

How to compute π to 10^{12} digits

A crash course in high precision arithmetics

Part I: Fast multiplication

Part II: Iteration schemes, the AGM, binary splitting

Jörg Arndt <arndt@jjj.de>

*”Why make things difficult, when it is possible to make them cryptic
and totally illogic, with just a little bit more effort?”*

– Aksel Peter Jørgensen

Multiplication is convolution

Multiplication of two numbers is essentially a convolution of the sequences of their digits. The *convolution* of the two sequences $a_k, b_k, k = 0 \dots N - 1$ is defined as the sequence c where

$$c_k := \sum_{i+j=k} a_i b_j \quad k = 0 \dots 2N - 2 \quad (1)$$

A (n -digit) number written in radix R as

$$a_{n-1} \ a_{n-2} \ \dots \ a_2 \ a_1 \ a_0 \quad (2)$$

denotes a quantity of

$$\sum_{i=0}^{n-1} a_i \cdot R^i = a_{n-1} \cdot R^{n-1} + a_{n-2} \cdot R^{n-2} + \dots + a_1 \cdot R + a_0 \quad (3)$$

For example, with decimal numbers one has $R = 10$, and the number 578 equals $5 \cdot 10^2 + 7 \cdot 10^1 + 8 \cdot 10^0$.

The product of two numbers is almost the polynomial product

$$\sum_{k=0}^{2N-2} c_k R^k := \sum_{i=0}^{N-1} a_i R^i \cdot \sum_{j=0}^{N-1} b_j R^j \quad (4)$$

The c_k are found by comparing coefficients: they must satisfy the convolution equation 1.

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 $c'_k = c_k \bmod R$ and add $(c_k - c'_k)/R$ to its left neighbor.

Multiplication is convolution [cont.]

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

$$\begin{array}{r}
 82 \times 34 \\
 \hline
 3 2 \\
 2 4 \\
 \hline
 = 2 8
 \end{array}$$

We have seen that the carries can be delayed to the end of the computation:

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

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

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

The value of the polynomial $24x^2 + 38x + 8$ for $x = 10$ is 2788.

2-way splitting (Karatsuba multiplication)

Split the numbers A and B (assumed to have approximately the same length) into two pieces

$$\begin{aligned} A &= x A_1 + A_0 \\ B &= x B_1 + B_0 \end{aligned} \tag{1}$$

where x is a power of the radix (for decimal numbers the radix is 10) close to the half length of A and B . The usual multiplication scheme needs 4 multiplications with half precision for one multiplication with full precision:

$$A B = A_0 \cdot B_0 + x (A_0 \cdot B_1 + B_0 \cdot A_1) + x^2 A_1 \cdot B_1 \tag{2}$$

Only the multiplications $A_i \cdot B_j$ need to be considered. The multiplications by x , a power of the radix are only shifts. If we use the relation

$$A B = (1 + x) A_0 \cdot B_0 + x (A_1 - A_0) \cdot (B_0 - B_1) + (x + x^2) A_1 \cdot B_1 \tag{3}$$

we need 3 multiplications with half precision for one multiplication with full precision. Applying the scheme recursively until the numbers to multiply are of machine size we obtain an algorithm whose asymptotic cost is $\sim N^{\log_2(3)} \approx N^{1.585}$.

For squaring use

$$A^2 = (1 + x) A_0^2 - x (A_1 - A_0)^2 + (x + x^2) A_1^2 \tag{4}$$

We compute $8231^2 = 67749361$:

$$\begin{aligned} 8231^2 &= \\ &= (100 \cdot 82 + 31)^2 \\ &= (1 + 100) \cdot 31^2 - 100 \cdot (82 - 31)^2 + (100 + 100^2) \cdot 82^2 \\ &= (1 + 100) \cdot [961] - 100 \cdot [2601] + (100 + 100^2) \cdot [6724] \\ &= 961 + 96100 - 260100 + 672400 + 67240000 \\ &= 67749361 \end{aligned}$$

And, yes, the scheme also works for polynomials.

