

What About Toom-Cook Matrices Optimality ?

Marco Bodrato and Alberto Zanoni

Centro “Vito Volterra” – Università di Roma Tor Vergata

Via Columbia 2 – 00133 Rome, Italy

{bodrato, zanoni}@volterra.uniroma2.it

**

October 19, 2006

Abstract. Karatsuba and Toom-Cook are well-known methods used to multiply efficiently two long integers. There have been different proposals about the interpolating values used to determine the matrix to be inverted and the sequence of operations to invert it. A definitive word about which is the optimal matrix (values) and the (number of) basic operations to invert it seems still not to have been said. In this paper we present some particular examples of useful matrices and a method to generate automatically, by means of optimised exhaustive searches on a graph, the best sequence of basic operations to invert them.

AMS Subject Classification: 11A05, 11A25, 11K65, 11Y70

Keywords and phrases: Polynomial multiplication, Karatsuba, Toom-Cook, interpolation

1 Introduction

Starting with the works of Karatsuba [4], Toom [6] and Cook [2], who found methods to lower asymptotic complexity for polynomial multiplication from $O(n^2)$ to $O(n^{1+\epsilon})$ with $0 < \epsilon \leq \log_2 3$, many efforts have been done in finding optimised implementations in arithmetic software packages [5], [8], [3].

The family of so-called Toom-Cook methods is an infinite set of algorithms (called Toom-3, Toom-4, etc. - Karatsuba may be identified with Toom-2). Each of them may be viewed as a polynomial interpolation problem, for which the base points are not specified a priori, from which a matrix to be inverted rises. We indicate the matrix related to Toom- n method with $A_n \in GL(\mathbb{Z}, 2n - 1)$.

Moreover, a set O of basic operations (typically sums, subtractions, bit shiftings, multiplication and division by small numbers, etc.) is given. Practically, this is a set of very efficiently implemented basic functions in a certain computer language, and the idea is to use them to invert A_n step by step.

A particular implementation of Toom- n method must then specify

1. The interpolation points v_i (Toom- n requires $2n - 1$ values), which determine the matrix A_n .
2. The sequence of operations in O needed to invert the matrix (we indicate it with *inversion sequence*, or IS, for short, when we want to emphasise the dimension, we will use IS_n).

In the scientific literature there is still not a definitive word on which is the best matrix to be used and which is the corresponding sequence of basic operations to invert it. For Karatsuba method, an exhaustive search on a “reasonable” set of interpolation points is very easy, while already for Toom-3 things are less clear. For example, only recently the GMP library (from version 4.2) changed implementation for Toom-3, choosing a different IS_3 , more efficient than the one in the precedent release.

It is not trivial at all to prove the optimality of a matrix A_n with respect to the related Toom- n method, because in principle there are infinite possibilities for the v_i values, and a combinatorial number of inversion sequences, and different optimality criteria. Some heuristics for v_i choice are present in Zuras’ paper [9], but, to the best of our knowledge, the final word has still not been said.

In this paper we introduce some optimality criteria to measure goodness of inversion sequences, and present an algorithm to automatically search for an optimal IS starting from a given matrix A_n . We considered a theoretical model following an approach based on:

- Minimal use of extra memory (no temporary variables)
- Only exact integer divisions
- No matrix support grow (a zero entry remains always zero)
- No multiplication of a line times an integer

** A copy of this paper can be downloaded from <http://bodrato.it/papers/>

2 A short review on Toom- n methods

Toom- n is a class of recursive methods to multiply two polynomials of degree m with complexity lower than $O(m^2)$. In particular, a very standard analysis tells that Toom- n method has asymptotic complexity $O(m^{\log_n(2n-1)})$, so that in principle it would be possible to approach linear complexity for $n \rightarrow \infty$.

In practise, only the methods with very small values of n (as 2, 3, 4) are used, because of the asymptotically better $O(n \log n \log \log n)$ -complexity Schönhage-Strassen method [7]. The thresholds indicating the convenience of one method in comparison with another one depend very much on the implementation, but also on the choice of A_n and of the performed IS_n .

All these methods can be used to multiply also integers, just fixing a number $1 < B \in \mathbb{N}$ (the *base*) and considering integers expansion in base B , so that to each integer u we can univocally associate a polynomial $u_B(x)$:

$$\mathbb{Z} \ni u = \sum_{k=0}^{\log_B i} u_k B^k \quad \Longrightarrow \quad u_B(x) = \sum_{k=0}^{\log_B i} u_k x^k \in \mathbb{Z}[x]$$

Integer multiplication reduces to polynomial multiplication, followed by an evaluation:

$$u \cdot v = z \quad \Longrightarrow \quad u_B(x) \cdot v_B(x) = z_B(x) \quad \Longrightarrow \quad z = z_B(x)|_{x=B}$$

2.1 The classical point of view

As indicated above, usually one makes the assumption that the two factors have the same degree. If it is not the case, padding the lower-degree polynomial with zero coefficients is supposed. To compute $c(x) = a(x) \cdot b(x) \in \mathbb{K}[x]$ such that

$$a(x) = \sum_{i=0}^d a_i x^i \quad ; \quad b(x) = \sum_{i=0}^d b_i x^i$$

we set up an interpolation problem. Of course, for every $v \in \mathbb{K}$, $c(v) = a(v) \cdot b(v)$. Provided \mathbb{K} is not too small, we may choose $2d + 1$ different values $v_i \in \mathbb{K}$, and compute the sequences $\mathbf{a} = (a(v_0), \dots, a(v_{2d}))$ and $\mathbf{b} = (b(v_0), \dots, b(v_{2d}))$, that let us to compute the product sequence

$$\mathbf{w} = (w(v_0), \dots, w(v_{2d})) = (a(v_0) \cdot b(v_0), \dots, a(v_{2d}) \cdot b(v_{2d}))$$

Recovering the set of coefficients $\mathbf{c} = (c_0, \dots, c_{2d})$ is now just a classical interpolation problem solving

$$A_n \mathbf{c} = \mathbf{w} \quad \Longrightarrow \quad \begin{pmatrix} v_0^{2d} & \dots & v_0 & 1 \\ v_1^{2d} & \dots & v_1 & 1 \\ \vdots & & \vdots & \\ v_{2d}^{2d} & \dots & v_{2d} & 1 \end{pmatrix} \begin{pmatrix} c_0 \\ c_2 \\ \vdots \\ c_{2d} \end{pmatrix} = \begin{pmatrix} w_0 \\ w_2 \\ \vdots \\ w_{2d} \end{pmatrix} \quad \Longrightarrow \quad \mathbf{c} = A_n^{-1} \mathbf{w}$$

For integer multiplication, suppose the factors have l digits if considered in a certain base B . Now consider them in the new base $B' = B^{l/n}$: what practically happens is to consider n parts of l/n B -digits each. The i^{th} part is nothing but the coefficient a_i (b_i).

Now perform the above indicated computations: the whole process requires $2d + 1 = 2n - 1$ multiplications of numbers with n times less digits¹, with the additional overhead of $a(v_i), b(v_i)$ values computation and the inversion of the resulting matrix A_n . Note that only matrices with odd order are obtained this way.

