



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ

ΕΠΙΣΤΗΜΗ ΔΕΔΟΜΕΝΩΝ ΚΑΙ ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

Παράλληλες Αρχιτεκτονικές Υπολογισμού για Μηχανική Μάθηση

ΕΡΓΑΣΙΑ

Επιτάχυνση εκπαίδευσης νευρωνικού δικτύου σε αρχιτεκτονικές
κοινής μνήμης με OpenMP και CUDA

Βλάχος Ιωάννης

AM: 03400087 ioannisvlahos@mail.ntua.gr

Ορφανουδάκης Φίλιππος Σκόβελεφ

AM: 03400107 philipposorfanoudakis@mail.ntua.gr

Πανόπουλος Ιωάννης

AM: 03003050 ioannispanop@mail.ntua.gr

Παπαγεωργίου Δημήτριος

AM: 03400088 dimitriospapageorgiou@mail.ntua.gr

Αθήνα, Ιούνιος 2021

Περιεχόμενα

1	Εισαγωγή	2
1.1	CPU	2
1.2	GPU	2
1.3	Κώδικας προς Βελτιστοποίηση	2
1.4	Έλεγχος Ορθότητας	3
2	Κώδικας και Υλοποίηση	4
2.1	OpenMP	4
2.2	OpenBLAS	5
2.3	CUDA	7
2.4	cuBLAS	11
3	Πειράματα και Αποτελέσματα	13
3.1	Παραλληλοποίηση σε CPU	13
3.2	Παραλληλοποίηση σε GPU	14
3.3	Σύγκριση CPU - GPU	15

1. Εισαγωγή

Στα πλαίσια της άσκησης μας ζητείται να εκπαιδεύσουμε ένα νευρωνικό δίκτυο με 3 κρυφά επίπεδα, με στόχο την κατηγοριοποίηση χειρόγραφων ψηφίων του συνόλου δεδομένων MNIST. Η εκπαίδευση ενός νευρωνικού δικτύου, περιλαμβάνει αρκετούς πολλαπλασιασμούς πινάκων. Όπως γνωρίζουμε η απλή υλοποίηση του πολλαπλασιασμού πινάκων έχει $O(n^3)$ πολυπλοκότητα, και αυτό οφείλεται στα εμφωλευμένα loops που περιλαμβάνει. Επίσης η απλή υλοποίηση οδηγεί σε αρκετά cache misses τα οποία επιβαρύνουν την απόδοση και τον χρόνο εκτέλεσης. Για αυτούς τους λόγους, βασικό ζητούμενο της άσκησης είναι η αποδοτικότερη υλοποίηση του αλγορίθμου πολλαπλασιασμού πινάκων (GEneral Matrix Multiply, GEMM) με κύριο εργαλείο, την παραλληλοποίηση. Πιο συγκεκριμένα, θα ασχοληθούμε με την βελτιστοποίηση του αλγορίθμου τόσο σε CPU όσο και σε GPU.

1.1 CPU

Η παραλληλοποίηση σε CPU γίνεται είτε με το προγραμματιστικό μοντέλο OpenMP, ή με τη χρήση της βιβλιοθήκης OpenBLAS. Το μηχάνημα που χρησιμοποιήθηκε αποτελείται από δύο επεξεργαστές Intel Xeon X5650 (6 πυρήνες + 2 hyper-threads ανά πυρήνα, συνολικά 12 πυρήνες + 24 hyper-threads), στο οποίο έχουμε πρόσβαση από την ουρά termis.

1.2 GPU

Το παρόν κομμάτι της αναφοράς αφορά στην παραλληλοποίηση της διαδικασίας εκπαίδευσης ενός νευρωνικού δικτύου σε GPUs. Το κομμάτι που παραλληλοποιείται μέσω CUDA ή cuBLAS είναι οι πολλαπλασιασμοί πινάκων, ενώ τα πειράματα εκτελούνται πάνω στο μηχάνημα dungani που αποτελείται από έναν επεξεργαστή Intel i7-4820K (4 πυρήνες + 2 hyper-threads ανά πυρήνα, συνολικά 8 hyper-threads) και μία κάρτα γραφικών NVIDIA Tesla K40, στο οποίο έχουμε πρόσβαση από την ουρά serial.

1.3 Κώδικας προς Βελτιστοποίηση

Μας δίνονται 3 τμήματα κώδικα, καθένα εκ των οποίων περιέχει μιας μορφής πολλαπλασιασμό πινάκων A και B με την naïve υλοποίηση τους. Η naïve υλοποίηση αποτελείται από 3 loops. Το πρώτο loop σκανάρει τις γραμμές του A, το δεύτερο σκανάρει τις στήλες του B και το τρίτο παίρνει κάθε στοιχείο της στήλης του B. Πιο συγκεκριμένα τα τμήματα κώδικα αναφέρονται στον φάκελο που μας δίνεται ως **dgemm**, **dgemm_ta** και **dgemm_tb**. Πιο συγκεκριμένα :

- **dgemm**: $A \cdot B$
- **dgemm_ta**: $A^T \cdot B$
- **dgemm_tb**: $A \cdot B^T + C$

1.4 Έλεγχος Ορθότητας

Πριν προχωρήσουμε στην ανάλυση και επεξήγηση της μεθοδολογίας που εφαρμόσαμε, να αναφέρουμε πως όλα τα προγράμματα μας, έχουν περάσει τον έλεγχο ορθότητας που έχει τεθεί και παράγουν σωστά αποτελέσματα.

