A,X) with t terms in n variables
// compute result in (A,X) and return the number of terms
```

To implement Algorithm 2 we first coded it by allocating space for the polynomials in Algorithm 2, so including space for $\sigma_1, \dots, \sigma_{df}$ for example. Then we make all polynomials parameters of `HenselLift1` so that Algorithm 2 does not allocate any new memory. This is possible because all polynomials have bounded degree. The resulting code `HenselLift1` will be called on many inputs in parallel and the temporary space can be reused.

2.2 Reduction from multivariate to bivariate Hensel lifting.

Algorithm 3 describes how we reduce Hensel lifting of x_j in f_{j-1}, g_{j-1} to many bivariate Hensel lifts of x_j . When we implemented Algorithm 3 we tested it on polynomials f and g with 100–8000 terms in $n = 6–15$ variables of degree 7. We observed that almost all the time was spent evaluating a_j at Y_k in step 8. In the next section we discuss how we implement evaluation and how we parallelized it. Here we point out that if instead of evaluating out x_2, \dots, x_{j-1} we evaluate out x_3, \dots, x_{j-1} , and thus interpolate the σ_i and τ_i from bivariate images in x_1, x_2 , then we we likely reduce the number of evaluations s thus leading to a speedup. We describe what we have implemented using a homomorphism diagram.

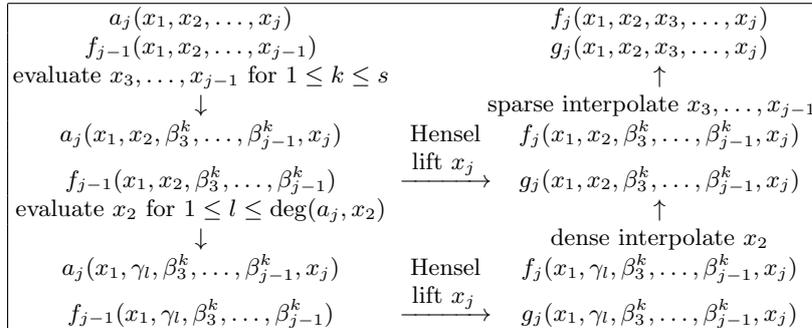


Fig. 1. Homomorphism diagram depicting our evaluation/interpolation strategy

In Figure 1 the reader will see two Hensel lifting steps which represent two possible ways of computing $f_j(x_1, x_2, \beta^k, x_j)$ and $g_j(x_1, x_2, \beta^k, x_j)$. In the first

Algorithm 3 Hensel Lift x_j

Input: Prime p , $\alpha_j \in \mathbb{Z}_p$, Monic polynomials $a_j \in \mathbb{Z}_p[x_1, \dots, x_j]$ $f_{j-1}, g_{j-1} \in \mathbb{Z}_p[x_1, \dots, x_{j-1}]$ with $j > 2$, s.t. $a_j(x_1, \dots, x_{j-1}, \alpha_j) = f_{j-1}g_{j-1}$.

- 1: Let $f_j = \sum_{i=0}^{df} \sigma_i(x_2, \dots, x_{j-1})x_1^i$ where $\sigma_i = \sum_{k=1}^{s_i} a_{ik}M_{ik}$ where $x_1^i M_{ik}$ are the monomials in $\text{supp}(f_{j-1})$ and $df = \deg(f_{j-1}, x_1)$.
- 2: Let $g_j = \sum_{i=0}^{dg} \tau_i(x_2, \dots, x_{j-1})x_1^i$ where $\tau_i = \sum_{k=1}^{t_i} b_{ik}N_{ik}$ where $x_1^i N_{ik}$ are the monomials in $\text{supp}(g_{j-1})$ and $dg = \deg(g_{j-1}, x_1)$.
- 3: Set $s = \max(s_i, t_i)$.
- 4: Pick $(\beta_2, \dots, \beta_{j-1}) \in \mathbb{Z}_p$ at random.
- 5: Compute monomial evaluation sets
 - $\{S_i = \{m_{ik} = M_{ik}(\beta_2, \dots, \beta_{j-1}) : 1 \leq k \leq s_i\} : 0 \leq i \leq df\}$ and
 - $\{T_i = \{n_{ik} = N_{ik}(\beta_2, \dots, \beta_{j-1}) : 1 \leq k \leq t_i\} : 0 \leq i \leq dg\}$.
 - If any $|S_i| \neq s_i$ or any $|T_i| \neq t_i$ try a different choice for $(\beta_2, \dots, \beta_{j-1})$.
 - If this fails **return** FAIL(1). (*p is not big enough*)
- 6: **for** k from 1 to s **in parallel do** (*Compute univariate images of σ_i and τ_i*)
- 7: Let $Y_k = (x_2 = \beta_2^k, \dots, x_{j-1} = \beta_{j-1}^k)$.
