

# Bit-Serial, VLSI Architecture for the Implementation of Maximum-Length Number-Theoretic Transforms Using Mixed Basis Representations

M.G.Parker, M.Benaissa

School of Eng., Univ.of Huddersfield, Huddersfield HD1 3DH, UK.

## Abstract

*Fermat and Mersenne NTTs are relatively easy to implement, but unsuitable for many DSP applications, due to small block length over wordlength. This paper presents VLSI design techniques appropriate for a wider range of NTTs, including maximum-length NTTs, and presents a systolic architecture exploiting block-length factorisation to decompose the architecture into sub-modules, themselves systolic NTTs. The design avoids the explicit implementation of modular addition, accumulation and multiplication by using systolic basis converters which, inherently, perform the aforementioned tasks. The implementation of a maximum-length, 60pt NTT is described as an example. It uses binary to ternary basis conversion to perform all addition and multiplication, and 'ternary basis compression' to perform accumulation.*

## 1 Introduction

A general 1-D NTT of a sequence  $x[n]$  is defined as,

$$X[k] = \left\langle \sum_{n=0}^{N-1} x[n] \cdot p^{nk} \right\rangle_M \quad 0 \leq n, k \leq N-1 \quad (1)$$

where  $p$  has order  $N \bmod M$ ,  $\langle N^{-1} \rangle_M$  exists and is unique, and an  $N$ -point NTT can be defined over all the moduli factors of  $M$  [1]. By choosing all NTTs with 2 as the kernel, over a modulus,  $M$ , where  $\langle 2^N \rangle_M = 1$ , a wide range of block lengths, for a given dynamic range, become available [2,3]. As with FNTs and MNTs, multiplication can be implemented as bit shifts but, in general, the problem of modulus overflow is non-trivial. This paper

shows how NTT block length factorisation enables a regular and highly-pipelined decomposition. Basis converters [4] distribute the addition over a series of simple cells, whilst converting the data basis to a form which eliminates multiplication. Finally, the concepts are illustrated by considering the specific implementation of a 60pt NTT, mod 61, using binary and ternary representations.

## 2 Factorisation of NTTs

An  $N$ -length NTT is factorised as follows,

$$X[k] = X_\mu[\theta] \text{ for } k = N_0 \cdot \theta + \mu$$

$$\text{where } 0 \leq \theta < N_1, 0 \leq k < N, \\ 0 \leq \mu < N_0 \text{ and } N = N_0 \cdot N_1,$$

and,

$$X_\mu[\theta] = \left\langle \sum_{n_0=0}^{N_0-1} \sum_{n_1=0}^{N_1-1} g^{n_1 \cdot \theta} \cdot [g^{D_{N_0}((n_1+N_1 \cdot n_0) \cdot \mu)}]_{p^{((n_1+N_1 \cdot n_0) \cdot \mu)_{N_0}}} \cdot x[n_1 + N_1 \cdot n_0] \right\rangle_M$$

where  $g = \langle p^{N_0} \rangle_M$ ,  $D_z() = () \text{ DIV } z$ , and  $n = n_1 + N_1 \cdot n_0$  where,  $0 \leq n_0 < N_0, 0 \leq n_1 < N_1$

We have expressed  $X_\mu[\theta]$  as a summation of  $N_0$  separate  $N_1$ pt NTTs, where  $x[n]$  is pre-processed to become,

$$g^{D_{N_0}((n_1+N_1 \cdot n_0) \cdot \mu)} \cdot p^{((n_1+N_1 \cdot n_0) \cdot \mu)_{N_0}} \cdot x[n_1 + N_1 \cdot n_0]$$

The architectures, developed, perform the pre-processing by pipelining the small powers of  $p$  multiplications, ( $p^0 \dots p^{N_0-1}$ ), followed by the small power of  $g$  multiplications, (implemented as a modification to the subsequent  $N_1$ pt NTT architectures).

Fig 1 shows the realisation of this decomposition. Gate rotators re-route the stream of partial

products, for successive  $x[n]$ , between each processing stage and are general to any block length. Fig 1 also shows the equivalent, unfactorised, architecture. The design combines  $N_0$ ,  $N_{1pt}$ , NTTs, together with a pre-processing stage, to form the complete NTT, and, if  $N_1$  is, itself, factorisable, the decomposition is repeated for each  $N_{1pt}$  NTT.

The  $N_{1pt}$  NTT modification, necessary for the extra small powers of  $g$  multiplications, manifests itself as a series of cumulative pulses to the  $GRN_1$  modules. These pulses update the re-routing pattern within the rotators accordingly.