3-way splitting (Toom-Cook multiplication)

The 3-way multiplication scheme of Bodrato and Zanoni:

$$A = a_2x^2 + a_1x + a_0$$

$$B = b_2x^2 + b_1x + b_0$$

$$S_0 = a_0 * b_0$$

$$S_1 = (a_2+a_1+a_0) * (b_2+b_1+b_0)$$

$$S_2 = (4*a_2+2*a_1+a_0) * (4*b_2+2*b_1+b_0)$$

$$S_3 = (a_2-a_1+a_0) * (b_2-b_1+b_0)$$

$$S_4 = a_2 * b_2$$

$$S_2 = (S_2 - S_3)/3 \quad \backslash\backslash \text{ division by 3}$$

$$S_3 = (S_1 - S_3)/2$$

$$S_1 = S_1 - S_0$$

$$S_2 = (S_2 - S_1)/2$$

$$S_1 = S_1 - S_3 - S_4$$

$$S_2 = S_2 - 2*S_4$$

$$S_3 = S_3 - S_2$$

$$P = S_4x^4 + S_2x^3 + S_1x^2 + S_3x + S_0$$

$$P - A*B \quad \backslash\backslash \quad == \text{ zero}$$

4-way splitting multiplication

The 4-way multiplication scheme of Bodrato and Zanoni:

$$A = a_3x^3 + a_2x^2 + a_1x + a_0$$

$$B = b_3x^3 + b_2x^2 + b_1x + b_0$$

$$S1 = a_3b_3$$

$$S2 = (8a_3+4a_2+2a_1+a_0)(8b_3+4b_2+2b_1+b_0)$$

$$S3 = (a_3+a_2+a_1+a_0)(b_3+b_2+b_1+b_0)$$

$$S4 = (-a_3+a_2-a_1+a_0)(-b_3+b_2-b_1+b_0)$$

$$S5 = (8a_0+4a_1+2a_2+a_3)(8b_0+4b_1+2b_2+b_3);$$

$$S6 = (-8a_0+4a_1-2a_2+a_3)(-8b_0+4b_1-2b_2+b_3)$$

$$S7 = a_0b_0$$

$$S2 += S5$$

$$S4 -= S3$$

$$S6 -= S5$$

$$S4 /= 2 \quad \backslash \backslash$$

$$S5 -= S1$$

$$S5 -= (64*S7)$$

$$S3 += S4$$

$$S5 *= 2; \quad S5 += S6$$

$$S2 -= (65*S3)$$

$$S3 -= S1$$

$$S3 -= S7$$

$$S4 = -S4 \quad \backslash \backslash$$

$$S6 = -S6 \quad \backslash \backslash$$

$$S2 += (45*S3)$$

$$S5 -= (8*S3)$$

$$S5 /= 24 \quad \backslash \backslash \text{ division by 24}$$

$$S6 -= S2$$

$$S2 -= (16*S4)$$

$$S2 /= 18 \quad \backslash \backslash \text{ division by 18}$$

$$S3 -= S5$$

$$S4 -= S2$$

$$S6 += (30*S2)$$

$$S6 /= 60 \quad \backslash \backslash \text{ division by 60}$$

$$S2 -= S6$$

$$P = S1x^6 + S2x^5 + S3x^4 + S4x^3 + S5x^2 + S6x + S7;$$

$$P - A*B \quad \backslash \backslash == \text{zero}$$

FFT-multiplication

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 $\sim N \cdot \log N$ operations. Instead of the direct convolution ($\sim N^2$) one proceeds as follows:

- Compute the FFTs of both factors.
- Multiply the transformed sequences element-wise.
- Compute inverse transform of the product.

Note that (1) the multiplication of two polynomials can be achieved by the (more complicated) scheme:

- evaluate both polynomials at sufficiently many points
- 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.

Re-launching our example ($82 \cdot 34 = 2788$), we use the fourth roots of unity ± 1 and $\pm i$:

$a = (8x + 2)$		\times	$b = (3x + 4)$		$c = a b$
$+1$	$+10$		$+7$		$+70$
$+i$	$+8i + 2$		$+3i + 4$		$+38i - 16$
-1	-6		$+1$		-6
$-i$	$-8i + 2$		$-3i + 4$		$-38i - 16$
<hr/>					
$c = (24x^2 + 38x + 8)$					

This table has to be read as follows: 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.

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.

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:

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*:

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

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

Radix-4 FFT algorithms

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:

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

Now we have expressed the length- $N = 2^n$ FFT as one length- $N/2$ and two length- $N/4$ FFTs.

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:

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

Cyclic vs. linear 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 C with elements C_τ as:

$$C = A \circledast B \quad (1a)$$

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

The last equation may be rewritten as

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

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

The FFT scheme

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

computes the cyclic convolution (the polynomial product is computed modulo z^n).

But we want the *linear* convolution:

$$C_\tau := \sum_{x=0}^{n-1} a_x b_{(\tau-x)} \bmod \text{nothing} \quad (4)$$

Solution: We (roughly) double the length of the input sequences by appending zeros to the end. The polynomial product is computed modulo z^{2n} , but $2n > \deg AB$.

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!