- 8: Evaluate: $a_k, f_0, g_0 \leftarrow a_j(x_1, Y_k, x_j), f_{j-1}(x_1, Y_k), g_{j-1}(x_1, Y_k)$.
- 9: **if** $\gcd(f_0, g_0) \neq 1$ **return** FAIL(2) (*an unlucky evaluation*)
- 10: Call HenselLift1($p, \alpha_j, a_k, f_0, g_0$) to compute $\sigma_{ik}(x_1)$ and $\tau_{ik}(x_1)$ such that $a_k - f_k g_k = 0$ where $f_k = \sum_{i=0}^{df} \sigma_{ik}(x_j - \alpha_j)^i$ and $g_k = \sum_{i=0}^{dg} \tau_{ik}(x_j - \alpha_j)^i$.
- 11: **end for**
- 12: **for** $i = 0$ to df **do**
- 13: Construct and solve the $s_i \times s_i$ linear system

$$\left\{ \sum_{k=1}^{s_i} a_{ik} m_{ik}^n = \text{coefficient of } x_1^i \text{ in } \sigma_{in}(x_1) \text{ for } 1 \leq n \leq s_i \right\}$$

for the coefficients a_{ik} of $\sigma_i(x_2, \dots, x_{j-1})$. Because it is a Vandermonde system in m_{ik} which are distinct by Step 5 it has a unique solution.

- 14: **end for**
- 15: Do the same for the $t_i \times t_i$ linear systems to solve for the coefficients b_{ik} of the τ_i .
- 16: Substitute the solutions for a_{ik} into f_j and b_{ik} into g_j and return(f_j, g_j).

way (the top Hensel lift) the diophantine equations $\sigma_i g_0 + \tau_i f_0 = c_i$ in Step 6 of Algorithm 2 are in $\mathbb{Z}_p[x_1, x_2]$ thus bivariate. One can solve these using dense evaluation and interpolation of x_2 in $O(d_1^2 d_2 + d_1 d_2^2)$ arithmetic operations in \mathbb{Z}_p where $d_2 = \deg(a_j, x_2)$. See Monagan and Tuncer [6].

We coded this approach in Maple as an experiment and found that the most expensive computation is the sum of products $\Sigma = \sum_{k=1}^{i-1} \sigma_k(x_1, x_2) \tau_{i-k}(x_1, x_2)$ in Step 5 of Algorithm 2 which are now bivariate multiplications which cost $O(id_1^2 d_2^2)$. To reduce this cost, we experimented with evaluating and interpolating x_2 which is described by the bottom Hensel lift in Figure 1. So the number of univariate images $f_j(x_1, \gamma_l, \beta^k, x_j), g_j(x_1, \gamma_l, \beta^k, x_j)$ needed to interpolate x_2 is $\max(\deg(f_j, x_2), \deg(g_j, x_2)) < \deg(a_j, x_2) = d_2$. In our current implementation we have parallelized the computation of the Hensel lifts of these images.

2.3 Parallelizing Evaluation

We describe how we parallelize the evaluations in Step 8 of Algorithm 3. Let $a_j = (A, X)$ where the monomials in X are sorted in lexicographical order with $x_1 > x_2 > \dots > x_j$. We first sort the monomials into $x_1 > x_j > x_2 > \dots > x_{j-1}$. Now when we evaluate $a_j(x_1, x_j, \beta_2^k, \dots, \beta_{j-1}^k)$ the remaining monomials will be sorted on $x_1 > x_j$. Let $a_j = \sum_{i=1}^t a_i x_1^{d_i} x_j^{e_i} M_i(x_2, \dots, x_{j-1})$. Let $A = [a_1, a_2, \dots, a_t]$ be the array of coefficients, $m_i = M_i(\beta_2, \dots, \beta_{j-1})$ and $B = [m_1, m_2, \dots, m_t]$ be the array of monomial evaluations and let Y be the array of monomials $[x_1^{d_1} x_j^{e_1}, \dots, x_1^{d_t} x_j^{e_t}]$. If we initialize $C_0 := A = [a_1, \dots, a_t]$ and define $C_k = [a_1 m_1^k, \dots, a_t m_t^k]$ then we have

$$a_j(x_1, x_j, \beta_2^k, \dots, \beta_{j-1}^k) = \sum_{i=1}^t a_i m_i^k x_1^{d_i} x_j^{e_i} = \sum_{i=1}^t C_{k,i} Y_i$$

and we can compute C_{k+1} using t multiplications with

$$C_{k+1} \leftarrow [B_1 \times C_{k1}, \dots, B_t \times C_{kt}] = [a_1 m_1^{k+1}, \dots, a_t m_t^{k+1}]$$

Then we assemble the result from $\sum_{i=1}^t C_{k+1,i} Y_i$ which requires adding coefficients of equal monomials in x_1, x_j . Since the monomials in Y_i are already sorted on $x_1 > x_j$ this is $O(t)$. Thus the total number of multiplications needed is st plus those needed to compute m_1, \dots, m_t .

