

# How to compute $\pi$ to $10^{12}$ digits

A crash course in high precision arithmetics

Jörg Arndt      <arndt@jjj.de>

## The radix-2 DIF FFT algorithm

Splitting of the Fourier sum into a left and right half leads to the *decimation in frequency* (DIF) FFT algorithm, also called *Sande-Tukey* FFT algorithm.

For even values of  $n$  the  $k$ -th element of the Fourier transform is

$$\mathcal{F}[a]_k = \sum_{x=0}^{n-1} a_x z^{xk} = \sum_{x=0}^{n/2-1} a_x z^{xk} + \sum_{x=n/2}^n a_x z^{xk} \quad (1a)$$

$$= \sum_{x=0}^{n/2-1} a_x z^{xk} + \sum_{x=0}^{n/2-1} a_{x+n/2} z^{(x+n/2)k} \quad (1b)$$

$$= \sum_{x=0}^{n/2-1} (a_x^{(left)} + z^{kn/2} a_x^{(right)}) z^{xk} \quad (1c)$$

where  $z = e^{\sigma i 2\pi/n}$ ,  $\sigma = \pm 1$  is the sign of the transform and  $k \in \{0, 1, \dots, n-1\}$ .

Here one has to distinguish the cases  $k$  even or odd, therefore we rewrite  $k \in \{0, 1, 2, \dots, n-1\}$  as  $k = 2j + \delta$  where  $j \in \{0, 2, \dots, \frac{n}{2} - 1\}$  and  $\delta \in \{0, 1\}$ :

$$\sum_{x=0}^{n-1} a_x z^{x(2j+\delta)} = \sum_{x=0}^{n/2-1} (a_x^{(left)} + z^{(2j+\delta)n/2} a_x^{(right)}) z^{x(2j+\delta)} \quad (2a)$$

$$= \begin{cases} \sum_{x=0}^{n/2-1} (a_x^{(left)} + a_x^{(right)}) z^{2xj} & \text{for } \delta = 0 \\ \sum_{x=0}^{n/2-1} z^x (a_x^{(left)} - a_x^{(right)}) z^{2xj} & \text{for } \delta = 1 \end{cases} \quad (2b)$$

$z^{(2j+\delta)n/2} = e^{\pm \pi i \delta}$  is equal to plus or minus one for  $\delta = 0$  or  $\delta = 1$  corresponding to  $k$  even or odd. The last two equations are, more compactly written, the key to the *radix-2 DIF FFT step*:

$$\mathcal{F}[a]^{(even)} \stackrel{n/2}{=} \mathcal{F}[a^{(left)} + a^{(right)}] \quad (3a)$$

$$\mathcal{F}[a]^{(odd)} \stackrel{n/2}{=} \mathcal{F}[\mathcal{S}^{1/2} (a^{(left)} - a^{(right)})] \quad (3b)$$

## The radix-2 DIT FFT algorithm

The following observation is the key to the (radix-2) *decimation in time* (DIT) FFT algorithm, also called *Cooley-Tukey* FFT algorithm:

For  $n$  even the  $k$ -th element of the Fourier transform is

$$\mathcal{F}[a]_k = \sum_{x=0}^{n-1} a_x z^{xk} = \sum_{x=0}^{n/2-1} a_{2x} z^{2xk} + \sum_{x=0}^{n/2-1} a_{2x+1} z^{(2x+1)k} \quad (1a)$$

$$= \sum_{x=0}^{n/2-1} a_{2x} z^{2xk} + z^k \sum_{x=0}^{n/2-1} a_{2x+1} z^{2xk} \quad (1b)$$

where  $z = e^{\sigma i 2\pi/n}$ ,  $\sigma = \pm 1$  is the sign of the transform and  $k \in \{0, 1, \dots, n-1\}$ .

The identity tells us how to compute the  $k$ -th element of the length- $n$  Fourier transform from the length- $n/2$  Fourier transforms of the even and odd indexed subsequences.

To actually rewrite the length- $n$  FT in terms of length- $n/2$  FTs one has to distinguish the cases  $0 \leq k < n/2$  and  $n/2 \leq k < n$ . In the expressions we rewrite  $k \in \{0, 1, 2, \dots, n-1\}$  as  $k = j + \delta \frac{n}{2}$  where  $j \in \{0, 1, \dots, n/2-1\}$  and  $\delta \in \{0, 1\}$ .

$$\sum_{x=0}^{n-1} a_x z^{x(j+\delta \frac{n}{2})} = \sum_{x=0}^{n/2-1} a_x^{(even)} z^{2x(j+\delta \frac{n}{2})} + z^{j+\delta \frac{n}{2}} \sum_{x=0}^{n/2-1} a_x^{(odd)} z^{2x(j+\delta \frac{n}{2})} \quad (2a)$$

$$= \begin{cases} \sum_{x=0}^{n/2-1} a_x^{(even)} z^{2xj} + z^j \sum_{x=0}^{n/2-1} a_x^{(odd)} z^{2xj} & \text{for } \delta = 0 \\ \sum_{x=0}^{n/2-1} a_x^{(even)} z^{2xj} - z^j \sum_{x=0}^{n/2-1} a_x^{(odd)} z^{2xj} & \text{for } \delta = 1 \end{cases} \quad (2b)$$

Observing that  $z^2$  is just the root of unity that appears in a length- $n/2$  transform one can rewrite the last two equations to obtain the *radix-2 DIT FFT step*:

$$\mathcal{F}[a]^{(left)} \stackrel{n/2}{=} \mathcal{F}[a^{(even)}] + \mathcal{S}^{1/2} \mathcal{F}[a^{(odd)}] \quad (3a)$$

$$\mathcal{F}[a]^{(right)} \stackrel{n/2}{=} \mathcal{F}[a^{(even)}] - \mathcal{S}^{1/2} \mathcal{F}[a^{(odd)}] \quad (3b)$$

## The radix-4 DIF FFT algorithm

The radix-2 DIF step in the new notation:

$$\begin{aligned}\mathcal{F}[a]^{(0\%2)} &\stackrel{n/2}{=} \mathcal{F}[\mathcal{S}^{0/2}(a^{(0/2)} + a^{(1/2)})] \\ \mathcal{F}[a]^{(1\%2)} &\stackrel{n/2}{=} \mathcal{F}[\mathcal{S}^{1/2}(a^{(0/2)} - a^{(1/2)})]\end{aligned}$$

The *radix-4 DIF FFT step*, applicable for  $n$  divisible by 4, is

$$\begin{aligned}\mathcal{F}[a]^{(0\%4)} &\stackrel{n/4}{=} \mathcal{F}[\mathcal{S}^{0/4}(a^{(0/4)} + a^{(1/4)} + a^{(2/4)} + a^{(3/4)})] \\ \mathcal{F}[a]^{(1\%4)} &\stackrel{n/4}{=} \mathcal{F}[\mathcal{S}^{1/4}(a^{(0/4)} + i\sigma a^{(1/4)} - a^{(2/4)} - i\sigma a^{(3/4)})] \\ \mathcal{F}[a]^{(2\%4)} &\stackrel{n/4}{=} \mathcal{F}[\mathcal{S}^{2/4}(a^{(0/4)} - a^{(1/4)} + a^{(2/4)} - a^{(3/4)})] \\ \mathcal{F}[a]^{(3\%4)} &\stackrel{n/4}{=} \mathcal{F}[\mathcal{S}^{3/4}(a^{(0/4)} - i\sigma a^{(1/4)} - a^{(2/4)} + i\sigma a^{(3/4)})]\end{aligned}$$

The radix-2 DIT step in the new notation:

$$\begin{aligned}\mathcal{F}[a]^{(0/2)} &\stackrel{n/2}{=} \mathcal{S}^{0/2}\mathcal{F}[a^{(0\%2)}] + \mathcal{S}^{1/2}\mathcal{F}[a^{(1\%2)}] \\ \mathcal{F}[a]^{(1/2)} &\stackrel{n/2}{=} \mathcal{S}^{0/2}\mathcal{F}[a^{(0\%2)}] - \mathcal{S}^{1/2}\mathcal{F}[a^{(1\%2)}]\end{aligned}$$

Note that  $\mathcal{S}^{0/2} = \mathcal{S}^0$  is the identity operator.