2.2 A temporary useful point of view

In this paper, for performance testing, we will also consider “intermediate” versions of Toom-Cook methods, which consider $a(x), b(x)$ with different degrees d_1, d_2 , without zero padding. Indicating with $n_1 = d_1 + 1, n_2 = d_2 + 1$ the number of necessary subdivisions, the principle remains basically the same, but now the obtained matrices – indicated with A_{n_1, n_2} – can have also even order.

Examples : (without loss of generality we suppose $n_1 > n_2$)

¹ Recursive application of the method to $a(v_i) \cdot b(v_i)$ products results in the asymptotic complexity estimate $O(m^{\log_n(2n-1)})$.

$$\boxed{n_2 = 1}$$

Here $b(x) = b_0$, and it is more convenient to compute directly the product coefficients $c_i = a_i \cdot b_0$.

$$\boxed{(n_1, n_2) = (3, 2)}$$

Here $c(x)$ has degree 3, and $A_{3,2}$ has order 4. In section 7 we will refer to this method as Toom-2.5

$$\boxed{(n_1, n_2) = (4, 3)}$$

Here $c(x)$ has degree 5, and $A_{4,3}$ has order 6. In section 7 we will refer to this method as Toom-3.5

$$\boxed{(n_1, n_2) = (4, 2)}$$

In this case $A_{4,2}$ has order 5, and \bar{A}_3 may well be chosen as an optimal choice $\bar{A}_{4,2}$

$$\boxed{\left. \begin{array}{l} n_1 - n_2 \\ d_1 - d_2 \end{array} \right\} \equiv 0 \pmod{2}}$$

Generalisation of the precedent cases: a “fall back” to the classical Toom- $(d_1 + d_2 - 1)$ method: $\bar{A}_{n_1+n_2-3}$ may well be chosen as an optimal choice for \bar{A}_{n_1, n_2}

While the (3,2) and (4,3) models have most a theoretical value, the (4,2) case suits very well in practical cases. Infact, a program/library implementing both Toom-3 and Toom-4 methods *should* use Toom-3 when one factor has 4 parts and the other one just 2.

3 The matrices A_n

Even if Toom- n works for whatever choice of the v_i , it is better to choose them in order to minimize the matrix inversion overhead as much as possible. The inverse is usually computed by a sequence of elementary row operations à la Gauss, and therefore we should search the “shortest” (least number of elementary operations) and easiest (fastest elementary operations) way to compute A_n^{-1} .

It is possible to consider rational v_i values even when working only with integers. Infact, if $v_i = N_i/D_i$, we have

$$a(v_i) \cdot b(v_i) = c(v_i) \implies \left(\sum_{k=0}^d a_k \left(\frac{N_i}{D_i} \right)^k \right) \cdot \left(\sum_{k=0}^d b_k \left(\frac{N_i}{D_i} \right)^k \right) = \sum_{k=0}^{2d} c_k \left(\frac{N_i}{D_i} \right)^k$$

and multiplying by D_i^{2d} we may get rid of the denominators, so that:

$$\left(\sum_{k=0}^d a_k N_i^k D_i^{2d-k} \right) \cdot \left(\sum_{k=0}^d b_k N_i^k D_i^{2d-k} \right) = \sum_{k=0}^{2d} c_k N_i^k D_i^{2d-k}$$

which means that the i^{th} line of A_n will be $(N_i^{2d}, N_i^{2d-1}D_i, \dots, N_i D_i^{2d-1}, D_i^{2d})$. In particular,

- integer values $v_i = N_i$ generate lines $(N_i^{2d}, N_i^{2d-1}, \dots, N_i, 1)$
- integer reciprocals $v_i = 1/D_i$ generate $(1, D_i, \dots, D_i^{2d-1}, D_i^{2d})$

We use ∞ as interpolation “value” to indicate $a_d \cdot b_d$ product computation, which in a certain sense represents interpolation on the reciprocal of zero, or, more precisely,

$$a(\infty) = \lim_{x \rightarrow \infty} \frac{a(x)}{x^d} = a_d \quad ; \quad b(\infty) = \lim_{x \rightarrow \infty} \frac{b(x)}{x^d} = b_d$$

Following the literature and for obvious efficiency, we will always consider the interpolating values $v_0 = 0$ and $v_{2d} = \infty$, so that A_n will have the shape

$$A_n = \begin{pmatrix} 1 & 0 & \dots & 0 & 0 \\ \text{X} & \text{X} & \dots & \text{X} & \text{X} \\ \vdots & \vdots & & \vdots & \\ \text{X} & \text{X} & \dots & \text{X} & \text{X} \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix} \quad \text{with X indicating not zero entries} \quad (1)$$

and $c_0 = w_0$, $c_{2d} = w_{2d}$. Our inversion analysis will then focus mainly on the inner lines.

The well-known result about Vandermonde determinants applied to Toom matrices gives the following:

Proposition 1. For a Vandermonde matrix A_n generated by the $r = 2n - 1$ values $\left\{ \infty, \frac{N_2}{D_2}, \dots, \frac{N_{r-1}}{D_{r-1}}, 0 \right\}$ (Toom matrix) we have

$$\det(A_n) = \left(\prod_{i=2}^{r-1} N_i D_i \right) \cdot \prod_{1 < i < j < r} (N_i D_j - N_j D_i)$$

Proof. Starting from the classical Vandermonde formula, and remembering that in order to have integer entries we cleared the denominators in every line of A_n multiplying it by D_i^{r-1} , we have

$$\begin{aligned} \det(A_n) &= \left(\prod_{i=1}^r D_i^{r-1} \right) \cdot \prod_{i < j} \left(\frac{N_i}{D_i} - \frac{N_j}{D_j} \right) = \left(\prod_{i=1}^r D_i^{r-1} \right) \cdot \prod_{i < j} \left(\frac{N_i D_j - N_j D_i}{D_i D_j} \right) = \\ &= \left(\prod_{i=1}^r D_i^{r-1} \right) \cdot \frac{\prod_{i < j} (N_i D_j - N_j D_i)}{\prod_{i < j} D_i D_j} \end{aligned}$$

In the denominator $\prod_{i < j} D_i D_j$ every term D_k appears $r - k$ times in the first place ($k = i, j = k + 1, \dots, r$) and $k - 1$ times in the second ($j = k, i = 1, \dots, k - 1$), so that the result is $\prod_{i=1}^r D_i^{r-1}$. Canceling out, we obtain

$$= \left(\prod_{i=1}^r D_i^{r-1} \right) \cdot \frac{\prod_{i < j} (N_i D_j - N_j D_i)}{\prod_{i=1}^r D_i^{r-1}} = \prod_{1 \leq i < j \leq r} (N_i D_j - N_j D_i)$$

The value 0 corresponds to the choice $N_r = 0, D_r = 1$, while ∞ to $N_1 = 1, D_1 = 0$ (it is easy to prove that the above formula is also valid in this latter case). Considering these two cases ($i = 1, j = r$) apart, we have

$$= \left(D_r \prod_{j=2}^{r-1} D_j \right) \left(\prod_{1 < i < j < r} (N_i D_j - N_j D_i) \right) \left(N_1 \prod_{i=2}^{r-1} N_i \right)$$

which gives the desired result.