2. Κώδικας και Υλοποίηση

2.1 OpenMP

Όπως αναλύσαμε και στην προηγούμενη ενότητα τα τμήματα κώδικα προς βελτιστοποίηση αποτελούνται από εμφωλευμένα for loops. Η οδηγία στο OpenMP που παραλληλοποιεί τα for loops είναι η **#pragma omp parallel for**. Επίσης, καθώς παραλληλοποιούμε τα for loops θέλουμε για κάθε στοιχείο του νέου πίνακά μας, που είναι το γινόμενο, να συσσωρεύουμε όλα τα επιμέρους αθροίσματα. Κάθε thread που έχει αναλάβει τον υπολογισμό ενός από τα πολλά αθροίσματα που χρειάζεται για τον υπολογισμό ενός στοιχείου γινομένου C_{ij} , προσθέτει το αποτέλεσμα σε μια κοινή μεταβλητή. Για αυτό το σκοπό το OpenMP προσφέρει την οδηγία `reduction(+:sum)`. Εμείς θα πειραματιστούμε με αυτές τις 2 εντολές για να βρούμε τελικά την υλοποίηση που μας επιστρέφει τους καλύτερους χρόνους.

Οι δοκιμές που πραγματοποιήθηκαν είναι οι εξής :

- `#pragma omp parallel collapse(3)`, δηλαδή unroll τα εμφωλευμένα loops και έπειτα παραλληλοποίηση τους
- `#pragma omp parallel for shared(i,j) private(k) reduction(+:sum)`, δηλαδή παραλληλοποίηση του 2ου εσωτερικού loop.
- `#pragma omp parallel for shared(A, B, C,i) private(j, k, sum)`, δηλαδή παραλληλοποίηση του 1ου εσωτερικού for loop. Πλέον κάθε thread κοιτάει μια στήλη του πίνακα B και επομένως αναλαμβάνει μια στήλη του γινομένου.
- `#pragma omp parallel for shared(A, B, C) private(i, j, k, sum)`, πλέον κάθε thread κοιτάει μια γραμμή του A και επομένως αναλαμβάνει μια γραμμή του γινομένου.

Από αυτές τις 4 δοκιμές τα καλύτερα αποτελέσματα από άποψη χρόνου, σημειώθηκαν στην τελευταία υλοποίηση, δηλαδή στην παραλληλοποίηση γραμμής. Οι αντίστοιχες υλοποιήσεις είναι οι εξής :

dgemm:

```
1 #pragma omp parallel for shared(A, B, C) private(i, j, k, sum)
2   for (i = 0; i < M; i++)
3   {
4       for (j = 0; j < N; j++)
5       {
6           sum = 0.;
7           for (k = 0; k < K; k++)
8               sum += A[i * K + k] * B[k * N + j];
9           C[i * N + j] = sum;
10      }
11  }
```

dgemm.ta:

```
1 #pragma omp parallel for shared(A, B, C) private(i, j, k, sum)
2   for (i = 0; i < M; i++)
3   {
```

```

4   for (j = 0; j < N; j++)
5   {
6       sum = 0.;
7       for (k = 0; k < K; k++)
8           sum += A[k * M + i] * B[k * N + j];
9       C[i * N + j] = sum;
10  }
11  }

```

dgemm.tb:

```

1  #pragma omp parallel for shared(A, B, C, D) private(i, j, k, sum)
2  for (i = 0; i < M; i++)
3  {
4      for (j = 0; j < N; j++)
5      {
6          sum = 0.;
7          for (k = 0; k < K; k++)
8              sum += A[i * K + k] * B[j * K + k];
9          D[i * N + j] = sum + C[i * N + j];
10     }
11 }

```

Ως shared θέτουμε μόνο τους πίνακες A,B,C,D, οι οποίοι πρέπει να υπάρχουν σε κάθε thread και στη συνέχεια τις υπόλοιπες μεταβλητές τις θέτουμε ως private έτσι ώστε να αποφύγουμε τα conflicts και κάθε thread να αναλάβει αποκλειστικά μια γραμμή του γινομένου.

2.2 OpenBLAS

Για την αξιοποίηση της βιβλιοθήκης OpenBLAS, χρησιμοποιούμε την συνάρτηση cblas_dgemm. Αυτή δέχεται τα εξής ορίσματα:

```

◆ cblas_dgemm()

void cblas_dgemm ( CBLAS_LAYOUT    layout,
                  CBLAS_TRANSPOSE TransA,
                  CBLAS_TRANSPOSE TransB,
                  const int         M,
                  const int         N,
                  const int         K,
                  const double      alpha,
                  const double *    A,
                  const int         lda,
                  const double *    B,
                  const int         ldb,
                  const double      beta,
                  double *          C,
                  const int         ldc
                )