The *radix-4 DIT FFT step*:

$$\begin{aligned}\mathcal{F}[a]^{(0/4)} &\stackrel{n/4}{=} +\mathcal{S}^{0/4}\mathcal{F}[a^{(0\%4)}] + \mathcal{S}^{1/4}\mathcal{F}[a^{(1\%4)}] + \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}] + \mathcal{S}^{3/4}\mathcal{F}[a^{(3\%4)}] \\ \mathcal{F}[a]^{(1/4)} &\stackrel{n/4}{=} +\mathcal{S}^{0/4}\mathcal{F}[a^{(0\%4)}] + i\sigma\mathcal{S}^{1/4}\mathcal{F}[a^{(1\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}] - i\sigma\mathcal{S}^{3/4}\mathcal{F}[a^{(3\%4)}] \\ \mathcal{F}[a]^{(2/4)} &\stackrel{n/4}{=} +\mathcal{S}^{0/4}\mathcal{F}[a^{(0\%4)}] - \mathcal{S}^{1/4}\mathcal{F}[a^{(1\%4)}] + \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}] - \mathcal{S}^{3/4}\mathcal{F}[a^{(3\%4)}] \\ \mathcal{F}[a]^{(3/4)} &\stackrel{n/4}{=} +\mathcal{S}^{0/4}\mathcal{F}[a^{(0\%4)}] - i\sigma\mathcal{S}^{1/4}\mathcal{F}[a^{(1\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}] + i\sigma\mathcal{S}^{3/4}\mathcal{F}[a^{(3\%4)}]\end{aligned}$$

In contrast to the radix-2 step that happens to be identical for forward and backward transform the sign of the transform  $\sigma = \pm 1$  appears here.

## Split radix FFT algorithms

The idea underlying the *split radix FFT* is to use both radix-2 and radix-4 decompositions at the same time.

From the radix-2 (DIF) decomposition we use the first, the one for the even indices. For the odd indices we use the radix-4 splitting: The radix-4 decimation in frequency (DIF) step for the split radix FFT:

$$\begin{aligned}\mathcal{F}[a]^{(0\%2)} &\stackrel{n/2}{=} \mathcal{F}\left[a^{(0/2)} + a^{(1/2)}\right] \\ \mathcal{F}[a]^{(1\%4)} &\stackrel{n/4}{=} \mathcal{F}\left[\mathcal{S}^{1/4}\left(a^{(0/4)} - a^{(2/4)}\right) + i\sigma\left(a^{(1/4)} - a^{(3/4)}\right)\right] \\ \mathcal{F}[a]^{(3\%4)} &\stackrel{n/4}{=} \mathcal{F}\left[\mathcal{S}^{3/4}\left(a^{(0/4)} - a^{(2/4)}\right) - i\sigma\left(a^{(1/4)} - a^{(3/4)}\right)\right]\end{aligned}$$

Now we have expressed the length- $N = 2^n$  FFT as one length- $N/2$  and two length- $N/4$  FFTs. Note that  $\mathcal{S}^{3/4} = \mathcal{S}^{-1/4}$  which means a saving in the trigonometric computations. The nice feature is that the operation count of the split radix FFT is actually lower than that of the radix-4 FFT. Using the introduced notation it is almost trivial to write down the DIT version of the algorithm:

The radix-4 decimation in time (DIT) step for the split radix FFT:

$$\begin{aligned}\mathcal{F}[a]^{(0/2)} &\stackrel{n/2}{=} \left(\mathcal{F}[a^{(0\%2)}] + \mathcal{S}^{1/2}\mathcal{F}[a^{(1\%2)}]\right) \\ \mathcal{F}[a]^{(1/4)} &\stackrel{n/4}{=} \left(\mathcal{F}[a^{(0\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}]\right) + i\sigma\mathcal{S}^{1/4}\left(\mathcal{F}[a^{(1\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(3\%4)}]\right) \\ \mathcal{F}[a]^{(3/4)} &\stackrel{n/4}{=} \left(\mathcal{F}[a^{(0\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(2\%4)}]\right) - i\sigma\mathcal{S}^{1/4}\left(\mathcal{F}[a^{(1\%4)}] - \mathcal{S}^{2/4}\mathcal{F}[a^{(3\%4)}]\right)\end{aligned}$$

## Radix- $\sqrt{n}$ FFT algorithms

The *matrix Fourier algorithm* (MFA) for the FFT:

1. Apply a (length  $R$ ) FFT on each column.
2. Multiply each matrix element (index  $r, c$ ) by  $\exp(\sigma 2 \pi i r c / n)$
3. Apply a (length  $C$ ) FFT on each row.
4. Transpose the matrix.

Note the elegance!

MFA = '*four step FFT*'

A trivial modification is obtained if the steps are executed in reversed order.

The *transposed matrix Fourier algorithm* (TMFA) for the FFT:

1. Transpose the matrix.
2. Apply a (length  $C$ ) FFT on each row of the transposed matrix.
3. Multiply each matrix element (index  $r, c$ ) by  $\exp(\sigma 2 \pi i r c / n)$ .
4. Apply a (length  $R$ ) FFT on each column of the transposed matrix.

## Cyclic convolution

The *cyclic convolution* (or *circular convolution*) of two length- $n$  sequences  $a = [a_0, a_1, \dots, a_{n-1}]$  and  $b = [b_0, b_1, \dots, b_{n-1}]$  is defined as the length- $n$  sequence  $h$  with elements  $h_\tau$  as:

$$h = a \circledast b \quad (1a)$$

$$h_\tau := \sum_{x+y \equiv \tau \pmod{n}} a_x b_y \quad (1b)$$

The last equation may be rewritten as

$$h_\tau := \sum_{x=0}^{n-1} a_x b_{(\tau-x) \bmod n} \quad (2)$$

That is, indices  $\tau - x$  wrap around, it is a cyclic convolution.

Pseudo code to compute the cyclic convolution of  $\mathbf{a}[]$  with  $\mathbf{b}[]$  using the definition, the result is returned in  $\mathbf{c}[]$ :

```
procedure convolution(a[],b[],c[],n)
{
  for tau:=0 to n-1
  {
    s := 0
    for x:=0 to n-1
    {
      tx := tau - x
      if tx<0 then tx := tx + n // modulo reduction
      s := s + a[x] * b[tx]
    }
    c[tau] := s
  }
}
```

For length- $n$  sequences this procedure involves proportional  $n^2$  operations, therefore it is slow for large values of  $n$ . The Fourier transform provides us with a more efficient way to compute convolutions that only uses proportional  $n \log(n)$  operations.