Lemma 1. Let $\alpha = \pm 2^a \neq \beta = \pm 2^b$, with $a, b \in \mathbb{Z}$. The number $\gamma = \alpha - \beta$ is a power of 2 only if $\alpha = -\beta$ or $\alpha = 2\beta$ or $\alpha = \beta/2$.

Proof. There are two cases:

1. $\alpha\beta > 0$: we have $\pm\gamma = 2^a - 2^b = 2^b(2^{a-b} - 1)$, and $2^{a-b} - 1$ is a power of 2 iff $|a - b| \leq 1$. Excluding the case $a = b$, corresponding to $\alpha = \beta$, only the cases $\alpha = 2\beta$, $\alpha = \beta/2$ remain.
2. $\alpha\beta < 0$: we have $\pm\gamma = 2^a + 2^b = 2^b(2^{a-b} + 1)$, and $2^{a-b} + 1$ is a power of 2 iff $a = b$, that is $\alpha = -\beta$

Proposition 2. For each Toom matrix A_n with $n \geq 3$, its determinant is not a power of 2.

Proof. By contradiction, we should have that $\det(A_n) = \left(\prod_{i=2}^{r-1} N_i D_i \right) \prod_{1 < i < j < r} (N_i D_j - N_j D_i)$ should be a power of 2, according to proposition 1. Looking at the first factor, this means that all the N_i, D_i should be powers of 2, say $N_i = \pm 2^{e_i}, D_i = 2^{f_i}$ and the same for all the factors $(N_i D_j - N_j D_i)$. Supposing that $(N_i, D_i) = 1$ ($e_i \cdot f_i = 0$), we have

$$(N_i D_j - N_j D_i) = D_i D_j \frac{(N_i D_j - N_j D_i)}{D_i D_j} = D_i D_j \left(\frac{N_i}{D_i} - \frac{N_j}{D_j} \right) = \pm 2^f (x_i - x_j)$$

We can fix $x_1 = 2^a$, any power of 2. Then, by lemma 1, we should choose the other x_i in the set $\{-2^a, 2^{a+1}, 2^{a-1}\}$, but for any two of them we will never have $x_i - x_j$ being a power of two.

4 Optimality criteria

In the following, for a square matrix M we indicate with $M[i, j]$ its entry in position (i, j) and with $M^{(i)}$ its i^{th} line. We begin with some definitions:

Definition 1. The **support** of $M^{(i)}$ is the set $s(M^{(i)})$ of column indexes $j \in \mathbb{N}$ such that $M[i, j] \neq 0$. The **support** of M is the set $s(M)$ of pairs $(i, j) \in \mathbb{N} \times \mathbb{N}$ such that $M[i, j] \neq 0$. The cardinality of a line support $s(M^{(i)})$ is indicated with $\#M^{(i)}$, and of the matrix support $s(M)$ with $\#M$.

Definition 2. We note with $\text{gcd}_i(M) = \text{gcd}(M^{(i)})$ the greatest common divisor of all coefficients in $M^{(i)}$.

Inversion sequences are sequences of basic operations changing the initial matrix and producing as last matrix the identity matrix. Gauss' method is e.g. a classical algorithm to produce effective IS, reducing at each step the *maximum* possible cardinality of still-to-be-analysed lines support. However, it uses as basic operations, apart from sums and subtractions, multiplications and division *at each step*, and from a computer point of view this is not quite optimal when looking for efficiency.

Gauss' method merit is that it is very efficient for *general* matrices, but in our case, where matrices have a very particular form, it may be possible to find different SI which are, in terms of nature and number of used basic operations, more efficient.

We propose here some optimality criteria, that guarantee termination of our exhausting algorithm. For A_n matrices as specified in equation (1), we have

$$\#A_n = 1 + (2n - 1 - 2)(2n - 1) + 1 = ((2n - 1)^2 - 2(2n - 1) + 1) + 1 = (2n - 1 - 1)^2 + 1 = 4(n - 1)^2 + 1$$

We now list the used criteria: we suppose that the basic row operations involve lines $M^{(i_1)}, M^{(i_2)}$, and that $M^{(i_1)}$ is overwritten with the result. We indicate with \widetilde{M} the matrix *after* the execution of the operation.

- A) Support reduction criterion :** $\#\widetilde{M}^{(i)} < \#M^{(i)}$. The resulting line must have at least one more coefficient with value 0.
- B) Alignment criterion :** $\widetilde{M}[i_1, j_1]/M[i_2, j_1] = \widetilde{M}[i_1, j_2]/M[i_2, j_2]$. The resulting line must have more entries differing from the corresponding ones in another line by an equal multiplicative factor.

Criterion **(A)** is sufficient to guarantee termination of the exhaustive research algorithm that we'll describe in the following section. The solution actually implemented by GMP applies once criterion **(B)**.

Gauss' method uses a support lexicographic ordering. Supposing no row and column permutation is necessary, at step k we have that $n_2 = k$ for all following lines whose k^{th} entry is not 0. Their support is reduced (coefficient is set to 0) starting from the leftmost entry. Note that Gauss' method can lead to intermediate step in which the matrix is more dense than the precedent one, while we do not permit this.

5 The Toom graph

We note the identity matrix of order s with I_s ; if the order is not relevant, simply with I .

Let an appropriate set of interpolation values $\{v_i\}$ be given, that is, let A_n be known from the beginning. We perform an exhaustive search on all the possible IS, according to the criteria introduced in section 4. This can be modeled by a oriented and weighted graph (N, E) (which we call Toom graph), in which each node $\nu \in N$ is represented by the matrix M_ν . We call α the distinguished "initial" node, such that $M_\alpha = A_n$, and similarly ω the corresponding "final" node for I . Edges are described below.

In order to have a sufficiently easy theoretical modelisation of the graph procedure analysis, we will consider as basic operations the following ones (i, i_1, i_2 are row indexes and $c, c_1, c_2 \in \mathbb{Z}$):

1. $M_{\nu'}^{(i_1)} = \widetilde{M}_\nu^{(i_1)} = c_1 M_\nu^{(i_1)} + c_2 M_\nu^{(i_2)}$
2. $M_{\nu'}^{(i)} = \widetilde{M}_\nu^{(i)} = \frac{M_\nu^{(i_1)}}{c} = \left(\frac{M_\nu[i, 1]}{c}, \frac{M_\nu[i, 2]}{c}, \dots, \frac{M_\nu[i, 2n-1]}{c} \right)$

where each of the divisions by c is exact. The edges $\varepsilon = (\nu_1, \nu_2) \in E$ are identified by the quadruple (i_1, i_2, c_1, c_2) or by the pair (i, c) , respectively. With this in mind, IS can simply be considered as edges sequences, or, in graph terminology, as paths joining α to ω .