```

Figure 2.1: OpenBlas Documentation

Αν ο πίνακας που δίνεται ως όρισμα δεν είναι ανάστροφος, τότε δίνεται το όρισμα CblasNoTrans, αλλιώς δίνεται το όρισμα CblasTrans. Το leading dimension κάθε ορίσματος είναι ο αριθμός των στηλών του πίνακα στη μορφή που δίνεται, δηλαδή είτε σε κανονική μορφή είτε ανάστροφος. Πιο συγκεκριμένα, οι υλοποιήσεις μας είναι οι παρακάτω :

dgemm:

```

1 cblas_dgemm(CblasRowMajor,
2             CblasNoTrans,
3             CblasNoTrans,
4             M,
5             N,
6             K,
7             1, // alpha = 1
8             A,
9             K, // Leading dimension (number of columns) of A = M x K
10            B,
11            N, // Leading dimension (number of columns) of B = K x N
12            0, // beta = 0
13            C,
14            N // Leading dimension (number of columns) of C = M x N
15 );

```

dgemm.ta:

```

1 cblas_dgemm(CblasRowMajor,
2             CblasTrans,
3             CblasNoTrans,
4             M,
5             N,
6             K,
7             1, // alpha = 1
8             A,
9             M, // Leading dimension (number of columns) of A = M x K
10            B,
11            N, // Leading dimension (number of columns) of B = K x N
12            0, // beta = 0
13            C,
14            N // Leading dimension (number of columns) of C = M x N
15 );

```

dgemm.tb:

```

1 cblas_dgemm(CblasRowMajor,
2             CblasNoTrans,
3             CblasTrans,
4             M,
5             N,
6             K,
7             1, // alpha = 1
8             A,
9             K, // Leading dimension (number of columns) of A = M x K
10            B,
11            K, // Leading dimension (number of columns) of B = K x N
12            1, // beta = 1
13            C,
14            N // Leading dimension (number of columns) of C = M x N
15 );
16 memcpy(D, C, M * N * sizeof(double));

```

Στην τελευταία υλοποίηση ο B^T έχει διαστάσεις $N \times K$, άρα το leading term θα είναι ο αριθμός των στηλών δηλαδή K . Συγκεκριμένα η τελευταία υλοποίηση απαιτεί την εξής πράξη πινάκων :

$D = A \cdot B^T + C$. Η `cblas_dgemm` αποθηκεύει το αποτέλεσμα της στον πίνακα C, και στο τέλος αντιγράφουμε το αποτέλεσμα μας στον πίνακα D.

2.3 CUDA

Η υλοποίηση της άσκησης αποτελείται από την χρήση κοινής μνήμης στο περιβάλλον CUDA. Παρακάτω εξετάζονται οι ιδέες που χρησιμοποιήθηκαν.

Στην υλοποίηση μας συμμετέχουν πολλά threads για τον υπολογισμό ενός αποτελέσματος του πίνακα C. Συγκεκριμένα, πολλά threads υπολογίζουν τα μερικά αθροίσματα ενός tile των πινάκων εισόδου και εξόδου. Όμως ένα thread τη φορά είναι επιφορτισμένο για την ανάθεση της τελικής τιμής σε κάθε κελί του C πίνακα που καλούμαστε να υπολογίσουμε.

```
1 int row = blockIdx.y * blockDim.y + threadIdx.y;
2 int col = blockIdx.x * blockDim.x + threadIdx.x;
3 int local_tid = threadIdx.x;
4 int local_tidy = threadIdx.y;
```

Κάθε block περιέχει μέρος των πινάκων A, B (tile) και εντός αυτού θα γίνεται ο υπολογισμός των μερικών αθροισμάτων για ένα μέρος του συνολικού αθροίσματος των θέσεων $C[i][j]$ του εν λόγω κομματιού.

Σε ό,τι αφορά το kernel launch configuration, έχουμε:

- BLOCK SIZE = 16
- grid μεγέθους $N \times M$ blocks
- κάθε block είναι διδιάστατο μεγέθους $BLOCK_SIZE \times BLOCK_SIZE$
- Η κοινή μνήμη κάθε block έχει μέγεθος $2 \times BLOCK_SIZE \times BLOCK_SIZE \times \text{sizeof}(\text{double})$, ώστε να χωράνε τα περιεχόμενα των πινάκων A, B για ένα tile τύπου double.

```
1 dim3 block(BLOCK_SIZE, BLOCK_SIZE);
2 dim3 grid((N + BLOCK_SIZE - 1) / BLOCK_SIZE,
3           (M + BLOCK_SIZE - 1) / BLOCK_SIZE);
4 size_t shmem_size = 2 * BLOCK_SIZE * BLOCK_SIZE * sizeof(double);
5 dgemm_tb_optimized<<<grid, block, shmem_size>>>(A, B, C, D, M, N, K,
6           BLOCK_SIZE);
7 checkCudaErrors(cudaPeekAtLastError());
8 checkCudaErrors(cudaDeviceSynchronize());
```