## Fast computation of convolutions

First we have to state the convolution property of the Fourier transform:

$$\mathcal{F}[a \circledast b] = \mathcal{F}[a] \mathcal{F}[b] \quad (1)$$

That is, convolution in original space is element-wise multiplication in Fourier space. The statement can be motivated as follows:

$$\mathcal{F}[a]_k \mathcal{F}[b]_k = \sum_x a_x z^{kx} \sum_y b_y z^{ky} \quad (2a)$$

$$\text{with } y := \tau - x \quad (2b)$$

$$= \sum_x a_x z^{kx} \sum_{\tau-x} b_{\tau-x} z^{k(\tau-x)} \quad (2c)$$

$$= \sum_x \sum_{\tau-x} a_x z^{kx} b_{\tau-x} z^{k(\tau-x)} \quad (2d)$$

$$= \sum_{\tau} \left( \sum_x a_x b_{\tau-x} \right) z^{k\tau} \quad (2e)$$

$$= \left( \mathcal{F} \left[ \sum_x a_x b_{\tau-x} \right] \right)_k \quad (2f)$$

$$= (\mathcal{F}[a \circledast b])_k \quad (2g)$$

Rewriting relation 1 as

$$a \circledast b = \mathcal{F}^{-1}[\mathcal{F}[a] \mathcal{F}[b]] \quad (3)$$

tells us how to proceed.



## Acyclic (linear) convolution

In the definition of the cyclic convolution one can distinguish between those summands where the  $x + y$  ‘wrapped around’ (i.e.  $x + y = n + \tau$ ) and those where simply  $x + y = \tau$  holds. These are denoted by  $h^{(1)}$  and  $h^{(0)}$  respectively. We have

$$h = h^{(0)} + h^{(1)} \quad \text{where} \quad (1a)$$

$$h^{(0)} = \sum_{x \leq \tau} a_x b_{\tau-x} \quad (1b)$$

$$h^{(1)} = \sum_{x > \tau} a_x b_{n+\tau-x} \quad (1c)$$

There is a simple way to separate  $h^{(0)}$  and  $h^{(1)}$  as the left and right half of a length- $2n$  sequence. This is just what the *acyclic convolution* (or *linear convolution*) does: Acyclic convolution of two (length- $n$ ) sequences  $a$  and  $b$  can be defined as that length- $2n$  sequence  $h$  which is the cyclic convolution of the *zero padded* sequences  $A$  and  $B$ :

$$A := [a_0, a_1, a_2, \dots, a_{n-1}, 0, 0, \dots, 0] \quad (2)$$

Same for  $B$ . Then the acyclic convolution is defined as

$$h = a \otimes_{ac} b \quad (3a)$$

$$h_\tau := \sum_{x=0}^{2n-1} A_x B_{\tau-x} \quad \tau = 0, 1, 2, \dots, 2n-1 \quad (3b)$$

As an illustration consider the convolution of the sequence  $[1, 1, 1, 1]$  with itself: its linear self convolution is the length-8 sequence  $[h_0][h_1] = [1, 2, 3, 4][3, 2, 1, 0]$ , its cyclic self convolution is  $[h_0 + h_1] = [4, 4, 4, 4]$ .

## Convolutions: semi-symbolic tables

A convenient way to illustrate the cyclic convolution of two sequences is the following semi-symbolical table:

| +-  | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 |
|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0:  | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 |
| 1:  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  |
| 2:  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  |
| 3:  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  |
| 4:  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  |
| 5:  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  |
| 6:  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  |
| 7:  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  |
| 8:  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  |
| 9:  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  |
| 10: | 10 | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
| 11: | 11 | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 |
| 12: | 12 | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 |
| 13: | 13 | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 |
| 14: | 14 | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 |
| 15: | 15 | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 |

The entries denote where in the convolution the products of the input elements can be found:

| +- | 0   | 1 | 2 | 3                            | ... |
|----|-----|---|---|------------------------------|-----|
| 0: | 0   | 1 | 2 | 4                            | ... |
| 1: | 1   | 3 | 5 | <--- h[5] contains a[2]*b[1] |     |
| 2: | 4   | 8 | 9 | <--- h[9] contains a[2]*b[b] |     |
| 3: | ... |   |   |                              |     |

Acyclic convolution (where there are 32 buckets 0..31) looks like:

| +-  | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 |
|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0:  | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 |
| 1:  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| 2:  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
| 3:  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| 4:  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 5:  | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| 6:  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| 7:  | 7  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 |
| 8:  | 8  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 9:  | 9  | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| 10: | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| 11: | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 12: | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 |
| 13: | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 |
| 14: | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 |
| 15: | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 |

## Number theoretic transforms (NTTs)

How to make a number theoretic transform out of your FFT:

‘Replace  $\exp(\pm 2\pi i/n)$  by a primitive  $n$ -th root of unity, done.’

We want to implement FFTs in  $\mathbb{Z}/m\mathbb{Z}$  (the ring of integers modulo some integer  $m$ ) instead of  $\mathbb{C}$ , the (field of the) complex numbers. These FFTs are called *number theoretic transforms* (NTTs), mod  $m$  FFTs or (if  $m$  is a prime) prime modulus transforms.

There is a restriction for the choice of  $m$ : For a length  $n$  NTT we need a primitive  $n$ -th root of unity. A number  $r$  is called an  $n$ -th root of unity if  $r^n = 1$ . It is called a *primitive  $n$ -th root* if  $r^k \neq 1 \forall k < n$ .

In  $\mathbb{C}$  matters are simple:  $e^{\pm 2\pi i/n}$  is a primitive  $n$ -th root of unity for arbitrary  $n$ .  $e^{2\pi i/21}$  is a 21-th root of unity.  $r = e^{2\pi i/3}$  is also 21-th root of unity but not a primitive root, because  $r^3 = 1$ . A primitive  $n$ -th root of 1 in  $\mathbb{Z}/m\mathbb{Z}$  is also called an *element of order  $n$* . The ‘cyclic’ property of the elements  $r$  of order  $n$  lies in the heart of all FFT algorithms:  $r^{n+k} = r^k$ .

In  $\mathbb{Z}/m\mathbb{Z}$  things are not that simple: for a given modulus  $m$  primitive  $n$ -th roots of unity do not exist for arbitrary  $n$ . They exist for some maximal order  $R$  only. Roots of unity of an order different from  $R$  are available only for the divisors  $d_i$  of  $R$ :  $r^{R/d_i}$  is a  $d_i$ -th root of unity because  $(r^{R/d_i})^{d_i} = r^R = 1$ .

Therefore  $n$  must divide  $R$ , the first condition for NTTs:

$$n \setminus R \iff \exists \sqrt[n]{1} \quad (1)$$

The operations needed in FFTs are addition, subtraction and multiplication. Division is not needed, except for division by  $n$  for the final normalization after transform and back-transform. Division by  $n$  is multiplication by the inverse of  $n$ . Hence  $n$  must be invertible in  $\mathbb{Z}/m\mathbb{Z}$ :  $n$  must be co-prime to  $m$  (i.e.  $\gcd(n, m) = 1$ ), the second condition for NTTs:

$$n \perp m \iff \exists n^{-1} \text{ in } \mathbb{Z}/m\mathbb{Z} \quad (2)$$

## Prime modulus

If the modulus is a prime  $p$  then  $\mathbb{Z}/p\mathbb{Z}$  is the field  $\mathbb{F}_p = GF(p)$ : All elements except 0 have inverses and ‘division is possible’ in  $\mathbb{Z}/p\mathbb{Z}$ . Thereby the second condition is trivially fulfilled for all FFT lengths  $n < p$ : a prime  $p$  is coprime to all integers  $n < p$ .