### 3 NTT Multiplications Using Basis Conversion

To simplify designs, we assume  $\alpha$ , the i/p basis parameter (see appendix) =  $p$ , the kernel of the transform. (Binary input data  $\Rightarrow \alpha = p = 2$ ). Thus small power of  $p$  multiplications become a few additions, mod  $M$ . For instance, if  $\langle p^{N_h} \rangle_M = h = p^0 + p^1$ , multiplications by  $p^0 \dots p^{N_h-1}$  become single additions of bit re-orderings of  $x[n]$ . As shown in the previous section, we require multiplications by  $p^0 \dots p^{N_0-1}$ . These multiplications become single additions if  $N_h \geq N_0$ . Ideally  $N_h = N_0$ . These additions can be incorporated into the basis conversion process.

The converter changes the data from a  $p$ -basis to a  $g$ -basis to simplify multiplication by small powers of  $g$ .  $g$  is the kernel of the  $N_{1pt}$  NTTs and, as  $\langle g^{N_1} \rangle_M = 1$ , power of  $g$  multiplications become bit re-orderings.

Each converter cell contains a number of digits, ('bits' for binary digits), which are related by a state transition table and identical cells are 'joined' together to form the complete systolic converter.

If  $N_1$  is, itself, factored, the same process is repeated and  $g$ -basis to 'power of  $g$ '-basis converters will be required, (and so on).

### 4 Compression Replaces Accumulation

Each accumulator sums  $N$  data words over a  $g$ -basis to calculate  $X[k]$ . By interpreting these  $N$  serial data words as one long data word with an identical basis flow, (apart from  $j'_N = N \cdot j'$ , where  $j'_N$  refers to the long word) (see appendix), this single data word =  $X[k - z]$ , (where  $g^{j'} = p^z \Rightarrow z = N_0 \cdot j'$ ). As the mapping from  $X[k]$  to  $X[k - z]$  is one-to-one, accumulation is not strictly necessary, and for many applications, (such as convolution), accumulators can be omitted from the design, giving substantial hardware savings.

If it is desired to reduce the value of  $j'_N$  to a minimum, basis compressors, similar in design to the basis converter, but with minimal feedback inserted between cells, are used.

## 5 Optimal NTTs

The design principles, developed in previous sections, are applicable to any NTT. One class of NTTs for which they are particularly suited is maximum-length NTTs. These occur when  $M$  is prime and  $N = M - 1$  (eqn 1). Table 1 list a few such NTTs for  $p = 2$ . When  $N$  is factorisable, a suitable  $N_h \geq N_0$  should be found so that  $h$  is a simple additive combination of small powers of 2, where  $h = \langle 2^{N_h} \rangle_M$ . As explained in the last section, these simple additions can then be incorporated into the basis conversion where the basis converters are required to convert from a binary basis to a  $g$ -basis. Suitable  $N_h$ ,  $N_0$ ,  $h$  and  $g$  are given in Table 1 for each example, along with the expression of  $h$  as a simple modular addition.

M	N	$N_h$	$N_0$	$h$	$h_{Add}$	$g$
11	10	8	5	3	$2^0 + 2^1$	32
13	12	9	6	5	$2^0 + 2^2$	64
29	28	5	4	3	$2^0 + 2^1$	16
61	60	6	6	3	$2^0 + 2^1$	3
4093	4092	12	12	3	$2^0 + 2^1$	3

Table 1: Possible First-Stage Factorisation of Selected Prime Moduli

## 6 High Thr'put VLSI Implementation of a 60pt NTT, Mod 61

This example of a maximum-length NTT implements eqn 1 when  $p = 2$ ,  $M = 61$  and  $N = 60$ , and is shown in Fig 2. The input basis flow is given by,

$$\leftarrow 2^5, 2^4, 2^3, 2^2, 2^1, 2^0$$

where  $j = j' = 6$  and  $d = 1$ ,

and, after binary to ternary conversion, by,

$$\leftarrow \begin{matrix} 3^4, & 3^3, & 3^2, & 3^1, & 3^0 \\ 3^9, & 3^8, & 3^7, & 3^6, & 3^5 \end{matrix}$$

where  $\langle 3^{10} \rangle_{61} = 1$ ,  $j = 5$ ,  $j' = 6$  and  $d = 2$ .

Both bases span the data range 0 - 60 and are therefore acceptable.

We factorise  $N = 60$  as  $2.3.5.2 = N_0.N_1.N_2.N_3$ . As  $\langle 2^6 \rangle_{61} = 3 = 2^0 + 2^1$ , multiplications by