Στη συνέχεια δουλεύουμε στην 3η διάσταση μεγέθους K, στην οποία υπολογίζονται τα αθροίσματα. Σε κάθε block φορτώνεται ένα μέρος του πίνακα A και πίνακα B κατά γραμμή και κατά στήλη αντίστοιχα. Διαπερνάμε με τόσα blocks (που ισοδυναμούν με tiles) όσα χρειάζεται για να υπολογιστεί το άθροισμα των γινομένων. Σε κάθε block φορτώνουμε στην κοινή μνήμη στοιχεία των σειρών του A και των στηλών του B που θέλουμε. Έπειτα, υπολογίζουμε τα μερικά αθροίσματα για τα περιεχόμενα που υπάρχουν σε ένα block μόνο για την **κοινή μνήμη**. Πριν και μετά από αυτόν τον υπολογισμό συγχρονίζουμε τα threads ώστε να μην έχουμε προβλήματα στην προσπέλαση της μνήμης από άλλα threads που ενδεχομένως δεν έχουν ολοκληρώσει τους υπολογισμούς τους. Τελικά καλύπτουμε επαναληπτικά όλη την κοινή διάσταση των πινάκων, κατά την οποία πολλαπλασιάζουμε.

```
1 #pragma unroll
2 for (int current_block = 0; current_block < (K - 1) / SHARED_MEM_DIM +
3     1; current_block++)
```



```

3 {
4     // Column j of matrix A (sliding offset per tile)
5     j = current_block * blockDim.x + threadIdx.x;
6     // Row i of matrix B (sliding offset per tile)
7     i = current_block * blockDim.y + threadIdx.y;
8     // Load A[row][j] to shared mem from global mem
9     partials_A[local_tidy * SHARED_MEM_DIM + local_tidix] = A[row * K + j];
10    // Load B[i][col] to shared mem from global mem
11    partials_B[local_tidy * SHARED_MEM_DIM + local_tidix] = B[i * N + col];
12    // Synchronize before computation so that shared memory is full !
13    __syncthreads();
14    // Calculate **partial** dot products ONLY from shared memory !
15    // In a single tile, SHARED_MEM_DIM * SHARED_MEM_DIM partial
16    // results will be computed.
17    // ( & unroll loop for performance)
18    #pragma unroll
19    for (int k = 0; k < SHARED_MEM_DIM; k++)
20    {
21        // No shared memory bank conflict!
22        temp += partials_A[local_tidy * SHARED_MEM_DIM + k] * partials_B[k *
SHARED_MEM_DIM + local_tidix];
23    }
24    // Synchronize before next iteration so that previous
25    // values of shared memory are not overwritten BEFORE
26    // another thread has not used them yet!
27    // (for whatever reason)
28    __syncthreads();
29 }

```

Τέλος, αφού έχουμε διατρέξει όλα τα απαραίτητα blocks για να υπολογίσουμε K γινόμενα και να τα αθροίσουμε, αποθηκεύουμε το αποτέλεσμα στο αντίστοιχο κελί του C στην global memory.

```

1 C[row * N + col] = temp;

```

Για τις παραλλαγές του πολλαπλασιασμού πινάκων, αρκεί να αλλάζουμε τον τρόπο που διατρέχουμε τους πίνακες A , B ανάλογα με τον πολλαπλασιασμό που θέλουμε να εκτελέσουμε. Για παράδειγμα, στην περίπτωση του $C = A^T B$, πρέπει να διατρέξουμε και τους 2 πίνακες ανά γραμμή, ώστε ο A να θεωρηθεί ανάστροφος. Αλλάζουμε έτσι τις πρώτες γραμμές του εξωτερικού for loop. Η υλοποίηση λοιπόν για το dgemm_ta είναι :

```

1 extern __shared__ double partials[];
2 // Shared storage for a single tile of A
3 double *partials_A = partials;
4 // Shared storage for a single tile of B
5 double *partials_B = partials_A + SHARED_MEM_DIM * SHARED_MEM_DIM;
6
7 // Get row and column of **output C**
8 // (We utilize one thread per output value of C)
9 int row = blockIdx.y * blockDim.y + threadIdx.y;
10 int col = blockIdx.x * blockDim.x + threadIdx.x;
11 int local_tidix = threadIdx.x;
12 int local_tidy = threadIdx.y;
13
14 if (row < M && col < N)
15 {
16     double temp = 0;
17     int i, j;
18
19     // For all "imaginary" blocks (tiles) that are required
20     // to traverse dimension K (dimension of dot product)
21     // ( & unroll loop for performance)

```

```

22 #pragma unroll
23   for (int current_block = 0; current_block < (K - 1) / SHARED_MEM_DIM +
24       1; current_block++)
25   {
26       // Row j of matrix A (sliding offset per tile)
27       j = current_block * blockDim.x + threadIdx.x;
28       // Row i of matrix B (sliding offset per tile)
29       i = current_block * blockDim.y + threadIdx.y;
30
31       // Load A[j][row] to shared mem from global mem
32       // row === column of transposed matrix A
33       partials_A[local_tidy * SHARED_MEM_DIM + local_tididx] = A[j * M + row];
34
35       // Load B[i][col] to shared mem from global mem
36       partials_B[local_tidy * SHARED_MEM_DIM + local_tididx] = B[i * N + col];
37
38       // Synchronize before computation
39       // so that shared memory is full !
40       __syncthreads();
41
42       // Calculate **partial** dot products ONLY from shared memory !
43       // In a single tile, SHARED_MEM_DIM * SHARED_MEM_DIM partial
44       // results will be computed.
45       // ( & unroll loop for performance)
46 #pragma unroll
47       for (int k = 0; k < SHARED_MEM_DIM; k++)
48       {
49
50           // No shared memory bank conflict!
51           temp += partials_A[local_tidy * SHARED_MEM_DIM + k] * partials_B[k *
52               SHARED_MEM_DIM + local_tididx];
53       }
54
55       // Synchronize before next iteration so that previous
56       // values of shared memory are not overwritten BEFORE
57       // another thread has not used them yet!
58       // (for whatever reason)
59       __syncthreads();
60   }
61
62   // Now that all tiles in the dimension of dot product have been
63   // **accumulated** on 'temp', we can store one result per thread to
64   // output array.
65
66   // Store result in C
67   C[row * N + col] = temp;