To consider the different computational cost of row elementary operations depending on c, c_1, c_2 values in graph analysis, we introduce some constant weights, named in table 1 (possibly swapping $|c_1|$ and $|c_2|$ values). COMBINATION_WEIGHT is practically the basic cost of an addition/subtraction of two long integers. An appropriated tuning of these constants will result in choosing a particular criterion of graph visit.

Definition 3. The weight $w(\varepsilon)$ of an edge $\varepsilon \in E$ is the constant defined by table 1 according to the values of c_1, c_2 (or simply c , in case 2). The weight w_{IS} of a IS $= (\varepsilon_1, \dots, \varepsilon_l)$, where $\varepsilon_i = (\nu_i, \nu_{i+1})$ and $M_{\nu_{l+1}} = I$, is $w_{IS} = w(\varepsilon_1) + \dots + w(\varepsilon_l)$. An IS is **minimal** if its weight is minimal among the weights of all possible IS connecting ν_1 to ν_{l+1} . The weight $w(M)$ of a matrix M is w_{IS} of a minimal IS for M . The weight $w(\nu)$ of a node ν is defined as $w(\nu) = w(M_\nu)$.

	$ c_1 $	$ c_2 $	Operation	Weight
	1	1	+ or -	COMBINATION_WEIGHT
1.	$= 2^k > 1$	1	\pm , shifting	'' + $_1_2_WEIGHT$
	$> 1, \neq 2^k$	1	\pm , multiplication	'' + $_1_X_WEIGHT$
	$= 2^k > 1$	$> 1, \neq 2^k$	\pm , shifting and multiplication	'' + $_2_X_WEIGHT$
	> 1	> 1	general linear combination	'' + $_X_Y_WEIGHT$

	$ c $	Operation	Weight
2.	1	(may be unitary -)	~ 0
	$= 2^k > 1$	shifting	SHIFT_WEIGHT
	$> 1, \neq 2^k$	/	DIVISION_WEIGHT

Table 1. Weight constants for elementary row operations

Obviously, $w(I) = 0$.

Example (Karatsuba graph): Let $(v_0 = \infty, v_1 = 1, v_2 = 0)$ identify Karatsuba matrix A_2 . We have

$$\begin{array}{ccc}
 A_2 & \xrightarrow{\varepsilon_1} & M_1 \\
 \varepsilon_2 \downarrow & & \downarrow \varepsilon_3 \\
 M_2 & \xrightarrow{\varepsilon_4} & I_2
 \end{array}
 \quad \text{with } A_2 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}, M_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}, M_2 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Example (Knuth graph): Let $(v_0 = \infty, v_1 = -1, v_2 = 0)$ identify the Knuth matrix A'_2 . We have

$$\begin{array}{ccc}
 A'_2 & \xrightarrow{\varepsilon_1} & M'_1 \\
 \varepsilon_2 \downarrow & & \downarrow \varepsilon_3 \\
 M'_2 & \xrightarrow{\varepsilon_4} & I_2
 \end{array}
 \quad \text{with } A'_2 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & -1 & 1 \\ 0 & 0 & 1 \end{pmatrix}, M'_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix}, M'_2 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

These almost trivial graphs show that in the Toom graph generated by a fixed matrix A_n there are many minimal IS, and, that different matrices may have equivalent minimal IS. We call A_n -graph the graph generated by a matrix A_n . In the following, we consider Toom graphs generated by A_n (or A_{n_1, n_2}) matrices.

6 Travelling through Toom graphs

What we're really trying to do is practically to solve a minimum weight path problem between α and ω .

Toom graphs are very big. Beginning with Toom-3, it is easy to work with dozens of thousands nodes, and, in order to cope both with memory limitation and time reduction, we developed some strategies not to consider/add some nodes to the graphs. Note that a node ν (and therefore all of its neighbors) might be visited many times, since in general there are many different sequences ("paths") of elementary row operations joining ν with another node ν' .

In general, two opposite approaches are possible to automatically visit the graph by means of a computer program: the *purely functional* one, which does not need the graph to be really constructed, and a second one for which the graph is actually built.

- It is clear that a recursing function f visiting a graph G would have no computer memory occupation problems (Gauss' method gives immediately a reasonable upper bound l for IS lengths), but this approach would be extremely lengthy in time. Infact, in this case G would be implicitly represented by the stack of f nested callings, with length $\leq K \times l$ for a certain fixed factor K , and the amount of needed memory can be precisely estimated from the very beginning (avoiding a computational "explosion" because of space lacking). But there is a "time" disadvantage. When an instance of f ends its work finishing analysing a node ν , all the information about the subgraph G' having ν as initial node is lost. If ν appears again, reached by another path, G' must be completely analysed again.
- On the other hand, f could effectively memorize the nodes of a graph while analysing it, in order to recognise if a node has already been inserted, avoiding thus to loose time. Unfortunately, as said, these graphs are huge, and without a tuned analysis management they grow so much as to rapidly fill all the available memory.

We try to stay someway in the middle, keeping only *some* nodes for *some* time, so as to avoid until possible – possibly before filling up all the available memory – to repeat analysing nodes, but we also apply the below explained considerations in order to skip “trivial” graph analysis, when the path(s) from the current node to the identity node are trivial or can be sufficiently easily estimated.

6.1 A guided tour of the graph

Let $f(\alpha)$ be, as indicated before, the program visiting the graph having α as its initial node. Keeping in mind that we want to obtain IS with lowest possible weight, we introduce a second argument for f , the DESIRED_WEIGHT (DW). Its meaning is related to the fact that we’re interested only in IS with $w_{IS} \leq DW$, marking as “not interesting” the nodes ν with $w(M_\nu) > DW$.²

We implement this marking not as a boolean variable updating, but as two weights (to be compared) we associate to every node ν . These weights are a lower and upper bound of $w(\nu)$, respectively, and are indicated with $m_w(\nu)$ and $\mathcal{M}_w(\nu)$ (for min and max). The function f returns as result the pair $(m_w(\nu), \mathcal{M}_w(\nu))$. The underlying idea is that making use of DW , it is not always really necessary to know precisely the exact weight of a node to classify a matrix as interesting or not.

In the analysis made by $f(\nu, t)$ – where t is the desired weight – for a node ν the following may happen:

$m_w(\nu) = \mathcal{M}_w(\nu)$: this means that ν has already been completely analysed, and that $w(\nu) = m_w(\nu) = \mathcal{M}_w(\nu)$. Then f immediately returns the pair $(m_w(\nu), \mathcal{M}_w(\nu))$, with no further computations.

$t < m_w(\nu)$: then f does not need to recurse over and over again, because, so to say, a sufficient number of operations to invert M_ν is not permitted. Then f immediately returns the pair $(m_w(\nu), \mathcal{M}_w(\nu))$. Note that this may happen even if $w(M_\nu)$ has still not been determined exactly, in particular even if ν has never been found before in the graph.

