

# Αρχιτεκτονική Προηγμένων Υπολογιστών και Επιταχυντών – Αναφορά 1<sup>ου</sup> εργαστηρίου

Δημήτριος Ορέστης Βαγενάς, 10595

Αγγελική Στρατάκη, 10523

## Ερώτημα 1:

Για το πρώτο εργαστήριο χρησιμοποιήσαμε το εργαλείο Vivado HLS με σκοπό την διαδικασία επιτάχυνσης του αλγορίθμου για πολλαπλασιασμό πινάκων MATRIX\_MUL, σε γλώσσα C. Σε ένα testbench έχουμε αρχικοποιήσει τους πίνακες A και B με τυχαία ψευδό-τυχαίες τιμές στο εύρος 0-255 για την S/W λύση φτιάξαμε την συνάρτηση matrix\_multiply\_sw και την H/W λύση την συνάρτηση matrix\_multiply\_hw, με σκοπό την διασφάλιση της σωστής εκτέλεσης της hardware accelerator λειτουργικότητας, που θα τροποποιήσουμε.

- Testbench:

```
#include <iostream>
#include <cstdlib> // For rand() and srand()
#include <ctime> // For time()
#include <iomanip> // For std::setw

#define LM 6
#define LN 6
#define LP 6

#define M (1 << LM)
#define N (1 << LN)
#define P (1 << LP)

void matrix_multiply_hw(const uint8_t A[N][M], const uint8_t B[M][P], uint32_t C[N][P]);

// Software matrix multiplication for verification
void matrix_multiply_sw(const uint8_t A[N][M], const uint8_t B[M][P], uint32_t C[N][P]) {
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            C[i][j] = 0; // Initialize element C[i][j]
            for (int k = 0; k < M; ++k) {
                C[i][j] += (uint32_t)A[i][k] * (uint32_t)B[k][j];
            }
        }
    }
}

int main() {
    // Initialize random seed
    std::srand(std::time(0));

    // Declare matrices A, B, and C for both HW and SW computation
    uint8_t A[N][M];
    uint8_t B[M][P];
    uint32_t C_hw[N][P];
    uint32_t C_sw[N][P];

    // Initialize matrices A and B with random values in the range 0-255
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < M; ++j) {
            A[i][j] = rand() % 256;
        }
    }
    for (int i = 0; i < M; ++i) {
        for (int j = 0; j < P; ++j) {
            B[i][j] = rand() % 256;
        }
    }

    // Perform matrix multiplication using both hardware and software functions
    matrix_multiply_hw(A, B, C_hw);
    matrix_multiply_sw(A, B, C_sw);

    // Verify the results
    bool success = true;
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            if (C_hw[i][j] != C_sw[i][j]) {
                success = false;
                break;
            }
        }
        if (!success) break;
    }

    // Print matrices and results
    std::cout << "Matrix A:" << std::endl;
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < M; ++j) {
            std::cout << std::setw(4) << static_cast<int>(A[i][j]) << " ";
        }
        std::cout << std::endl;
    }

    std::cout << "Matrix B:" << std::endl;
    for (int i = 0; i < M; ++i) {
        for (int j = 0; j < P; ++j) {
            std::cout << std::setw(4) << static_cast<int>(B[i][j]) << " ";
        }
        std::cout << std::endl;
    }

    std::cout << "Matrix C (Hardware Result):" << std::endl;
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            std::cout << std::setw(6) << C_hw[i][j] << " ";
        }
        std::cout << std::endl;
    }

    std::cout << "Matrix C (Software Result):" << std::endl;
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            std::cout << std::setw(6) << C_sw[i][j] << " ";
        }
        std::cout << std::endl;
    }

    // Print test result
    if (success) {
        std::cout << "Test Passed" << std::endl;
    } else {
        std::cout << "Test Failed" << std::endl;
    }

    return 0;
}
```

Έχουμε ορίσει lm=ln=lp=6 για το επόμενο ερώτημα.

## Ερώτημα 2:

- Η συνάρτηση matrix\_multiply\_hw ορίζεται αρχικά ως εξής:

```
#include <stdint.h>

#define LM 6 // Define  $\underline{l}_m$  (for example,  $\underline{l}_m = 6 \rightarrow m = 2^6 = 64$ )
#define LN 6 // Define  $\underline{l}_n$  (for example,  $\underline{l}_n = 6 \rightarrow n = 2^6 = 64$ )
#define LP 6 // Define  $\underline{l}_p$  (for example,  $\underline{l}_p = 6 \rightarrow p = 2^6 = 64$ )

#define M (1 << LM) //  $m = 2^{\underline{l}_m}$ 
#define N (1 << LN) //  $n = 2^{\underline{l}_n}$ 
#define P (1 << LP) //  $p = 2^{\underline{l}_p}$ 

void matrix_multiply_hw(
    const uint8_t A[N][M],
    const uint8_t B[M][P],
    uint32_t C[N][P]
) {
    // Perform matrix multiplication:  $C = A * B$ 
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            C[i][j] = 0; // Initialize element C[i][j]
            for (int k = 0; k < M; ++k) {
                C[i][j] += (uint32_t)A[i][k] * (uint32_t)B[k][j];
            }
        }
    }
}
```