```

Και για το dgemm_tb :

```

1 extern __shared__ double partials[];
2 // Shared storage for a single tile of A
3 double *partials_A = partials;
4 // Shared storage for a single tile of B
5 double *partials_B = partials_A + SHARED_MEM_DIM * SHARED_MEM_DIM;
6
7 // Get row and column of **output C**
8 // (We utilize one thread per output value of C)
9 int row = blockIdx.y * blockDim.y + threadIdx.y;
10 int col = blockIdx.x * blockDim.x + threadIdx.x;
11 int local_tididx = threadIdx.x;

```

```

12  int local_tidy = threadIdx.y;
13
14  if (row < M && col < N)
15  {
16      double temp = 0;
17      int i, j;
18
19      // For all "imaginary" blocks (tiles) that are required
20      // to traverse dimension K (dimension of dot product)
21      // ( & unroll loop for performance)
22  #pragma unroll
23      for (int current_block = 0; current_block < (K - 1) / SHARED_MEM_DIM +
24          1; current_block++)
25      {
26          // Column j of matrix A (sliding offset per tile)
27          j = current_block * blockDim.x + threadIdx.x;
28          // Row i of matrix B (sliding offset per tile)
29          i = current_block * blockDim.y + threadIdx.y;
30
31          // Load A[row][j] to shared mem from global mem
32          partials_A[local_tidy * SHARED_MEM_DIM + local_tidix] = A[row * K + j];
33
34          // Load B[i][col] to shared mem from global mem
35          partials_B[local_tidy * SHARED_MEM_DIM + local_tidix] = B[col * K + i];
36
37          // Synchronize before computation
38          // so that shared memory is full !
39          __syncthreads();
40
41          // Calculate **partial** dot products ONLY from shared memory !
42          // In a single tile, SHARED_MEM_DIM * SHARED_MEM_DIM partial
43          // results will be computed.
44          // ( & unroll loop for performance)
45  #pragma unroll
46          for (int k = 0; k < SHARED_MEM_DIM; k++)
47          {
48
49              // No shared memory bank conflict!
50              temp += partials_A[local_tidy * SHARED_MEM_DIM + k] * partials_B[k *
51                  SHARED_MEM_DIM + local_tidix];
52
53              // Synchronize before next iteration so that previous
54              // values of shared memory are not overwritten BEFORE
55              // another thread has not used them yet!
56              // (for whatever reason)
57              __syncthreads();
58          }
59
60          // Now that all tiles in the dimension of dot product have been
61          // **accumulated** on 'temp', we can store one result per thread to
62          // output array.
63
64          // Store result in D
65          D[row * N + col] = C[row * N + col] + temp;
66      }

```

Για να τρέξουμε τις optimized versions, αλλάζουμε την μεταβλητή GEMM_OPTIMIZED σε 1 στο αρχείο Makefile.gpu, ώστε να γίνεται compile και run με τις βέλτιστες υλοποιήσεις.

2.4 cuBLAS

Εξετάζουμε στη συνέχεια το documentation του cuBLAS, όπως επίσης και την χρήση των έτοιμων συναρτήσεων.

```
cublasStatus_t cublasDgemm(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n, int k,
                           const double *alpha,
                           const double *A, int lda,
                           const double *B, int ldb,
                           const double *beta,
                           double *C, int ldc)
```

Figure 2.2: cublasDgemm

Στην βιβλιοθήκη cuBLAS οι πίνακες είναι αποθηκευμένοι κατά στήλες, ενώ οι πράξεις γίνονται σε ανάστροφους πίνακες. Έτσι πρέπει να φέρουμε τους πίνακες σε κατάλληλη μορφή για να δοθούν σωστά τα ορίσματα στην cublasDgemm.

Για τον απλό πολλαπλασιασμό $A \times B$, η ζητούμενη μορφή είναι η

$$C = A \times B = (A^T)^T \times (B^T)^T = (A')^T \times (B')^T = (B' \times A')^T$$

Έτσι λοιπόν σε αυτή την περίπτωση, τα ορίσματα θα είναι:

```
1 cublasHandle_t handle;
2 cublasCreate(&handle);
3 const double a1 = 1;
4 const double b1 = 0;
5 cublasDgemm(cublas_handle, CUBLAS_OP_N, CUBLAS_OP_N,
6             N, M, K, //m, n, k (because B*A)
7             &a1, //*alpha
8             B, N, //lda, rows of B^T if no trans, cols of B^T otherwise
9             A, K, //ldb
10            &b1, //*beta
11            C, N); //ldc, no transform on C^T
12 cublasDestroy(handle);
```

Δίνουμε δηλαδή τους πίνακες χωρίς να είναι ανεστραμμένοι, αλλά πρώτα τον B και μετά τον A και με ld τις στήλες των πινάκων.