Roots of unity are available for the maximal order  $R = p - 1$  and its divisors: Therefore the first condition on  $n$  for a length- $n$  mod  $p$  FFT being possible is that  $n$  divides  $p - 1$ . This restricts the choice for  $p$  to primes of the form  $p = v n + 1$ : For length- $n = 2^k$  FFTs one will use primes like  $p = 3 \cdot 5 \cdot 2^{27} + 1$  (31 bits),  $p = 13 \cdot 2^{28} + 1$  (32 bits),  $p = 3 \cdot 29 \cdot 2^{56} + 1$  (63 bits) or  $p = 27 \cdot 2^{59} + 1$  (64 bits). Primes of that form are not ‘exceptional’. The elements of maximal order in  $\mathbb{Z}/p\mathbb{Z}$  are called *primitive elements*, *generators* or *primitive roots* modulo  $p$ . If  $r$  is a generator, then every element in  $\mathbb{F}_p$  different from 0 is equal to some power  $r^e$  ( $1 \leq e < p$ ) of  $r$  and its order is  $R/e$ . To test whether  $r$  is a primitive  $n$ -th root of unity in  $\mathbb{F}_p$  one does not need to check  $r^k \neq 1$  for all  $k < n$ . It suffices to do the check for exponents  $k$  that are prime factors of  $n$ . This is because the order of any element divides the maximal order.

To find a primitive root in  $\mathbb{F}_p$  proceed as indicated by the following pseudo code:

```
function primroot(p)
{
    if p==2 then return 1
    f[] := distinct_prime_factors(p-1)
    for r:=2 to p-1
    {
        x := TRUE
        foreach q in f[]
        {
            if r**((p-1)/q)==1 then x:=FALSE
        }
        if x==TRUE then return r
    }
    error("no primitive root found") // p cannot be prime !
}
```

The algorithm is a simple search and might seem ineffective. In practice the root is found after only several tries.

An element of order  $n$  in  $\mathbb{F}_p$  is returned by this function:

```
function element_of_order(n,p)
{
    R := p-1 // maxorder
    if (R/n)*n != R then error("order n must divide maxorder p-1")
    r := primroot(p)
    x := r**(R/n)
    return x
}
```

## Multiplication and polynomial products

A number written in radix  $r$  as

$$a_P \ a_{P-1} \ \dots \ a_2 \ a_1 \ a_0 \cdot a_{-1} \ a_{-2} \ \dots \ a_{-p+1} \ a_{-p} \quad (1)$$

denotes a quantity of

$$\sum_{i=-p}^P a_i \cdot r^i = a_P \cdot r^P + a_{P-1} \cdot r^{P-1} + \dots + a_{-p} \cdot r^{-p}. \quad (2)$$

That means, the digits can be considered as coefficients of a polynomial in  $r$ . For example, with decimal numbers one has  $r = 10$  and  $123.4 = 1 \cdot 10^2 + 2 \cdot 10^1 + 3 \cdot 10^0 + 4 \cdot 10^{-1}$ . The product of two numbers is almost the polynomial product

$$\sum_{k=0}^{2N-2} c_k r^k := \sum_{i=0}^{N-1} a_i r^i \cdot \sum_{j=0}^{N-1} b_j r^j \quad (3)$$

The  $c_k$  are found by comparing coefficients. One easily checks that the  $c_k$  must satisfy the convolution equation.

As the  $c_k$  can be greater than ‘nine’ (that is,  $r - 1$ ), the result has to be ‘fixed’ using *carry* operations: Go from right to left, replace  $c_k$  by  $c_k \% r$  and add  $(c_k - c_k \% r)/r$  to its left neighbor.

An example: usually one would multiply the numbers 82 and 34 as follows:

$$\begin{array}{r} 82 \times 34 \\ \hline 3 \ 32 \ 8 \\ 2 \ 24 \ 6 \\ \hline = 2 \ 7 \ 8 \ 8 \end{array}$$

We just said that the carries can be delayed to the end of the computation:

$$\begin{array}{r} 82 \times 34 \\ \hline 32 \ 8 \\ 24 \ 6 \\ \hline 24 \ 38 \ 8 \\ \hline = 2 \ 2 \ 7 \ 3 \ 8 \ 8 \end{array}$$

... which is really polynomial multiplication (which in turn is a convolution of the coefficients):

$$\begin{array}{r} (8x + 2) \times (3x + 4) \\ \hline 32x \quad 8 \\ 24x^2 \quad 6x \\ \hline = 24x^2 + 38x + 8 \end{array}$$

## Multiplication via FFTs

Convolution can be done efficiently using the Fast Fourier Transform (FFT): Convolution is a simple (element-wise) multiplication in Fourier space. The FFT itself takes  $N \cdot \log N$  operations. Instead of the direct convolution ( $\sim N^2$ ) one proceeds like this:

- compute the FFTs of multiplicand and multiplier
- multiply the transformed sequences elementwise
- compute inverse transform of the product

To understand why this actually works note that (1) the multiplication of two polynomials can be achieved by the (more complicated) scheme:

- evaluate both polynomials at sufficiently many points (at least one more point than the degree of the product polynomial  $c$ :  $\deg c = \deg a + \deg b$ )
- element-wise multiply the values found
- find the polynomial corresponding to those (product-)values

and (2) that *the FFT is an algorithm for the parallel evaluation of a given polynomial at many points, namely the roots of unity*. (3) the inverse FFT is an algorithm to find (the coefficients of) a polynomial whose values are given at the roots of unity.

## FFT multiplication: example

Relaunching our example we use the fourth roots of unity  $\pm 1$  and  $\pm i$ :

| $a = (8x + 2) \quad \times \quad b = (3x + 4)$ |         |         | $c = ab$                |
|--|---------|---------|-------------------------|
| +1   | +10     | +7      | +70                     |
| +i   | +8i + 2 | +3i + 4 | +38i - 16               |
| -1   | -6      | +1      | -6                      |
| -i   | -8i + 2 | -3i + 4 | -38i - 16               |
|  |         |         | $c = (24x^2 + 38x + 8)$ |

This table has to be read like this: first the given polynomials  $a$  and  $b$  are evaluated at the points given in the left column, thereby the columns below  $a$  and  $b$  are filled. Then the values are multiplied to fill the column below  $c$ , giving the values of  $c$  at the points. Finally, the actual polynomial  $c$  is found from those values, resulting in the lower right entry. You may find it instructive to verify that a 4-point FFT really evaluates  $a$ ,  $b$  by transforming the sequences 0, 0, 8, 2 and 0, 0, 3, 4 by hand. The backward transform of 70,  $38i - 16$ ,  $-6$ ,  $-38i - 16$  should produce the final result given for  $c$ .

The operation count is dominated by that of the FFTs (the elementwise multiplication is of course  $\sim N$ ), so the whole fast convolution algorithm takes  $\sim N \cdot \log N$  operations. The following carry operation is also  $\sim N$  and can therefore be neglected when counting operations.