In the above cases, no further analysis is carried beyond ν , and this make the graph analysis faster. Suppose now $f(\nu, t)$ is executing, and let ν' be a neighbor of ν , joined to ν by the edge ε , and $t' = t - w(\varepsilon)$. The recursive call on ν' will be $f(\nu', t')$, and we call $w = (m_w(\nu'), \mathcal{M}_w(\nu'))$ the returned value.

At this point, the “parent” instance of f receiving w must try the other neighbors of ν , but before doing this, t is updated: $t = \min\{t, m_w(\nu') + w(\varepsilon)\}$. This is a very important step, because if f succeeded in inverting $M_{\nu'}$ with a IS having weight $m_w(\nu') = \mathcal{M}_w(\nu') < t'$, the new value t for DW will be strictly smaller than the precedent one, and all the paths that will be analyzed from now will have a stronger restriction, so that pruning is likely to become more and more effective.

6.2 Heuristics for weight estimates

It is clear that the better estimate one is able to give for a matrix weight $w(M)$, the shorter and more efficient the graph analysis will be, both in space and time requirements. we present here some of the ideas we implemented in order to give a as best as possible estimate.

A first trivial estimate for $m_w(\nu)$ and $\mathcal{M}_w(\nu)$ can be obtained from the support of M_ν and of the single lines and columns, considering that $\#I_m = m$.

Definition 4. A matrix row $M^{(i)}$ is a **singleton** line if all its entries but one are zero, and the not zero entry has value $v = 1$ or $v = -1$. It is **almost singleton** if $|v| \neq 1$. Similarly for columns. A singleton line is indicated with SL_j , an almost singleton line with ASL_j , where j is the position of v .

Proposition 3. Let $M_\nu \in GL(m, \mathbb{Z})$ with $\#M_\nu = s$, ns_l the number of its non-singleton lines, and ns_c the number of its non-singleton columns. Then $m_w(\nu) \geq \max\{ns_l, ns_c\} \times \text{COMBINATION_WEIGHT}$ and $\mathcal{M}_w(\nu) \leq (s - m) \times (\text{COMBINATION_WEIGHT} + _X_Y_WEIGHT) + ns_l \times \text{DIVISION_WEIGHT}$

Proof. The common idea is that every non-singleton line must be made singleton. For m_w it is sufficient to consider the best possible case (just one operation for each line to update), while for \mathcal{M}_w consider the worst possible case, in which each operation of type 1 produces just one more zero, so there are $s - m$ of them, followed by a final division for each non-singleton line, to make it monic.

² In a certain sense, DW works as a threshold, useful to prune Toom graphs analysis.

Obviously, the more accurate the estimates are, the less nodes are considered/build in the graph, because the function f has more probabilities to obtain sufficient information from the estimates. Weight estimation can be performed line by line, just summing up the corresponding values. There is an interesting particular case:

Proposition 4. *Let $M_\nu \in GL(m, \mathbb{Z})$, and $M_\nu^{(i)} = (a, 0, \dots, 0, b, 0, \dots, 0, c)$, with b in j^{th} position, different from the first and the last one. Then this line-weight estimate can be computed exactly.*

Proof. The entries a and c can be set to zero only using respectively the first and last line (both singleton). The minimal (partial) IS depends on a, b and c values. A default strategy could then be:

- Set a to zero using the first line.
- Set c to zero using the last line.
- If $b \neq 1$, divide the i^{th} line by b .

but the three operations might be executed in a different order, depending on the values of weights associated to the basic operations (table 1): for example, it could be more convenient first to divide and then to subtract twice

$$(3, 0, \dots, 0, 3, 0, \dots, 0, 3) \Rightarrow (1, 0, \dots, 0, 1, 0, \dots, 0, 1) \Rightarrow (0, \dots, 0, 1, 0, \dots, 0, 1) \Rightarrow S_j$$

or division could be better done in the middle:

$$(1, 0, \dots, 0, 3, 0, \dots, 0, 3) \Rightarrow (0, 0, \dots, 0, 3, 0, \dots, 0, 3) \Rightarrow (0, \dots, 0, 1, 0, \dots, 0, 1) \Rightarrow S_j$$

The various sequence possibilities correspond to different coefficients in basic operations (different weight constants are used). A very small number of coefficient values comparisons is sufficient to obtain the result, without having to recurse more on this line entries.

Another interesting case: exactly one not singleton line. Here the global analysis is precise, too.

Proposition 5. *Let $M_\nu \in GL(m, \mathbb{Z})$ with $\#\{i \mid M_\nu^{(i)} \text{ is not singleton}\} = 1$. Then the weight estimate can be computed exactly.*

Proof. In this case $m - 1$ basic operation of type 1 (and possibly 1 division) are surely needed, because at most one new 0 can appear at each step. Similar considerations as the ones appearing in the precedent case apply here, in order to recognise which is the best sequence order.

Similar considerations may be done, *mutatis mutandis*, reasoning with columns.

7 Implementation and results

The authors developed code in C++ based on the STL library to make experiments. In this section we expose some of the obtained results. The basic operations are indicated between matrices using the following notation:

$i \pm = j$: j^{th} line is destructively added to or subtracted from i^{th} line
$i -= (c)j$: j^{th} line multiplied by c is destructively subtracted from i^{th} line
$(c)i \pm = j$: j^{th} line is destructively added or subtracted from i^{th} line multiplied by c
$i /= (c)$: the entries in i^{th} line are all destructively divided by c
$i \gg (c)$: the entries in i^{th} line are all destructively divided by 2^c

7.1 Toom-2.5

This case could also be treated by hand, and we present it as a checking test. Starting from the matrix defined by the reasonably good values $\{\infty, -1, 1, 0\}$ we obtained the (expected) IS indicated below. Note that $\det(A_{3,2}) = 2$, and therefore at least a division by 2 (shifting) is necessary. The nodes that were really inserted in the Toom-graph are 44.

$$A_{3,2} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2=-3} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & -1 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[3=-1]{2 \gg (1)} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & -1 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[2=-4]{3+=2} I_4$$

7.2 Toom-3

GMP solution : We first present the IS used by GMP library starting from the matrix obtained by using $\{\infty, 2, -1, 1, 0\}$, and compute its weight w_{GMP} according to our definitions.

$$\begin{aligned}
 A_3 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 16 & 8 & 4 & 2 & 1 \\ 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2+=(2)3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 18 & 6 & 6 & 0 & 3 \\ 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[3+=-4]{2/=(3)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 6 & 2 & 2 & 0 & 1 \\ 2 & 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2+=-5} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 6 & 2 & 2 & 0 & 2 \\ 2 & 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
 &\xRightarrow[3\gg(1)]{2\gg(1)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 3 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2-=(2)1} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{4-=-2} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2-=-3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[3-=-5]{3-=-1} I_5
 \end{aligned}$$

$$w_{GMP} = 8 \cdot \text{COMBINATION_WEIGHT} + \text{DIVISION_WEIGHT} + 2 \cdot \text{SHIFT_WEIGHT} + 2 \cdot \text{_1_2_WEIGHT}$$

Surprisingly enough, just by applying criterion **(A)** our program was able to compute solutions having a weight less than GMP's. We present here two of them, both obtained from the matrix defined by the interpolating values $\{\infty, -1, 1, \frac{1}{2}, 0\}$. Because of matrix central symmetry, the same IS (with adapted line and column indexes) can be applied to the matrix used by GMP. The general weight we obtain is

$$\begin{aligned}
 &8 \cdot \text{COMBINATION_WEIGHT} + \text{DIVISION_WEIGHT} + \text{SHIFT_WEIGHT} + \\
 &\min(\text{_1_X_WEIGHT}, \text{SHIFT_WEIGHT}) + \text{_1_2_WEIGHT}
 \end{aligned}$$

and depending on `SHIFT_WEIGHT` and `_1_X_WEIGHT` values, the below indicated solutions are obtained.