Κάνοντας σύνθεση τη παραπάνω σχεδίαση (C Synthesis) με default settings έχουμε:

Target	Estimated	Uncertainty		Summary			
10.00 ns	6.304 ns	2.70 ns					

Latency (cycles)		Latency (absolute)		Interval (cycles)		
min	max	min	max	min	max	type
131091	131091	1.311 ms	1.311 ms	131092	131092	no

Summary					
Name	BRAM_18K	DSP	FF	LUT	URAM
DSP	-	32	-	-	-
Expression	-	-	0	2080	-
FIFO	-	-	-	-	-
Instance	-	0	0	1280	-
Memory	-	-	-	-	-
Multiplexer	-	-	-	1060	-
Register	-	-	1491	-	-
Total	0	32	1491	4420	0
Available	4320	6840	2364480	1182240	960
Available SLR	1440	2280	788160	394080	320
Utilization (%)	0	~0	~0	~0	0
Utilization SLR (%)	0	1	~0	1	0

### Ερώτημα 3:

Για το ερώτημα αυτό τρέξαμε C/RTL cosimulation

```
$finish called at time : 1311095 ns :
#0: 5977397 1013998
#1: 6 Test Passed
#2: 7 INFO: [SIM 1]
```

Τις πληροφορίες για το latency τις βλέπουμε στο προηγούμενο ερώτημα.

### Ερώτημα 4:

Για την βελτιστοποίηση της matrix\_multiply\_hw χρησιμοποιήσαμε την εντολή #pragma ARRAY\_PARTITION pragmas διαμερίζοντας τους πίνακες A και B κατά συγκεκριμένες διαστάσεις. Ο πλήρης διαμερισμός χωρίζει αυτούς τους πίνακες σε μεμονωμένα registers, κάτι που επιταχύνει τους υπολογισμούς επιτρέποντας παραλληλία. Επιπλέον, την εντολή #pragma HLS ARRAY\_PARTITION variable=C complete dim=2: διαμερίζει μερικώς τον πίνακα C για να επιτρέψει παράλληλη πρόσβαση στις στήλες. Τέλος, την εντολή #pragma HLS UNROLL μέσα στον εσωτερικό βρόχο ξετυλίγει τον βρόχο για το k, αυξάνοντας την παραλληλία, καθώς επιτρέπει ταυτόχρονους υπολογισμούς για στοιχεία των A και B.

```
#include <stdint.h>

#define LM 6 // Define  $\underline{l}_m$  (for example,  $\underline{l}_m = 6 \rightarrow m = 2^6 = 64$ )
#define LN 6 // Define  $\underline{l}_n$  (for example,  $\underline{l}_n = 6 \rightarrow n = 2^6 = 64$ )
#define LP 6 // Define  $\underline{l}_p$  (for example,  $\underline{l}_p = 6 \rightarrow p = 2^6 = 64$ )

#define M (1 << LM) //  $m = 2^{\underline{l}_m}$ 
#define N (1 << LN) //  $n = 2^{\underline{l}_n}$ 
#define P (1 << LP) //  $p = 2^{\underline{l}_p}$ 

void matrix_multiply_hw(
    const uint8_t A[N][M],
    const uint8_t B[M][P],
    uint32_t C[N][P]
) {
    #pragma HLS ARRAY_PARTITION variable=A complete dim=2
    #pragma HLS ARRAY_PARTITION variable=B complete dim=1
    #pragma HLS ARRAY_PARTITION variable=C complete dim=2

    // Perform matrix multiplication:  $C = A * B$ 
    for (int i = 0; i < N; ++i) {
        for (int j = 0; j < P; ++j) {
            C[i][j] = 0; // Initialize element C[i][j]
            for (int k = 0; k < M; ++k) {
                #pragma HLS UNROLL
                C[i][j] += (uint32_t)A[i][k] * (uint32_t)B[k][j];
            }
        }
    }
}
```

ii) Για την βέλτιστη λύση έχουμε:

iii) Για την επιτάχυνση έχουμε:

```

$finish called at time : 1205 ns

```