## Radix and precision with FFT multiplication

Restrictions are due to the fact that the components of the convolution must be representable as integer numbers with the data type used for the FFTs: The cumulative sums have to be represented precisely enough to distinguish every (integer) quantity from the next bigger (or smaller) value. The highest possible value for a will appear in the middle of the product and when multiplicand and multiplier consist of ‘nines’ (that is  $R-1$ ) only. For radix  $R$  and a precision of  $N$  LIMBs Let the maximal possible value be  $C$ , then

$$C = N(R-1)^2 \quad (1)$$

The number of bits to represent  $C$  exactly is the integer greater or equal to

$$\log_2(N(R-1)^2) = \log_2 N + 2 \log_2(R-1) \quad (2)$$

Due to numerical errors there must be a few more bits for safety. If computations are made using double-precision floating point numbers (C-type `double`) one typically has a mantissa of 53 bits. then we need to have

$$M \geq \log_2 N + 2 \log_2(R-1) + S \quad (3)$$

where  $M :=$ mantissa-bits and  $S :=$ safety-bits. Using  $\log_2(R-1) < \log_2(R)$ :

$$N_{max}(R) = 2^{M-S-2 \log_2(R)} \quad (4)$$

| Radix $R$        | max # LIMBs   | max # hex digits | max # bits |
|------------------|---------------|------------------|------------|
| $2^{10} = 1024$  | 1048, 576 $k$ | 2621, 440 $k$    | 10240 $M$  |
| $2^{11} = 2048$  | 262, 144 $k$  | 720, 896 $k$     | 2816 $M$   |
| $2^{12} = 4096$  | 65, 536 $k$   | 196, 608 $k$     | 768 $M$    |
| $2^{13} = 8192$  | 16384 $k$     | 53, 248 $k$      | 208 $M$    |
| $2^{14} = 16384$ | 4096 $k$      | 14, 336 $k$      | 56 $M$     |
| $2^{15} = 32768$ | 1024 $k$      | 3840 $k$         | 15 $M$     |
| $2^{16} = 65536$ | 256 $k$       | 1024 $k$         | 4 $M$      |

For decimal numbers:

| Radix $R$ | max # LIMBs | max # digits | max # bits |
|-----------|-------------|--------------|------------|
| $10^2$    | 110 $G$     | 220 $G$      | 730 $G$    |
| $10^3$    | 1100 $M$    | 3300 $M$     | 11 $G$     |
| $10^4$    | 11 $M$      | 44 $M$       | 146 $M$    |
| $10^5$    | 110 $k$     | 550 $k$      | 1826 $k$   |
| $10^6$    | 1 $k$       | 6, 597       | 22 $k$     |
| $10^7$    | 11          | 77           | 255        |

Do the sum of digits test!



## Division: Inversion

The ordinary division algorithm is far too expensive for numbers of extreme precision. Instead one replaces the division  $\frac{a}{b}$  by the multiplication of  $a$  with the inverse of  $b$ . The inverse of  $b$  is computed by finding a starting approximation  $x_0 \approx \frac{1}{b}$  and then iterating

$$x_{k+1} = x_k + x_k (1 - b x_k) \quad (1)$$

until the desired precision is reached. The convergence is quadratic (second order), which means that the number of correct digits is doubled with each step: if  $x_k = \frac{1}{b}(1 + e)$  then  $x_{k+1} = \frac{1}{b} (1 - e^2)$ .

Moreover, each only requires computations with twice the number of digits that were correct at its beginning. Still better: the multiplication  $x_k(\dots)$  needs only to be done with half of the current precision as it computes the correcting digits (which alter only the less significant half of the digits). Thus, at each step we have 1.5 multiplications of the current precision. The total work<sup>1</sup> amounts to  $1.5 \cdot \sum_{n=0}^N \frac{1}{2^n}$  which is less than 3 full precision multiplications. Together with the final multiplication a division costs as much as 4 multiplications. Another nice feature of the algorithm is that it is self-correcting. The following numerical example shows the first two steps of the computation of an inverse starting from a two-digit initial approximation:

$$b := 3.1415926 \quad (2)$$

$$x_0 = 0.31 \quad \text{initial 2 digit approximation for } 1/b \quad (3)$$

$$b \cdot x_0 = 3.141 \cdot 0.3100 = 0.9737 \quad (4)$$

$$y_0 := 1.000 - b \cdot x_0 = 0.02629 \quad (5)$$

$$x_0 \cdot y_0 = 0.3100 \cdot 0.02629 = 0.0081(49) \quad (6)$$

$$x_1 := x_0 + x_0 \cdot y_0 = 0.3100 + 0.0081 = 0.3181 \quad (7)$$

$$b \cdot x_1 = 3.1415926 \cdot 0.31810000 = 0.9993406 \quad (8)$$

$$y_1 := 1.0000000 - b \cdot x_1 = 0.0006594 \quad (9)$$

$$x_1 \cdot y_1 = 0.31810000 \cdot 0.0006594 = 0.0002097(5500) \quad (10)$$

$$x_2 := x_1 + x_1 \cdot y_1 = 0.31810000 + 0.0002097 = 0.31830975 \quad (11)$$

---

<sup>1</sup>The asymptotics of the multiplication is set to  $\sim N$  (instead of  $N \log(N)$ ) for the estimates made here, this gives a realistic picture for large  $N$ .

## Root extraction

Computation of square roots can be done using a similar scheme: first compute  $\frac{1}{\sqrt{d}}$  then a final multiply with  $d$  gives  $\sqrt{d}$ . Find a starting approximation  $x_0 \approx \frac{1}{\sqrt{b}}$  then iterate

$$x_{k+1} = x_k + x_k \frac{(1 - d x_k^2)}{2} \quad (1)$$

until the desired precision is reached. Convergence is again 2nd order: if  $x_k = \frac{1}{\sqrt{b}}(1 + e)$  then

$$x_{k+1} = \frac{1}{\sqrt{b}} \left( 1 - \frac{3}{2}e^2 - \frac{1}{2}e^3 \right) \quad (2)$$

Similar considerations as above (with squaring considered as expensive as multiplication<sup>2</sup>) give an operation count of 4 multiplications for  $\frac{1}{\sqrt{d}}$  or 5 for  $\sqrt{d}$ .

Note that this algorithm is considerably better than the one where  $x_{k+1} := \frac{1}{2}(x_k + \frac{d}{x_k})$  is used as iteration, because no long divisions are involved.

In `hfloat`, when the achieved precision is below a certain limit a third order correction is used to assure maximum precision at the last step:

$$x_{k+1} = x_k + x_k \frac{(1 - d x_k^2)}{2} + x_k \frac{3(1 - d x_k^2)^2}{8} \quad (3)$$

---

<sup>2</sup>Indeed it costs about  $\frac{2}{3}$  of a multiplication.

## Inverse $n$ -th root, a general expression

There is a nice general formula that allows to build iterations with arbitrary order of convergence for  $\sqrt[a]{d} = d^{-1/a}$  that involve no long division.

One uses the identity

$$d^{-1/a} = x (1 - (1 - x^a d))^{-1/a} \quad (1)$$

$$= x (1 - y)^{-1/a} \quad \text{where} \quad y := (1 - x^a d) \quad (2)$$

Taylor expansion gives

$$d^{-1/a} = x \sum_{k=0}^{\infty} (1/a)^{\bar{k}} y^k \quad (3)$$

where  $z^{\bar{k}} := z(z+1)(z+2)\dots(z+k-1)$  (and  $z^{\bar{0}} = 1$ ). Written out:

$$\begin{aligned} d^{-1/a} = x \frac{1}{\sqrt[a]{1-y}} = x \left( 1 + \frac{y}{a} + \frac{(1+a)y^2}{2a^2} + \frac{(1+a)(1+2a)y^3}{6a^3} + \right. \\ \left. + \frac{(1+a)(1+2a)(1+3a)y^4}{24a^4} + \dots + \frac{\prod_{k=1}^{n-1} (1+ka)}{n! a^n} y^n + \dots \right) \end{aligned} \quad (4)$$