1st solution : `SHIFT_WEIGHT > _1_X_WEIGHT`

Note that $\det(A_3) = 12$, and the division is “split” into a shifting and a division by 6. By using estimates, the nodes that were really inserted in the Toom-graph are only 11205. Considering also criterion **(B)** – applied in the IS used by GMP – we have that the graph get bigger, reaching 123213 nodes, but the optimal found solution \overline{IS} does not change.

Moreover, *independently from weight values*, one has $w(\overline{IS}) < w_{GMP}$.

$$\begin{aligned}
 A_3 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 & 16 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{4-=-2} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 3 & 3 & 9 & 15 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2-=-3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & -2 & 0 & -2 & 0 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 3 & 3 & 9 & 15 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[2/=(-1)]{3-=-1} \\
 \tilde{A}_3 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 2 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 3 & 3 & 9 & 15 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[2\gg(1)]{4-=(3)3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 6 & 12 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[3-=-2]{4/=(6)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow[4-=(2)5]{3-=-5} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xRightarrow{2-=-4} I_5
 \end{aligned}$$

2nd solution : SHIFT_WEIGHT < .1 X_WEIGHT

Here only 11862 nodes are generated. The initial operations are the same as above: the different behavior begins from the fourth matrix \tilde{A}_3 .

$$\tilde{A}_3 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 2 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 3 & 3 & 9 & 15 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[2 \gg (1)]{4 / = (3)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 3 & 5 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{4 - = 3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 2 & 4 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[3 - = 2]{4 \gg (1)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[4 - = (2)5]{3 - = 5} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{2 - = 4} I_5$$

With our program we tested all the matrices obtained from interpolating points $\{\infty, a, b, c, 0\}$, where $a, b, c \in \{\pm 1, \pm 2, \dots, \pm 16, \pm \frac{1}{2}, \pm \frac{1}{3}, \dots, \pm \frac{1}{16}\}$. No IS with weight smaller than the two presented above was found. Paths with the same weight are substantially equivalent to these ones and were found starting from interpolating points $\{\infty, 1, -1, \pm 2, 0\}$ and $\{\infty, 1, -1, \pm \frac{1}{2}, 0\}$.

7.3 Toom-3.5

Starting from the matrix defined by the values $\{\infty, 2, -2, 1, -1, 0\}$ we obtained the sequence of operations indicated below, with weight

$$12 \cdot \text{COMBINATION_WEIGHT} + 2 \cdot \text{DIVISION_WEIGHT} + 2 \cdot \text{SHIFT_WEIGHT} + 2 \cdot \text{.1.2_WEIGHT}$$

$$A_{4,3} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 32 & 16 & 8 & 4 & 2 & 1 \\ -32 & 16 & -8 & 4 & -2 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ -1 & 1 & -1 & 1 & -1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[5 - = 4]{3 - = 2} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 32 & 16 & 8 & 4 & 2 & 1 \\ -64 & 0 & -16 & 0 & -4 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ -2 & 0 & -2 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[3 \gg (1)]{2 - = 6} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 32 & 16 & 8 & 4 & 2 & 0 \\ -32 & 0 & -8 & 0 & -2 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ -2 & 0 & -2 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\xrightarrow[5 / = (-1)]{2 + = 3, 4 - = 6} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 16 & 0 & 4 & 0 & 0 \\ -32 & 0 & -8 & 0 & -2 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 \\ 2 & 0 & 2 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[5 \gg (1)]{3 + = 5} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 16 & 0 & 4 & 0 & 0 \\ -30 & 0 & -6 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\xrightarrow[4 - = 5]{3 / = (-6)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 16 & 0 & 4 & 0 & 0 \\ 5 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[2 - = (4)4]{3 - = (4)1} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 12 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[2 / = (12)]{5 - = 3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \xrightarrow[4 - = 2]{3 - = 1} I_6$$

8 Conclusions

We presented a method to determine an optimal sequence of basic operations inverting a particular matrix appearing in Toom multiplication methods through an optimized exhaustive research on graphs. Weights for each basic operation are introduced, which determine the inversion cost for the matrices appearing in the intermediate steps. An interval-like estimate analysis permit to know in advance that a path in the graph is not optimal, saving thus time and space.

Intermediate version of Toom methods are also presented, to show more results. As an interesting result, we were able to determine two new inversion sequences for Toom-3 matrix, which seem, to the best of our knowledge, to be new, and good alternatives to the one proposed in the actual version of GMP library.

References

1. Chung, J., Anwar Hasan, M. *Asymmetric squaring formulae* Technical Report, CACR 2006-24, University of Waterloo (2006) http://ln.bodrato.it/cacr2006-24_pdf
2. Cook, S.A. *On the Minimum Computation Time of Functions* Thesis, Harvard University, pp. 51-77 (1966)
3. The GNU Multiplication Library (GMP) documentation <http://www.swox.com/gmp/#DOC>
4. Karatsuba, A., Ofman, Yu. *Multiplication of multidigit numbers on automata* Soviet Physics-Doklady, 7, 595-596 (1963); translation from Dokl. Akad. Nauk SSSR, 145:2, 293-294, (1962)
5. Knuth, D.E. *The Art of Computer Programming, Vol. 2, Second Edition* Addison-Wesley, Reading Mass., Chapter 4, Section 3.3, pp. 278-301 (1981)
6. Toom, A.L. *The Complexity of a Scheme of Functional Elements Realizing the Multiplication of Integers* Soviet Mathematics, Vol. 3, pp. 714-716 (1963)
7. Schonhage, A., Strassen, V. *Schnelle Multiplikation großer Zahlen* Computing 7 pp. 281-292 (1971)
8. St Denis, T., Rasmussen, M., Rose, G. *Multi-precision math* (tommath library documentation) <http://math.libtomcrypt.com/files/tommath.pdf>
9. Zuras, D. *On Squaring and Multiplying Large Integers* 11th IEEE Symposium on Computer Arithmetic pp. 260-271 (1993)

Appendix A: Toom-4, Toom-4.5 and Toom-5

We present here some found IS for the indicated methods, that were obtained not in a completely automatic way. Because of limited computing resources, some initial steps were done manually, while the general automatical analysis starts some matrix later.