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
}
```

Division: Inversion

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

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

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

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 work¹ 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:

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

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

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

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

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

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

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

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

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

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

¹The 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 $x_0 \approx \frac{1}{\sqrt{d}}$ 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 $x_k = \frac{1}{\sqrt{d}}(1 + e)$ then

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

Similar considerations as above (with squaring considered as expensive as multiplication²) 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:

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

²Indeed it costs about $\frac{2}{3}$ of a multiplication.

Inverse a -th root, a general expression

There is a nice general formula that allows to build iterations with arbitrary order of convergence for $1/\sqrt[a]{d} = d^{-1/a}$ 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

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

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:

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

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

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

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

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

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

There seems to be no standard term for emphasizing the number of derivatives vanishing at the fixed point: *super-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,

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

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.

Write

$$S = x - U_1 \frac{f}{1! f'} - U_2 \frac{f^2}{2! f'^3} - U_3 \frac{f^3}{3! f'^5} - \dots - U_n \frac{f^n}{n! f'^{2n-1}} - \dots \quad (4)$$

then $U_1 = 1$, $U_2 = f''$, $U_3 = 3f''^2 - f' f'''$, and we have the recursion

$$U_n = (2n - 3) f'' U_{n-1} - f' U'_{n-1} \quad (5)$$

Alternatively write

$$S = x - Y_1 \left(\frac{f}{f'} \right) - Y_2 \left(\frac{f}{f'} \right)^2 - Y_3 \left(\frac{f}{f'} \right)^3 - \dots - Y_t \left(\frac{f}{f'} \right)^t - \dots \quad (6)$$

then $Y_1 = 1$ and

$$Y_t = \frac{1}{t} \left(2(t-1) \frac{f''}{2f'} Y_{t-1} - Y'_{t-1} \right) \quad (7)$$

A simple derivation of Schröder's formula

The starting point is the Taylor series of a function f around x_0 :

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + \dots \quad (1)$$

Now let $f(x_0) = y_0$ and r be the zero of f (that is, $f(r) = 0$). We expand the inverse $g = f^{-1}$ around y_0 :

$$g(0) = g(y_0) + g'(y_0)(0 - y_0) + \frac{1}{2}g''(y_0)(0 - y_0)^2 + \dots \quad (2)$$

Using $x_0 = g(y_0)$ and $g(0) = r$ we obtain

$$r = x_0 - g'(y_0)f(x_0) + \frac{1}{2}g''(y_0)f(x_0)^2 - \frac{1}{6}g'''(y_0)f(x_0)^3 + \dots \quad (3)$$

Remains to express the derivatives of the inverse g in terms of (derivatives of) f . Set

$$f \circ g = \text{id}, \quad \text{that is: } f(g(x)) = x \quad (4)$$

and derive the equation (chain rule) to obtain $g'(f(x))f'(x) = 1$, so $g'(y) = \frac{1}{f'(x)}$. Derive $f(g(x)) = x$ multiple times to obtain (arguments y of g and x of f are omitted for readability):

$$1 = f'g' \quad (5a)$$

$$0 = g'f'' + f'^2g'' \quad (5b)$$

$$0 = g'f''' + 3f'f''g'' + f'^3g''' \quad (5c)$$

$$0 = g'f'''' + 4f'g''f''' + 3f''^2g'' + 6f'^2f''g''' + f'^4g'''' \quad (5d)$$

This system of linear equations in the derivatives of g can be solved successively for g' , g'' , g''' , \dots :

$$g' = \frac{1}{f'} \quad (6a)$$

$$g'' = -\frac{f''}{f'^3} \quad (6b)$$

$$g''' = \frac{1}{f'^5} (3f''^2 - f'f''') \quad (6c)$$

Thereby equation 3 can be written as

$$r = x - \frac{1}{f'}f + \frac{1}{2}\left(-\frac{f''}{f'^3}\right)f^2 - \frac{1}{6}\left(\frac{1}{f'^5}(3f''^2 - f'f''')\right)f^3 + \dots \quad (7)$$

which is Schröder's iteration.

Householder's formula

For $n \geq 2$ the expression

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

gives a n -th order iteration for a (simple) root r of f .

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

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

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

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

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:

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

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

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

For multiple roots use iterations for f/f' , for example (Schröder):

$$S^{\%} = x - \frac{ff'}{(f'^2 - ff'')} - \frac{f^2f'(ff'f''' - 2ff''^2 + f'^2f'')}{2(ff'' - f'^2)^3} - \dots \quad (5)$$

The AGM

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

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

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

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

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

$$\pi \approx \frac{2 a_n^2}{t_n} \tag{2}$$

This is a special case of Legendre's relation (a relation between complete elliptic integrals) and was discovered (independently) 1976 by Salamin and Brent. However, the relations has been given 200 years earlier by Gauss.

Superlinear iterations for π

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 π to 4 million decimal digits³ 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.

³using radix 10,000 and 1 million LIMBs.

More iterations for π

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 (for products)

We motivate the binsplit algorithm by giving the analogue for the fast computation of the factorial. Define $f_{m,n} := m \cdot (m+1) \cdot (m+2) \cdots (n-1) \cdot n$, then $n! = f_{1,n}$. We compute $n!$ by recursively using the relation $f_{m,n} = f_{m,x} \cdot f_{x+1,n}$ where $x = \lfloor (m+n)/2 \rfloor$:

```
indent(i)=for(k=1,8*i,print1(" "));  \\ aux: print 8*i spaces
F(m, n, i=0)=
{ /* Factorial, self documenting */
  local(x, ret);
  indent(i); print( "F(", m, ", ", " ", n, ")");
  if ( m==n, /* then: */
      ret = m;  \\ == F(m,m)
  , /* else: */
      x = floor( (m+n)/2 );
      ret = F(m, x, i+1) * F(x+1, n, i+1);
  );
  indent(i); print( "^== ", ret);
  return( ret );
}
```

The function prints the intermediate values occurring in the computation. The additional parameter `i` keeps track of the calling depth, used with the auxiliary function `indent()`.

The intermediate quantities with the computation of $8! = F(1,8)$ are

```

F(1, 8)
  F(1, 4)
    F(1, 2)
      F(1, 1)
        ^== 1
      F(2, 2)
        ^== 2
    ^== 2
    F(3, 4)
      F(3, 3)
        ^== 3
      F(4, 4)
        ^== 4
    ^== 12
  ^== 24
  F(5, 8)
    F(5, 6)
      F(5, 5)
        ^== 5
      F(6, 6)
        ^== 6
    ^== 30
    F(7, 8)
      F(7, 7)
        ^== 7
      F(8, 8)
        ^== 8
    ^== 56
  ^== 1680
^== 40320
```

A fragment like

```

F(5, 6)
  F(5, 5)
    ^== 5
  F(6, 6)
    ^== 6
  ^== 30
```

says “`F(5,6)` called `F(5,5)` [which returned 5], then called `F(6,6)` [which returned 6]. Then `F(5,6)` returned 30.” For the computation of other products modify the line `ret=m`; as indicated in the code.

The binary splitting algorithm (for sums)

For the evaluation of a sum $\sum_{k=0}^{N-1} a_k$ we use the ratios r_k of consecutive terms:

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

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

$$\sum_{k=0}^{N-1} a_k =: r_0 (1 + r_1 (1 + r_2 (1 + r_3 (1 + \dots (1 + r_{N-1}) \dots)))) \quad (2)$$

Now define

$$r_{m,n} := r_m (1 + r_{m+1} (\dots (1 + r_n) \dots)) \quad \text{where } m < n \quad (3a)$$

$$r_{m,m} := r_m \quad (3b)$$

then

$$r_{m,n} = \frac{1}{a_{m-1}} \sum_{k=m}^n a_k \quad (4)$$

and especially

$$r_{0,n} = \sum_{k=0}^n a_k \quad (5)$$

We have

$$r_{m,n} = r_m + r_m \cdot r_{m+1} + r_m \cdot r_{m+1} \cdot r_{m+2} + \dots \quad (6a)$$

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

$$= r_{m,x} + \prod_{k=m}^x r_k \cdot r_{x+1,n} \quad (6b)$$

The product telescopes, one gets (for $m \leq x < n$)

$$r_{m,n} = r_{m,x} + \frac{a_x}{a_{m-1}} \cdot r_{x+1,n} \quad (7)$$

The binary splitting algorithm (for sums) [cont.]

```

R(m, n)=
{ /* Rational binsplit */
  local(x, ret);
  if ( m==n, /* then: */
    ret = A(m)/A(m-1);
  , /* else: */
    x = floor( (m+n)/2 );
    ret = R(m, x) + A(x) / A(m-1) * R(x+1, n);
  );
  return( ret );
}

```

The intermediate values with the computation of $\sum_{k=0}^6 2^{-(k+1)}$ are

```

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

```

Reducing the intermediate fractions to lowest terms is necessary with certain sums. For example, computing $\arctan(1/10)$ without reduction shows the explosive growth of intermediate quantities:

```

Q(0, 6)
  Q(0, 3)
    Q(0, 1)
      Q(0, 0)
        ^== [1, 10]
      Q(1, 1)
        ^== [-10, 3000]
    ^== [29900, 300000]
    Q(2, 3)
      Q(2, 2)
        ^== [3000, -500000]
      Q(3, 3)
        ^== [-500000, 70000000]
    ^== [-1042500000000000000, 1750000000000000000]
  ^== [1569781275000000000000000000, 1575000000000000000000000000]
[--snip--]

```

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>