A  $n$ -th order iteration for  $d^{-1/a}$  is obtained by truncating the above series after the  $(n-1)$ -th term:

$$\Phi_n(x) := x \sum_{k=0}^{n-1} (1/a)^{\bar{k}} y^k \quad (5)$$

$$x_{k+1} = \Phi_n(x_k) \quad (6)$$

Convergence is  $n$ -th order:

$$\Phi_n(d^{-1/a}(1+e)) = d^{-1/a}(1+O(e^n)) \quad (7)$$

## Iterations for the inversion of a function

An *iteration* for a zero  $r$  (or root,  $f(r) = 0$ ) of a function  $f(x)$  are themselves functions  $\Phi(x)$  that, when ‘used’ like

$$x_{k+1} = \Phi(x_k) \quad (1)$$

will make  $x_k$  converge towards the root:  $x_\infty = r$ . Convergence is subject to the condition that  $x_0$  was chosen not too far away from  $r$ . The function  $\Phi(x)$  must (and can) be constructed so that it has an attracting fixed point where  $f(x)$  has a zero:

$$\Phi(r) = r \quad (\text{fixed point}) \quad (2)$$

$$|\Phi'(r)| < 1 \quad (\text{attracting}) \quad (3)$$

This type of iteration is a so-called *one-point iteration*. There are also *multi-point iterations*, these are of the form  $x_{k+1} = \Phi(x_k, x_{k-1}, \dots, x_{k-j}), j \geq 1$ . The best known example is the two-point iteration

$$\Phi(x_k, x_{k-1}) = \frac{x_{k-1} f(x_k) - x_k f(x_{k-1})}{f(x_k) - f(x_{k-1})} \quad (4)$$

We are mainly concerned with one-point iterations in what follows.

Order of convergence: linear vs. *super-linear*.

The number of correct digits grows exponentially (to the base  $n$ ) at each step. Iterations of second order ( $n = 2$ ) are often called *quadratic* (or *quadratically convergent*), those of third order *cubic* iterations. Fourth, fifth and sixth order iterations are called *quartic*, *quintic* and *sextic* and so on.

For  $n \geq 2$  the function  $\Phi$  has a *super-attracting fixed point* at  $r$ :  $\Phi'(r) = 0$ . For an iteration of order  $n$  one has

$$\Phi'(r) = 0, \quad \Phi''(r) = 0, \quad \dots, \quad \Phi^{(n-1)}(r) = 0 \quad (5)$$

There seems to be no standard term for emphasizing the number of derivatives vanishing at the fixed point: *attracting of order  $n$*  might be appropriate.

## Schröder's formula

Let  $n \geq 2$  then the expression

$$S_n(x) := x + \sum_{t=1}^{n-1} (-1)^t \frac{f(x)^t}{t!} \left( \frac{1}{f'(x)} \partial \right)^{t-1} \frac{1}{f'(x)} \quad (1)$$

gives a  $n$ -th order iteration for a (simple) root  $r$  of  $f$ . This is, explicitly,

$$\begin{aligned} S = x & - \frac{f}{1! f'} - \frac{f^2}{2! f'^3} \cdot f'' - \frac{f^3}{3! f'^5} \cdot (3f''^2 - f' f''') \\ & - \frac{f^4}{4! f'^7} \cdot (15f''^3 - 10f' f'' f''' + f'^2 f'''' ) \\ & - \frac{f^5}{5! f'^9} \cdot (105f''^4 - 105f' f''^2 f''' + 10f'^2 f''''^2 + 15f'^2 f'' f'''' - f'^3 f''''') - \dots \end{aligned} \quad (2)$$

The third order iteration obtained upon truncation after the third term on the right hand side, written as

$$S_3 = x - \frac{f}{f'} \left( 1 + \frac{f f''}{2 f'^2} \right) \quad (3)$$

is sometimes referred to as 'Householder's method'. Approximating the second term on the rhs. as  $\frac{f}{f'} \left( 1 - \frac{f f''}{2 f'^2} \right)^{-1}$  gives Halley's formula.

Cite Schroeder (p.16, translation has a typo in the first formula):

If we denote the general term by

$$- \frac{f^a}{a!} \frac{\chi_a}{f'^{2a-1}} \quad (4a)$$

the numbers  $\chi_a$  can be easily computed by the recurrence

$$\chi_{a+1} = (2a-1)f''\chi_a - f'\partial\chi_a \quad (4b)$$

.

## Householder's formula

For  $n \geq 2$  the expression

$$H_n(x) := x + (n-1) \frac{\left(\frac{g(x)}{f(x)}\right)^{(n-2)}}{\left(\frac{g(x)}{f(x)}\right)^{(n-1)}} \quad (1)$$

gives a  $n$ -th order iteration for a (simple) root  $r$  of  $f$ . The function  $g(x)$  must be analytic near the root and is set to 1 in what follows.

$$H_2(x) = x - \frac{f}{f'} \quad (2a)$$

$$H_3(x) = x - \frac{2ff'}{2f'^2 - ff''} \quad (2b)$$

$$H_4(x) = x - \frac{3f(ff'' - 2f'^2)}{6ff'f'' - 6f'^3 - f^2f'''} \quad (2c)$$

$$H_5(x) = x + \frac{4f(6f'^3 - 6ff'f'' + f^2f''')}{f^3f'''' - 24f'^4 + 36ff'^2f'' - 8f^2f'f''' - 6f^2f''^2} \quad (2d)$$

The second order variant is Newton's formula, the third order iteration is called Halley's formula.

The well-known derivation of Halley's formula by applying Newton's formula to  $f/\sqrt{f'}$  can be generalized to produce  $m$ -order iterations as follows: Let  $F_1(x) = f(x)$  and for  $m \geq 2$  let

$$F_m(x) = \frac{F_{m-1}(x)}{F'_{m-1}(x)^{1/m}} \quad (3a)$$

$$H_m(x) = x - \frac{F_{m-1}(x)}{F'_{m-1}(x)} \quad (3b)$$

An alternative recursive formulation:

$$Q_2(x) = 1 \quad (4a)$$

$$Q_{m+1} = f'(x) Q_m(x) - \frac{1}{m-1} f(x) Q'_m(x) \quad (4b)$$

$$H_m = x - f(x) \frac{Q_m(x)}{Q_{m-1}(x)} \quad (4c)$$

## The AGM

The AGM (arithmetic geometric mean) plays a central role in the high precision computation of logarithms and  $\pi$ .

The  $AGM(a, b)$  is defined as the limit of the iteration

$$a_{k+1} = \frac{a_k + b_k}{2} \quad (1a)$$

$$b_{k+1} = \sqrt{a_k b_k} \quad (1b)$$

starting with  $a_0 = a$  and  $b_0 = b$ . Both of the values converge quadratically to a common limit. The related quantity  $c_k$  (used in many AGM based computations) is defined as

$$c_k^2 = a_k^2 - b_k^2 \quad (2)$$

$$= (a_{k-1} - a_k)^2 \quad (3)$$

One further defines

$$R'(k) := 1 - \frac{1}{2} \sum_{n=0}^{\infty} 2^n c_n^2 \quad (4)$$

corresponding to  $AGM(1, k)$ , that is,  $a_0 = 1, b_0 = k, c_0 = \sqrt{1 - k^2}$ .