Toom-4

Starting from the matrix defined by the values $\{\infty, 2, 1, -1, \frac{1}{2}, -\frac{1}{2}, 0\}$ we obtained the sequence of operations indicated below, with weight

$$18 \cdot \text{COMBINATION_WEIGHT} + 3 \cdot \text{DIVISION_WEIGHT} + \text{SHIFT_WEIGHT} + \min(.1_X_WEIGHT, \text{SHIFT_WEIGHT}) + 2 \cdot .1_X_WEIGHT + 4 \cdot (.1_2_WEIGHT)$$

Depending on SHIFT_WEIGHT and .1_X_WEIGHT values, we present the two corresponding solutions.

1st solution : SHIFT_WEIGHT > .1_X_WEIGHT

$$\begin{aligned}
 A_4 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 & 1 & -1 & 1 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 \\ 1 & -2 & 4 & -8 & 16 & -32 & 64 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[4^--=3]{\begin{matrix} 2^+=5 \\ 6^--=5 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 65 & 34 & 20 & 16 & 20 & 34 & 65 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 & 0 & -2 & 0 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 \\ 0 & -4 & 0 & -16 & 0 & -64 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[4^{\gg}1]{\begin{matrix} 5^--=1 \\ 5^--=(64)7 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 65 & 34 & 20 & 16 & 20 & 34 & 65 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & -1 & 0 & -1 & 0 & -1 & 0 \\ 0 & 2 & 4 & 8 & 16 & 32 & 0 \\ 0 & -4 & 0 & -16 & 0 & -64 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
 & \xrightarrow[(2)5^+=6]{\begin{matrix} 3^+=4 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 65 & 34 & 20 & 16 & 20 & 34 & 65 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & -1 & 0 & -1 & 0 & -1 & 0 \\ 0 & 0 & 8 & 0 & 32 & 0 & 0 \\ 0 & -4 & 0 & -16 & 0 & -64 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[3^--=1]{\begin{matrix} 2^--=(65)3 \\ 4^/=(\leftarrow)1 \\ 6^/=(\leftarrow)1 \\ 3^--=7 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 34 & -45 & 16 & -45 & 34 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 8 & 0 & 32 & 0 & 0 \\ 0 & 4 & 0 & 16 & 0 & 64 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[5^--=(8)3]{\begin{matrix} 2^+= (45)3 \\ 5^--=(8)3 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 34 & 0 & 16 & 0 & 34 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 24 & 0 & 0 & 0 \\ 0 & 4 & 0 & 16 & 0 & 64 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[2^--=(16)4]{\begin{matrix} 5^/=(24) \\ 6^--=2 \end{matrix}} \\
 \tilde{A}_4 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 18 & 0 & 0 & 18 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & -30 & 0 & 0 & 30 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \xrightarrow[3^--=5]{\begin{matrix} 2^/=(18) \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & -30 & 0 & 0 & 30 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \xrightarrow[6^+= (30)2]{\begin{matrix} 4^--=2 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 60 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \xrightarrow[6^/=(60)]{\begin{matrix} 6^/=(60) \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \xrightarrow[2^--=6]{\begin{matrix} 2^--=6 \end{matrix}} I_7
 \end{aligned}$$

2nd solution : SHIFT_WEIGHT < .1_X_WEIGHT

The initial operations are the same as above: the difference begins from the matrix \tilde{A}_4 .

$$\tilde{A}_4 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 18 & 0 & 0 & 18 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & -30 & 0 & 0 & 0 & 30 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[6/=(30)]{2/=(18) \atop 3/=-5} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[6+=2]{4/=-2 \atop 7+=4} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{6 \gg (1)} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{2/=-6} I_7$$

Toom-4.5

Starting from the matrix defined by the values $\{\infty, -1, -2, \frac{1}{2}, 1, 2, -\frac{1}{2}, 0\}$ we obtained the sequence of operations indicated below, with weight

$$22 \cdot \text{COMBINATION_WEIGHT} + 4 \cdot \text{DIVISION_WEIGHT} + \text{SHIFT_WEIGHT} + 3 \cdot \text{.1_X_WEIGHT} + 6 \cdot \text{.1_2_WEIGHT}$$

$$A_{5,4} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 1 & -1 & 1 & -1 & 1 & 1 \\ -128 & 64 & -32 & 16 & -8 & 4 & -2 & 1 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 & 128 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 1 & -2 & 4 & -8 & 16 & -32 & 64 & -128 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[3/=-6]{2+=5 \atop 7+=4} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 2 & 0 & 2 & 0 & 2 \\ -256 & 0 & -64 & 0 & -16 & 0 & -4 & 0 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 & 128 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 2 & 0 & 8 & 0 & 32 & 0 & 128 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\xrightarrow[2 \gg (1)]{4/=-2 \atop (2)6+=3} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ -256 & 0 & -64 & 0 & -16 & 0 & -4 & 0 \\ 1 & 0 & 4 & 6 & 16 & 30 & 64 & 126 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 128 & 0 & 32 & 0 & 8 & 0 & 2 \\ 2 & 0 & 8 & 0 & 32 & 0 & 128 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[(2)3+=7]{5/=-2 \atop 6/=-128} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ -510 & 0 & -120 & 0 & 0 & 0 & 120 & 0 \\ 1 & 0 & 4 & 6 & 16 & 30 & 64 & 126 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -96 & 0 & -120 & 0 & -126 \\ 2 & 0 & 8 & 0 & 32 & 0 & 128 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\xrightarrow[2 \gg (1)]{4+=6 \atop 2/=-8} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ -510 & 0 & -120 & 0 & 0 & 0 & 120 & 0 \\ 1 & 0 & 4 & -90 & 16 & -90 & 64 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -96 & 0 & -120 & 0 & -126 \\ 2 & 0 & 8 & 0 & 32 & 0 & 128 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[7/=-32]{(2)4/=-7 \atop 6+=126} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ -510 & 0 & -120 & 0 & 0 & 0 & 120 & 0 \\ 0 & 0 & 0 & -180 & 0 & -180 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -96 & 0 & -120 & 0 & 0 \\ -30 & 0 & -24 & 0 & 0 & 0 & 96 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\xrightarrow[3+=510]{(5)7/=-3 \atop 4/=-180 \atop 6/=-24} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & -120 & 0 & 0 & 120 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 4 & 0 & 5 \\ 360 & 0 & 0 & 0 & 0 & 360 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[2/=-4]{3/=-120 \atop 7/=-360 \atop 6/=-4} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[7/=-1]{5/=-7 \atop 4/=-6} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[5/=-3]{3+=7} I_8$$