Στην περίπτωση του dgemm_tα η μετατροπή των εισόδων είναι η εξής :

$$C = A^T \times B = A' \times (B^T)^T = A' \times (B')^T = (B' \times (A')^T)^T$$

και τα ορίσματα είναι τα εξής :

```
1 cublasHandle_t handle;
2 cublasCreate(&handle);
3 const double a2 = 1;
4 const double b2 = 0;
5 cublasDgemm(cublas_handle, CUBLAS_OP_N, CUBLAS_OP_T,
6             N, M, K, //m, n, k
7             &a2, //*alpha
8             B, N, //lda
9             A, M, //ldb
```

```

10         &b2,      /*beta
11         C, N);    //ldc
12 cublasDestroy(handle);

```

Για την περίπτωση του dgemm_tb αρχικά έχουμε την εξής μετατροπή :

$$D = A \times B^T + C = (A^T)^T \times B^T + C = (A')^T \times B' + C = ((B')^T \times A')^T + C$$

Για το πρώτο όρισμα της πρόσθεσης χρησιμοποιούμε την εξής υλοποίηση :

```

1 cublasHandle_t handle;
2 cublasCreate(&handle);
3 const double a3 = 1;
4 const double b3 = 0;
5 cublasDgemm(handle, CUBLAS_OP_T, CUBLAS_OP_N,
6             N, M, K, //m, n, k
7             &a3,     /*alpha
8             B, K,     //lda
9             A, K,     //ldb
10            &b3,      /*beta
11            D, N);    //ldc
12 cublasDestroy(handle);

```

Και τέλος για να προσθέσουμε τον D στο αποτέλεσμα χρησιμοποιούμε την cublasDgeam, η οποία είναι παρόμοια με την παραπάνω αλλά κάνει πρόσθεση με αντίστοιχους συντελεστές. Έτσι, έχοντας αποθηκεύσει το αποτέλεσμα στον D, προσθέτουμε σε αυτόν τον C.

```

1 cublasHandle_t handle1;
2 cublasCreate(&handle1);
3 const double a4 = 1;
4 const double b4 = 1;
5 cublasDgeam(handle1, CUBLAS_OP_N, CUBLAS_OP_N,
6             M, N, //m, n
7             &a4,   /*alpha
8             C, M, //lda
9             &b4,   /*beta
10            D, M, //ldb
11            D, M); //ldc
12 cublasDestroy(handle1);

```

3. Πειράματα και Αποτελέσματα

3.1 Παραλληλοποίηση σε CPU

Έχοντας τρέξει τις παραπάνω υλοποιήσεις, προκύπτουν οι παρακάτω χρόνοι αναλυτικά ανά μέθοδο και αριθμό νημάτων.

	n.threads	DGEMM	DGEMM_TA	DGEMM_TB	TOTAL
Serial	1	1046.15	1377.70	453.41	2891.63
OpenMP	1	664.15	1129.59	414.92	2224.86
	2	342.80	597.26	210.09	1167.09
	4	249.55	311.82	104.26	683.28
	6	242.27	222.67	72.22	554.93
	12	130.32	132.86	45.55	333.53
	24	250.90	145.63	41.76	466.15
OpenBLAS	1	67.06	83.58	79.02	243.45
	2	35.51	44.61	42.11	140.87
	4	18.80	24.21	22.84	84.72
	6	16.62	18.58	21.40	75.68
	12	11.46	12.73	15.16	59.55
	24	10.62	14.80	14.15	66.48