It can be shown that

$$F\left(\begin{matrix} \frac{1}{2}, \frac{1}{2} \\ 1 \end{matrix} \middle| 1 - \frac{b^2}{a^2}\right) = \frac{a}{AGM(a, b)} \quad (5)$$

$$F\left(\begin{matrix} \frac{1}{2}, \frac{1}{2} \\ 1 \end{matrix} \middle| x\right) = \frac{1}{AGM(1, \sqrt{1 - x})} \quad (6)$$

An alternative way for the computation for the AGM iteration is

$$c_{k+1} = \frac{a_k - b_k}{2} \quad (7a)$$

$$a_{k+1} = \frac{a_k + b_k}{2} \quad (7b)$$

$$b_{k+1} = \sqrt{a_{k+1}^2 - c_{k+1}^2} \quad (7c)$$

## The AGM, Schönhage's variant

Schönhage gives the most economic variant of the AGM, which, apart from the square root, only needs one squaring per step:

$$A_0 = a_0^2 \tag{1a}$$

$$B_0 = b_0^2 \tag{1b}$$

$$t_0 = 1 - (A_0 - B_0) \tag{1c}$$

$$S_k = \frac{A_k + B_k}{4} \tag{1d}$$

$$b_k = \sqrt{B_k} \quad [\text{square root}] \tag{1e}$$

$$a_{k+1} = \frac{a_k + b_k}{2} \tag{1f}$$

$$A_{k+1} = a_{k+1}^2 \quad [\text{squaring}] \tag{1g}$$

$$= \left( \frac{\sqrt{A_k} + \sqrt{B_k}}{2} \right)^2 = \frac{A_k + B_k}{4} + \frac{\sqrt{A_k B_k}}{2} \tag{1h}$$

$$B_{k+1} = 2(A_{k+1} - S_k) = b_{k+1}^2 \tag{1i}$$

$$c_{k+1}^2 = A_{k+1} - B_{k+1} = a_{k+1}^2 - b_{k+1}^2 \tag{1j}$$

$$t_{k+1} = t_k - 2^{k+1} c_{k+1}^2 \tag{1k}$$

Starting with  $a_0 = A_0 = 1$ ,  $B_0 = 1/2$  one has  $\pi \approx (2 a_n^2)/t_n$ .



## Superlinear iterations for $\pi$

The number of full precision multiplications (FPM) are an indication of the efficiency of the algorithm. The approximate number of FPMs that were counted with a computation of  $\pi$  to 4 million decimal digits<sup>3</sup> is indicated like this: #FPM=123.4.

AGM as in [hfloat: src/pi/piagm.cc], #FPM=98.4 (#FPM=149.3 for the quartic variant):

$$a_0 = 1 \quad (1a)$$

$$b_0 = \frac{1}{\sqrt{2}} \quad (1b)$$

$$p_n = \frac{2 a_{n+1}^2}{1 - \sum_{k=0}^n 2^k c_k^2} \rightarrow \pi \quad (1c)$$

$$\pi - p_n = \frac{\pi^2 2^{n+4} e^{-\pi 2^{n+1}}}{AGM^2(a_0, b_0)} \quad (1d)$$

Borwein's quartic (fourth order) iteration, variant  $r = 4$  as in [hfloat: src/pi/pi4th.cc], #FPM=170.5:

$$y_0 = \sqrt{2} - 1 \quad (2a)$$

$$a_0 = 6 - 4\sqrt{2} \quad (2b)$$

$$y_{k+1} = \frac{1 - (1 - y_k^4)^{1/4}}{1 + (1 - y_k^4)^{1/4}} \rightarrow 0 + \quad (2c)$$

$$= \frac{(1 - y_k^4)^{-1/4} - 1}{(1 - y_k^4)^{-1/4} + 1} \quad (2d)$$

$$a_{k+1} = a_k (1 + y_{k+1})^4 - 2^{2k+3} y_{k+1} (1 + y_{k+1} + y_{k+1}^2) \rightarrow \frac{1}{\pi} \quad (2e)$$

$$= a_k ((1 + y_{k+1})^2)^2 - 2^{2k+3} y_{k+1} ((1 + y_{k+1})^2 - y_{k+1}) \quad (2f)$$

$$0 < a_k - \pi^{-1} \leq 16 \cdot 4^n 2 e^{-4^n 2 \pi} \quad (2g)$$

Identities 2d and 2f show how to save operations.

---

<sup>3</sup>using radix 10,000 and 1 million LIMBs.

## More iterations for $\pi$

Derived AGM iteration (second order) as in [hfloat: src/pi/pideriv.cc], #FPM=276.2:

$$x_0 = \sqrt{2} \quad (1a)$$

$$p_0 = 2 + \sqrt{2} \quad (1b)$$

$$y_1 = 2^{1/4} \quad (1c)$$

$$x_{k+1} = \frac{1}{2} \left( \sqrt{x_k} + \frac{1}{\sqrt{x_k}} \right) \quad (k \geq 0) \rightarrow 1 + \quad (1d)$$

$$y_{k+1} = \frac{y_k \sqrt{x_k} + \frac{1}{\sqrt{x_k}}}{y_k + 1} \quad (k \geq 1) \rightarrow 1 + \quad (1e)$$

$$p_{k+1} = p_k \frac{x_k + 1}{y_k + 1} \quad (k \geq 1) \rightarrow \pi + \quad (1f)$$

$$p_k - \pi = 10^{-2^{k+1}} \quad (1g)$$

Cubic AGM as in [hfloat: src/pi/picubagm.cc], #FPM=182.7:

$$a_0 = 1 \quad (2a)$$

$$b_0 = \frac{\sqrt{3} - 1}{2} \quad (2b)$$

$$a_{n+1} = \frac{a_n + 2b_n}{3} \quad (2c)$$

$$b_{n+1} = \sqrt[3]{\frac{b_n (a_n^2 + a_n b_n + b_n^2)}{3}} \quad (2d)$$

$$p_n = \frac{3a_n^2}{1 - \sum_{k=0}^n 3^k (a_k^2 - a_{k+1}^2)} \quad (2e)$$

Quintic (5th order) iteration as in [hfloat: src/pi/pi5th.cc], #FPM=353.2:

$$s_0 = 5(\sqrt{5} - 2) \quad (3a)$$

$$a_0 = \frac{1}{2} \quad (3b)$$

$$s_{n+1} = \frac{25}{s_n(z + x/z + 1)^2} \rightarrow 1 \quad (3c)$$

$$\text{where } x = \frac{5}{s_n} - 1 \rightarrow 4 \quad (3d)$$

$$\text{and } y = (x - 1)^2 + 7 \rightarrow 16 \quad (3e)$$

$$\text{and } z = \left( \frac{x}{2} \left( y + \sqrt{y^2 - 4x^3} \right) \right)^{1/5} \rightarrow 2 \quad (3f)$$

$$a_{n+1} = s_n^2 a_n - 5^n \left( \frac{s_n^2 - 5}{2} + \sqrt{s_n (s_n^2 - 2s_n + 5)} \right) \rightarrow \frac{1}{\pi} \quad (3g)$$

$$a_n - \frac{1}{\pi} < 16 \cdot 5^n e^{-\pi 5^n} \quad (3h)$$

## High order = fast?

Nonic (9th order) iteration as in [hfloat: src/pi/pi9th.cc], #FPM=273.7:

$$a_0 = \frac{1}{3} \quad (1a)$$

$$r_0 = \frac{\sqrt{3} - 1}{2} \quad (1b)$$

$$s_0 = (1 - r_0^3)^{1/3} \quad (1c)$$

$$t = 1 + 2 r_k \quad (1d)$$

$$u = (9 r_k (1 + r_k + r_k^2))^{1/3} \quad (1e)$$

$$v = t^2 + t u + u^2 \quad (1f)$$

$$m = \frac{27 (1 + s_k + s_k^2)}{v} \quad (1g)$$

$$a_{k+1} = m a_k + 3^{2k-1} (1 - m) \rightarrow \frac{1}{\pi} \quad (1h)$$

$$s_{k+1} = \frac{(1 - r_k)^3}{(t + 2 u) v} \quad (1i)$$

$$r_{k+1} = (1 - s_k^3)^{1/3} \quad (1j)$$

Summary of operation count vs. algorithms:

#FPM - algorithm name in hfloat

|         |   |                                |
|---------|---|--------------------------------|
| 78.424  | - | pi_agm_sch()                   |
| 98.424  | - | pi_agm()                       |
| 99.510  | - | pi_agm3(fast variant)          |
| 108.241 | - | pi_agm3(slow variant)          |
| 149.324 | - | pi_agm(quartic)                |
| 155.265 | - | pi_agm3(quartic, fast variant) |
| 164.359 | - | pi_4th_order(r=16 variant)     |
| 169.544 | - | pi_agm3(quartic, slow variant) |
| 170.519 | - | pi_4th_order(r=4 variant)      |
| 182.710 | - | pi_cubic_agm()                 |
| 200.261 | - | pi_3rd_order()                 |
| 255.699 | - | pi_2nd_order()                 |
| 273.763 | - | pi_9th_order()                 |
| 276.221 | - | pi_derived_agm()               |
| 353.202 | - | pi_5th_order()                 |

## The binary splitting algorithm

Even if a series is as ‘good’ as Chudnovsky’s famous series for  $\pi$ :

$$\frac{1}{\pi} = \frac{6541681608}{\sqrt{640320}^3} \sum_{k=0}^{\infty} \left( \frac{13591409}{545140134} + k \right) \left( \frac{(6k)!}{(k!)^3 (3k)!} \frac{(-1)^k}{640320^{3k}} \right) \quad (1a)$$

$$= \frac{12}{\sqrt{640320}^3} \sum_{k=0}^{\infty} (-1)^k \frac{(6k)!}{(k!)^3 (3k)!} \frac{13591409 + k 545140134}{(640320)^{3k}} \quad (1b)$$

... the total work is proportional  $N^2$ , which makes it impossible to compute billions of digits from linearly convergent series

There is an alternative way to evaluate a sum  $\sum_{k=0}^{N-1} a_k$  of rational summands. One looks at the ratios  $r_k$  of consecutive terms:

$$r_k := \frac{a_k}{a_{k-1}} \quad (2)$$

Set  $a_{-1} := 1$  to avoid a special case for  $k = 0$ . One has

$$\sum_{k=0}^{N-1} a_k =: r_0 (1 + r_1 (1 + r_2 (1 + r_3 (1 + \dots (1 + r_{N-1}) \dots)))) \quad (3)$$

Now define

$$r_{m,n} := r_m (1 + r_{m+1} (\dots (1 + r_n) \dots)) \quad \text{where } m < n \quad (4a)$$

$$r_{m,m} := r_m \quad (4b)$$

then

$$r_{m,n} = \frac{1}{a_{m-1}} \sum_{k=m}^n a_k \quad (5)$$

and especially

$$r_{0,n} = \sum_{k=0}^n a_k \quad (6)$$

With

$$r_{m,n} = r_m + r_m \cdot r_{m+1} + r_m \cdot r_{m+1} \cdot r_{m+2} + \dots \quad (7a)$$

$$\dots + r_m \cdot \dots \cdot r_x + r_m \cdot \dots \cdot r_x \cdot [r_{x+1} + \dots + r_{x+1} \cdot \dots \cdot r_n]$$

$$= r_{m,x} + \prod_{k=m}^x r_k \cdot r_{x+1,n} \quad (7b)$$

The product telescopes, one gets

$$r_{m,n} = r_{m,x} + \frac{a_x}{a_{m-1}} \cdot r_{x+1,n} \quad (8)$$

(where  $m \leq x < n$ ).

Now we can formulate the binary splitting algorithm by giving a binsplit function `r`:

```
function r(function a, int m, int n)
{
  rational ret
  if m==n then
  {
    ret := a(m)/a(m-1)
  }
  else
  {
    x := floor( (m+n)/2 )
    ret := r(a,m,x) + a(x) / a(m-1) * r(a,x+1,n)
  }
  return ret
}
```

Here `a(k)` must be a function that returns the `k`-th term of the series we wish to compute, in addition one must have `a(-1)=1`. To compute `arctan(1/10)` one would use

```
function a(k)
{
  if k<0 then return 1
  else return (-1)^k/((2*k+1)*10^(2*k+1))
}
```

Calling `r(a,0,N)` returns  $\sum_{k=0}^N a_k$ .

It is instructive to modify the binsplit function so it's self-documenting. Simple enhancements will print out the calling tree in the process:

```
function r(function a, int m, int n, int i)
{
  rational ret
  indent(i)
  print( "r(", m, n, ")" )
  if m==n then
  {
    ret := a(m)/a(m-1)
  }
  else
  {
    x := floor( (m+n)/2 )
    ret := r(a,m,x, i+1) + a(x) / a(m-1) * r(a,x+1,n, i+1)
  }
  indent(i)
  print( "^== ", ret)
  return ret
}
```

The additional parameter `i` is the level of depth in the calling tree and the function `indent(i)` shall indent the following print output (that is, it prints `8·i` whitespaces that are not followed by a newline). Using a function  $a(k) = 2^{-(k+1)}$  one gets, upon calling `r(a, 0, 6, 0)`, the following output

```

r(0, 6)
  r(0, 3)
    r(0, 1)
      r(0, 0)
      ^== 1/2
      r(1, 1)
      ^== 1/2
    ^== 3/4
    r(2, 3)
      r(2, 2)
      ^== 1/2
      r(3, 3)
      ^== 1/2
    ^== 3/4
  ^== 15/16
  r(4, 6)
    r(4, 5)
      r(4, 4)
      ^== 1/2
      r(5, 5)
      ^== 1/2
    ^== 3/4
    r(6, 6)
    ^== 1/2
  ^== 7/8
^== 127/128

```

A fragment like

```

r(2, 3)
  r(2, 2)
  ^== 1/2
  r(3, 3)
  ^== 1/2
^== 3/4

```

means: “ $r(2,3)$  first called  $r(2,2)$  [which returned  $1/2$ ], then called  $r(3,3)$  [which returned  $1/2$ ]. Then  $r(2,3)$  returned  $3/4$ .”

The involved work is only  $O((\log N)^2 M(N))$ , where  $M(N)$  is the complexity of one  $N$ -bit multiplication. This means that sums of linear but sufficient convergence are again candidates for high precision computations.

In addition, the ratio  $r_{0,N-1}$  (that is, the sum of the first  $N$  terms) can be reused if one wants to evaluate the sum to a higher precision than before.

If one wants to stare at zillions of decimal digits of the floating point expansion then one division is also needed which costs not more than 4 multiplications.

Note that this algorithm can trivially be extended (or rather simplified) to infinite products, e.g. matrix products as Bellard’s

$$\prod_{k=0}^{\infty} \begin{bmatrix} \frac{2(k-\frac{1}{2})(k+2)}{27(k+\frac{2}{3})(k+\frac{4}{3})} & 10 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & \pi + 6 \\ 0 & 1 \end{bmatrix} \quad (9)$$

More details are given in the online draft of my book:

`http://www.jjj.de/fxt/#fxtbook`

Thanks for your feedback!

Jörg Arndt      <arndt@jjj.de>