Toom-5

Starting from the matrix defined by the values $\{\infty, -2, \frac{1}{2}, 4, 2, -1, 1, -\frac{1}{2}, 0\}$ we obtained the sequence of operations indicated below, with weight

$$32 \cdot \text{COMBINATION_WEIGHT} + 5 \cdot \text{DIVISION_WEIGHT} + 2 \cdot \text{SHIFT_WEIGHT} + 6 \cdot \text{.1_X_WEIGHT} + 8 \cdot \text{.1_2_WEIGHT}$$

$$\begin{aligned}
 A_5 &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 256 & -128 & 64 & -32 & 16 & -8 & 4 & -2 & 1 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 & 128 & 256 \\ 4^8 & 4^7 & 4^6 & 4^5 & 256 & 64 & 16 & 4 & 1 \\ 256 & 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -2 & 4 & -8 & 16 & -32 & 64 & -128 & 256 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow[8=-3]{\begin{matrix} 6=-7 \\ 2=-5 \\ 4=-9 \\ 4=(2^{16})1 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -256 & 0 & -64 & 0 & -16 & 0 & -4 & 0 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 & 128 & 256 \\ 0 & 4^7 & 4^6 & 4^5 & 256 & 64 & 16 & 4 & 0 \\ 256 & 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 0 & -2 & 0 & -2 & 0 & -2 & 0 & -2 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & -4 & 0 & -16 & 0 & -64 & 0 & -256 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
&\xrightarrow[8/=(-1)]{\begin{matrix} 6 \gg (1) \\ (2)5+=2 \\ 2/=(-1) \\ 8/=(-1) \end{matrix}} \xrightarrow[6/=(-1)]{7+=6} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 256 & 0 & 64 & 0 & 16 & 0 & 4 & 0 \\ 1 & 2 & 4 & 8 & 16 & 32 & 64 & 128 & 256 \\ 0 & 4^7 & 4^6 & 4^5 & 256 & 64 & 16 & 4 & 0 \\ 512 & 0 & 128 & 0 & 32 & 0 & 8 & 0 & 2 \\ 0 & -1 & 0 & -1 & 0 & -1 & 0 & -1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 4 & 0 & 16 & 0 & 64 & 0 & 256 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
&\xrightarrow[5-=(512)7]{\begin{matrix} 3=-7 \\ 5-=(512)7 \end{matrix}} \xrightarrow[8+=2]{\begin{matrix} (2)3=-8 \\ 7=-1 \\ 7=-9 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 256 & 0 & 64 & 0 & 16 & 0 & 4 & 0 \\ 0 & 0 & 6 & 0 & 30 & 0 & 126 & 0 & 510 \\ 0 & 4^7 & 4^6 & 4^5 & 256 & 64 & 16 & 4 & 0 \\ 0 & 0 & -384 & 0 & -480 & 0 & -504 & 0 & -510 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 4 & 0 & 16 & 0 & 64 & 0 & 256 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
&\xrightarrow[4=-2]{\begin{matrix} 5+=3 \\ 8-=(80)6 \\ 3-=(510)9 \end{matrix}} \xrightarrow[2 \gg (2)]{\begin{matrix} (3)3+=5 \\ 8/=(180) \\ 5+=(378)7 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 256 & 0 & 64 & 0 & 16 & 0 & 4 & 0 \\ 0 & 0 & 6 & 0 & 30 & 0 & 126 & 0 & 0 \\ 0 & 16128 & 4096 & 960 & 256 & 48 & 16 & 0 & 0 \\ 0 & 0 & -378 & 0 & -450 & 0 & -378 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 180 & 0 & 0 & 0 & 0 & 0 & 180 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
&\xrightarrow[2=-8]{\begin{matrix} 6=-2 \\ 5/=(72) \\ 3/=(-360) \end{matrix}} \xrightarrow[3=-5]{\begin{matrix} 7=-3 \\ 4-=(256)5 \end{matrix}} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 63 & 0 & 16 & 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 16128 & 4096 & 960 & 256 & 48 & 16 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & -63 & 0 & -15 & 0 & -3 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}
 \end{aligned}$$

$$\begin{array}{l}
4-=(4096)3 \\
4-=(16)7 \\
4+=(256)6 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 00 & 00 & 0000 \\
0 & 630 & 160 & 4000 \\
0 & 01 & 00 & 0000 \\
0 & 00 & -2880 & 0 -720 000 \\
0 & 00 & 01 & 0000 \\
0 & -630 & -150 & -3000 \\
0 & 00 & 00 & 0100 \\
0 & 10 & 00 & 0010 \\
0 & 00 & 00 & 0001
\end{pmatrix}
\begin{array}{l}
6+=2 \\
\Rightarrow \\
(180)2+=4
\end{array}
\begin{pmatrix}
1 & 00 & 00 & 0000 \\
0 & 11340 & 00 & 0000 \\
0 & 01 & 00 & 0000 \\
0 & 00 & -2880 & 0 -720 000 \\
0 & 00 & 01 & 0000 \\
0 & 00 & 10 & 1000 \\
0 & 00 & 00 & 0100 \\
0 & 10 & 00 & 0010 \\
0 & 00 & 00 & 0001
\end{pmatrix}$$

$$\begin{array}{l}
2/=(11340) \\
4+=(720)6 \\
\Rightarrow
\end{array}
\begin{pmatrix}
100 & 000000 \\
010 & 000000 \\
001 & 000000 \\
000 & -2160 00000 \\
000 & 010000 \\
000 & 101000 \\
000 & 000100 \\
010 & 000010 \\
000 & 000001
\end{pmatrix}
\begin{array}{l}
4/=(-2160) \\
\Rightarrow \\
6-=-4 \\
8-=-2
\end{array}
I_9$$

Appendix B: Application to asymmetrical squaring

Chung and Anwar Hasan showed different linear systems which can be built, but only for squaring. We refer to their report [1] for all the details. Here we just present the result given by our software applied to one of their matrices, although not optimized for this case.

The report proposed an inversion algorithm for the 5-way squaring method using temporary variables with a total cost of $18 \cdot \text{COMBINATION_WEIGHT} + 7 \cdot \text{SHIFT_WEIGHT}$. Starting from the same matrix, our program found the following IS, with no temporary variables use and smaller weight:

$$16 \cdot \text{COMBINATION_WEIGHT} + 3 \cdot \text{SHIFT_WEIGHT}$$

$$A_5^s = \begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 \\
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 & 1 \\
0 & 1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 \\
1 & 1 & 0 & -1 & -1 & 0 & 1 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
3+=4 \\
\Rightarrow \\
7-=-1 \\
7-=-8 \\
7-=-2
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
2 & 0 & 2 & 0 & 2 & 0 & 2 & 0 & 2 \\
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 & 1 \\
0 & 1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
3\gg(1) \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 & 1 \\
0 & 1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}$$

$$\begin{array}{l}
5+=3 \\
4-=-3 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\
2 & 0 & 0 & 0 & 2 & 0 & 0 & 0 & 2 \\
0 & 1 & 0 & -1 & 0 & 1 & 0 & -1 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
5\gg(1) \\
6+=4 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\
1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
0 & 2 & 0 & 0 & 0 & 2 & 0 & 0 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
4-=-8 \\
6\gg(1) \\
3-=-5 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\
0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\
1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}$$

$$\begin{array}{l}
4-=-6 \\
5-=-1 \\
5-=-9 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & -1 & -1 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
6-=-2 \\
7+=4 \\
7+=5 \\
\Rightarrow
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
\end{pmatrix}
\begin{array}{l}
3-=-7 \\
\Rightarrow
\end{array}
I_9$$