Our first attempt to parallelize this for N cores was to do N evaluations at a time as done by Hu and Monagan in [4]. First compute C_1, C_2, \dots, C_N and the array $\Gamma = [m_1^N, m_2^N, \dots, m_t^N]$. To obtain the next N evaluations in parallel, on the k 'th core compute $C_{k+N} \leftarrow [C_{k1} \times \Gamma_1, \dots, C_{kt} \times \Gamma_t] = [a_1 m_1^{k+N}, \dots, a_t m_t^{k+N}]$. One problem with this approach is that we require $\#a$ words of memory for each C_1, \dots, C_N . For one of our benchmark problems where $\#a = 64,000,000$ this is about a half a gigabyte per core. Another problem is that we did not obtain full parallel speedup on our 16 core computer as the computation becomes memory bound when $\#a$ is this large. The following works.

Split $a_j = (A, X)$ into N blocks of size t/N terms. Each core evaluates a block of a_j at β^{k+1} which must be combined later. Numbering the cores $0, 1, \dots, N-1$ core c computes $C_{kl} \times B_l$ for $c(t/N) < l \leq c(t/N+1)$. We found that we obtained a 20% improvement by also computing the evaluation β^{k+2} at the same time so that we compute two evaluations at a time.

3 Experimental Results

We give two sets of experimental results. The first set (see Table 1) is for polynomials in many variables with relatively low degree. Here, evaluation of a_j is the bottleneck in our method – the time spent Hensel lifting images is negligible. The second set (see Table 2) is for polynomials with higher degree where Hensel lifting becomes the bottleneck. All experiments were performed on a server with

two Intel E5-2660 8 core CPUs running at 2.2GHz (base) and 3.0GHz (turbo) hence the maximum theoretical parallel speedup is a factor of $16.2/3.0 = 11.7$.

In Tables 1 and 2 the factors f and g are of the form $x_1^d + \sum_{i=2}^t a_i \prod_{j=1}^n x_j^{e_{ji}}$ where the coefficients a_i are chosen randomly from $[1, 999]$ and the exponents e_{ji} randomly from $[0, d - 1]$. The timings are for Hensel lifting x_n the last variable only, which is always most of the time. The quantity s in column 4 is the number of images needed to interpolate x_3, \dots, x_n in Figure 1. Table 1 shows we achieve very good parallel speedup for evaluations. We are still working on improving the parallel speedup for our Hensel lifting (Table 2).

For Maple we report two timings. The first is the best case of Wang's method where the evaluation points $\alpha_2, \dots, \alpha_n$ are all 0. To obtain this timing we forced Maple to use x_1 as the main variable (by default, it chooses a variable of least degree) and we added a constant to f and g as Maple requires that the leading and trailing coefficient in x_1 not vanish at α . The second timing is the worst case for Wang's method where all evaluation points are non-zero. It is the actual timing for Maple on these inputs.

n	d	t	s	New times (1 core)			New times (16 cores)			Maple 2018	
				total	(hensel)	(eval)	total	(hensel)	(eval)	best	worst
6	7	500	18	0.100	(0.016)	(0.042)	0.084	(0.023)	(0.010 - 4.2x)	0.411	28.84
6	7	1000	30	0.387	(0.026)	(0.221)	0.219	(0.033)	(0.026 - 8.5x)	1.140	58.46
6	7	2000	47	1.470	(0.041)	(1.102)	0.285	(0.036)	(0.082 - 13.4x)	3.066	99.88
6	7	4000	81	4.694	(0.074)	(3.661)	0.702	(0.060)	(0.293 - 12.5x)	7.173	162.49
6	7	8000	145	11.53	(0.130)	(9.765)	1.710	(0.107)	(0.793 - 12.3x)	15.61	
9	7	500	16	0.100	(0.014)	(0.035)	0.081	(0.016)	(0.007 - 5.0x)	1.171	7564.9
9	7	1000	29	0.490	(0.025)	(0.267)	0.219	(0.024)	(0.028 - 9.5x)	3.704	10010.4
9	7	2000	50	2.625	(0.044)	(1.777)	0.426	(0.038)	(0.148 - 12.0x)	13.43	
9	7	4000	93	16.71	(0.082)	(13.33)	2.125	(0.070)	(1.199 - 11.1x)	51.77	
9	7	8000	164	106.2	(0.145)	(92.45)	10.68	(0.121)	(7.247 - 12.8x)		

Table 1: Timings (real time in seconds) for increasing n and t . NA = not attempted.