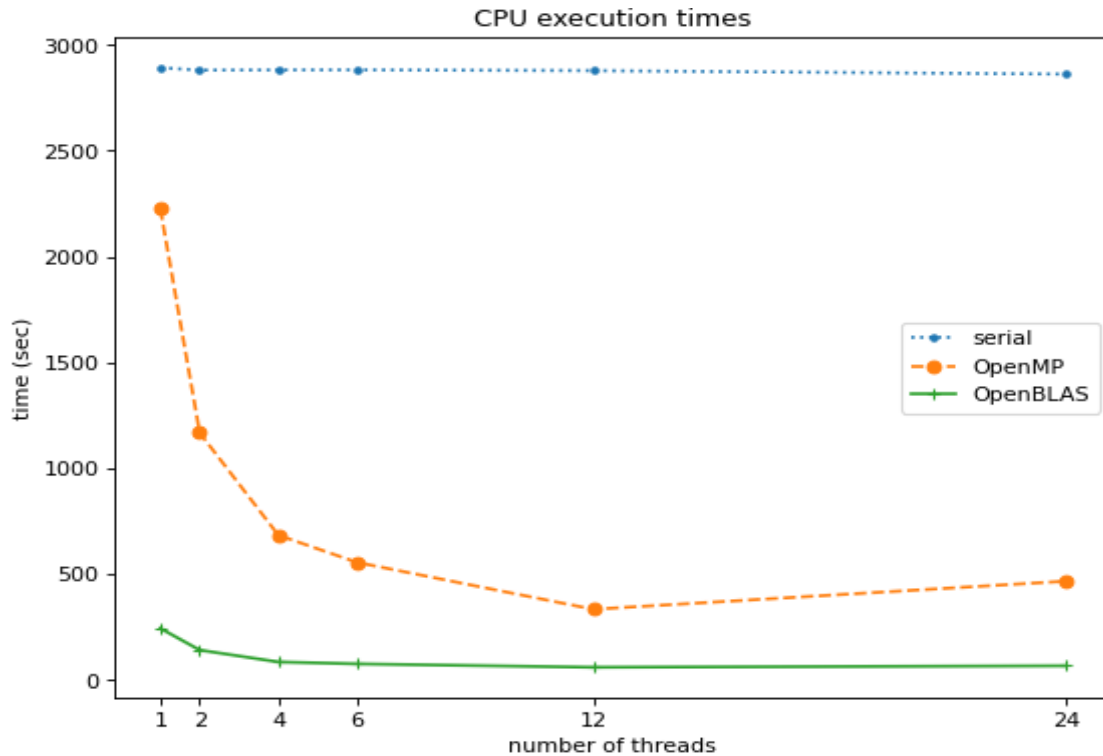


Figure 3.1: CPU χρόνοι εκτέλεσης

Διαιρώντας τους χρόνους της σειριακής υλοποίησης με τους χρόνους της κάθε μεθόδου λαμβάνουμε τα speed ups.

	n_threads	speed-up
Serial	1	1
OpenMP	1	1.299
	2	2.477
	4	4.231
	6	5.21
	12	8.669
	24	6.203
OpenBLAS	1	11,877
	2	20.526
	4	34.131
	6	38.20
	12	48.55
	24	43.496

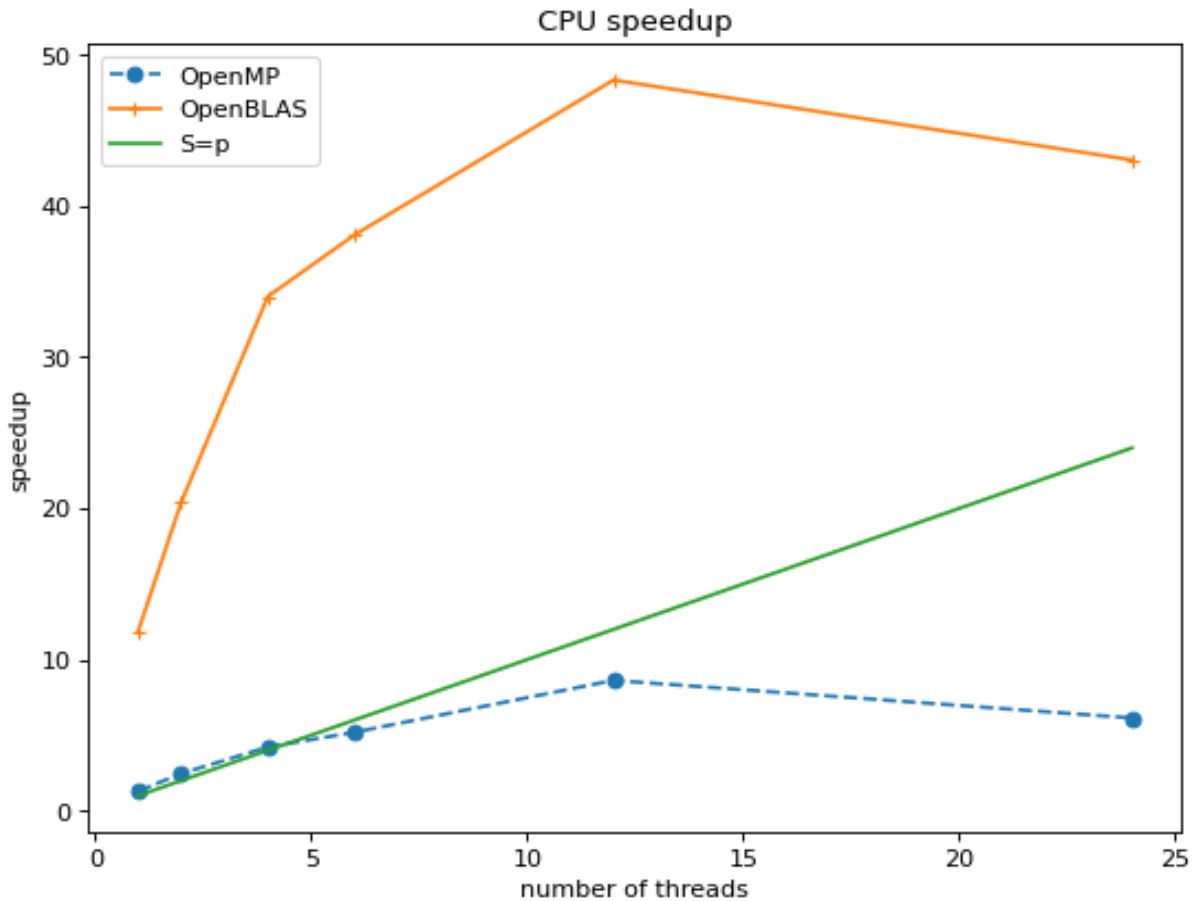


Figure 3.2: CPU speedups

3.2 Παραλληλοποίηση σε GPU

Οι παραπάνω υλοποιήσεις έδωσαν τους εξής χρόνους εκτέλεσης για τα ζητούμενα της εργασίας. Οι κατηγορίες είναι η naïve υλοποίηση, η υλοποίηση optimized CUDA και αυτή της βιβλιοθήκης cuBLAS.

	DGEMM	DGEMM_TA	DGEMM_TB	TOTAL
Naive	463.83	625.11	306.96	1406.92
Opt. CUDA	513.09	154.53	57.53	741.56
cuBLAS	512.17	154.64	57.56	740.72

Και το αντίστοιχο SpeedUp που λαμβάνουμε είναι :

- CUDA : 0.527
- cuBLAS : 0.526

3.3 Σύγκριση CPU - GPU

Αναλυτικότερα τα αποτελέσματα των πινάκων φαίνονται παρακάτω. Οι χρόνοι κάθε σεναρίου παρουσιάζονται ανά στήλη και ανά κατηγορία.

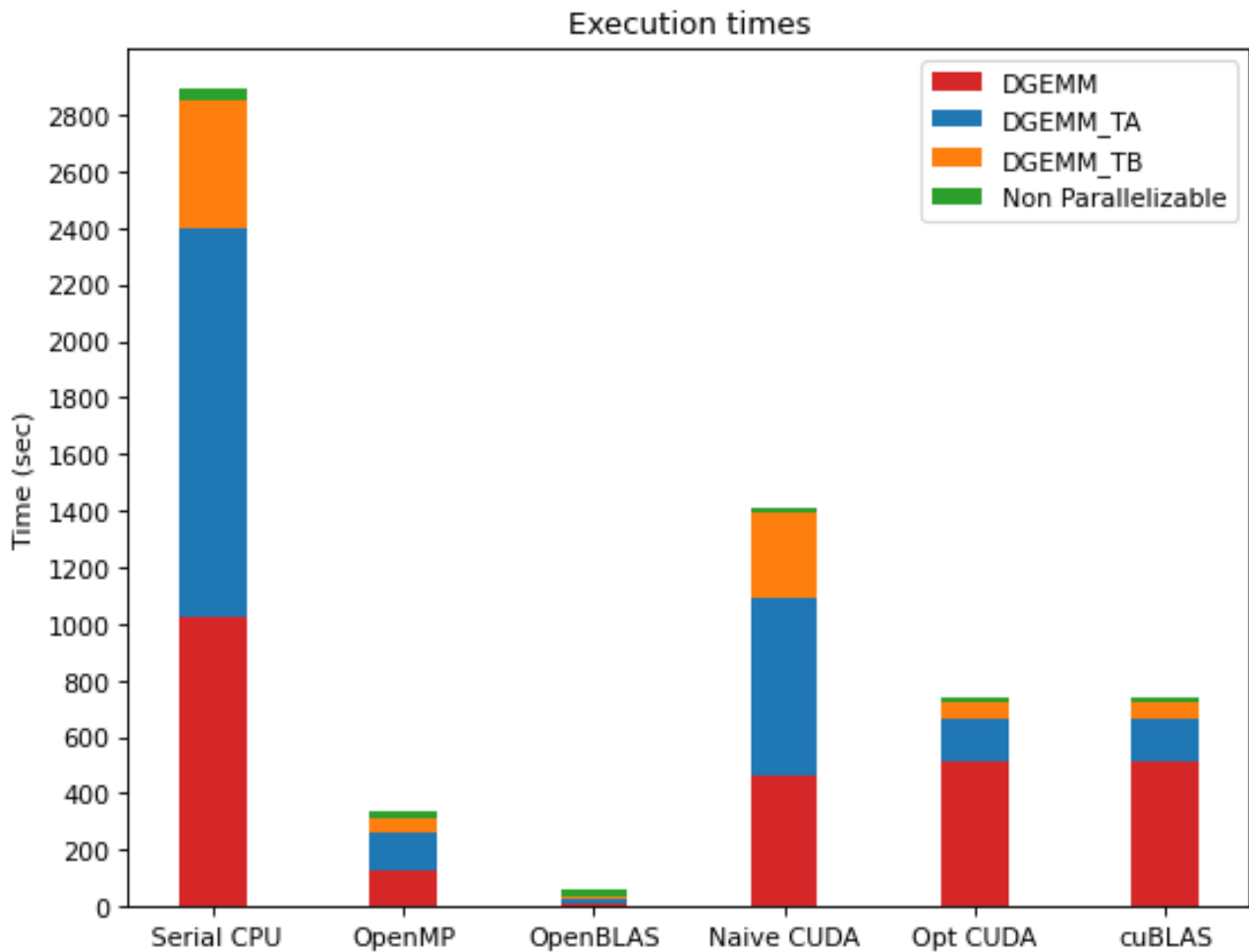


Figure 3.3: Χρόνοι εκτέλεσης

Συνολικά, παρατηρούμε ότι οι βελτιστοποιημένες εκτελέσεις σημειώνουν πολύ καλύτερους χρόνους συγκριτικά με τις αντιστοιχες Serial και Naive, επομένως ο στόχος της άσκησης επιτεύχθηκε. Όσον αφορά τις διαφορετικές πράξεις πινάκων, παρατηρούμε ότι ο DGEMM_TB και ακολούθως ο DGEMM_TA έχουν τους καλύτερους χρόνους εκτέλεσης και στις 4 περιπτώσεις παραλληλοποίησης. Επιπρόσθετα, ο DGEMM φαίνεται να αποδίδει χειρότερα στις

βελτιστοποιήσεις με CUDA και cuBLAS.

Τέλος, η βιβλιοθήκη OpenMP και OpenBLAS πετυχαίνουν καλύτερους χρόνους, με την υλοποίηση της έτοιμης βιβλιοθήκης να ξεχωρίζει.