n	d	t	s	New time (1 core)			New time (16 cores)			Maple 2018	
				total	(hensel)	(eval)	total	(hensel)	(eval)	best	worst
6	10	500	10	0.101	(0.031)	(0.025)	0.075	(0.020 - 1.5x)	(0.005)	0.571	92.49
6	15	500	6	0.139	(0.075)	(0.016)	0.094	(0.034 - 2.2x)	(0.004)	0.751	7956.5
6	20	500	5	0.249	(0.181)	(0.017)	0.115	(0.052 - 3.5x)	(0.005)	0.919	48610.1
6	40	500	3	1.292	(1.214)	(0.015)	0.335	(0.253 - 5.0x)	(0.011)	1.615	NA
6	60	500	3	4.932	(4.838)	(0.015)	1.161	(1.059 - 4.6x)	(0.011)	3.343	NA
6	80	500	3	14.81	(14.70)	(0.016)	3.436	(3.315 - 4.4x)	(0.012)	4.485	NA
6	10	2000	30	1.785	(0.097)	(1.068)	0.384	(0.048 - 2.0x)	(0.087)	5.237	976.94
6	15	2000	18	1.635	(0.238)	(0.707)	0.409	(0.080 - 3.0x)	(0.055)	7.166	23128.5
6	20	2000	12	1.670	(0.484)	(0.480)	0.456	(0.116 - 4.2x)	(0.040)	9.195	NA
6	40	2000	6	4.237	(3.217)	(0.260)	1.005	(0.607 - 5.3x)	(0.033)	15.98	NA
6	60	2000	6	15.26	(14.15)	(0.292)	3.030	(2.540 - 5.6x)	(0.053)	42.32	NA
6	80	2000	4	28.37	(27.27)	(0.217)	5.735	(5.181 - 5.3x)	(0.054)	57.33	NA

Table 2: Timings (real time in seconds) for increasing degree.

4 Implementation notes and Cilk C

We end with some comments about programming in Cilk C. Cilk has a very simple task model. One starts a new task using the `spawn` directive. Typically one creates several tasks in a `C` for loop inside a `C` function. One may wait for all the tasks started inside the function to complete using the `Cilk sync;` directive. And that's essentially it! We had few problems with Cilk. But

Coding in Cilk C basically means we are coding in C where we must manage the memory needed for every polynomial operation. Naively calling `malloc` and `free` in every subroutine will ruin parallel performance and degrade serial performance. Having to manage memory greatly increases coding effort. To reduce this difficulty we spent time designing many in-place algorithms, that is, algorithms requiring no additional memory.

It was very hard work getting an algorithm that took about two days to code in Maple to work in Cilk C. In C there is no array bounds checking. Incorrect memory references result in corrupted data which is difficult to track down. Maintaining an identical version of the code in Maple is helpful here. What we would find helpful is to code in C++ using the array data type, which does not support bounds checking, but have some tool for automatically converting arrays to C++ vectors where array bounds checking is available.

The data structure we use for multivariate polynomials assumes monomials can be packed into a 64 bit integer which limits the degree and number of variables that our software can handle. To accommodate more variables we plan to use the 128 bit integer type available in `gcc`, thus doubling the number of variables of a given degree that we can handle, and redesign our library to allow 32 bit, 64 bit and 128 bit monomials.

References

1. Matteo Frigo, Charles E. Leiserson and Keith H. Randall. The Implementation of the Cilk-5 Multithreaded Language. *Proc. PLDI 1988*, ACM, 212–223, 1998.
2. K.O. Geddes, S.R. Czapor, G. Labahn, *Algorithms for Computer Algebra*, Kluwer Acad. Publ. (1992).
3. The Cilk Project, <http://supertech.csail.mit.edu/cilk/>.
4. Jiaxiong Hu and Michael Monagan. A fast parallel sparse polynomial GCD algorithm. In *Proceedings of ISSAC 2016*, ACM Press, pp. 271–278, 2016.
5. Kaltofen, E., Sparse Hensel lifting. *Proc. EUROCAL '85*, Springer Verlag LNCS, vol 204, pp. 4–17, (1985).
6. Michael Monagan and Baris Tuncer, Using Sparse Interpolation in Hensel Lifting. *Proceedings of CASC 2016*, Springer-Verlag LNCS **9890**, pp. 381–400, 2016.
7. Wang, P.S., Rothschild, L.P. Factoring multivariate polynomials over the integers. *Mathematics of Computation*, vol 29, NUMBER 131, pp. 935–950, (1975).
8. Wang, P.S. An improved Multivariate Polynomial Factoring Algorithm, *Mathematics of Computation*, **32**, (1978).
9. Zippel, R.E. Newton's iteration and the sparse Hensel algorithm. *Proc. ACM Symp. Symbolic Algebraic Comp.*, 68–72, (1981).
10. Zippel, R.E. Interpolating polynomials from their values. *J. Symbolic Comput.*, **9**(3):375–403, (1990).