

Τμήμα Μηχανικών Ηλεκτρονικών
Υπολογιστών και Πληροφορικής,
Πανεπιστήμιο Πατρών

Project

Λογισμικό και προγραμματισμός συστημάτων υψηλής επίδοσης

Σπύρος Κανκιάς 5317
Ταλάκ Μπαιράμ 5362

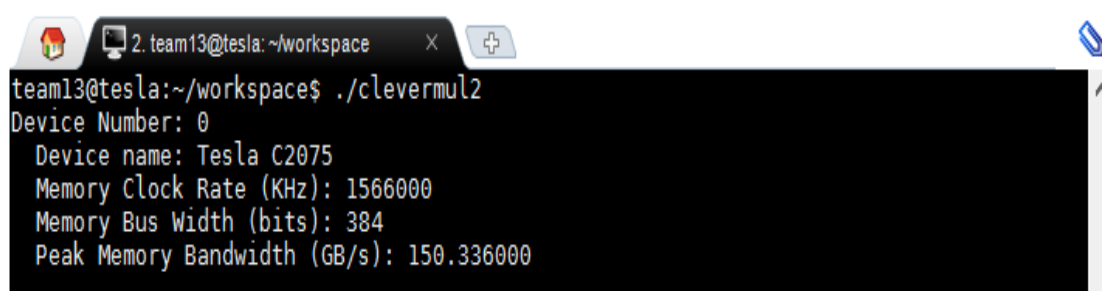
Εισαγωγή

Η nVidia θέλοντας να εκμεταλλευθεί την υπολογιστική ισχύ των επεξεργαστών γραφικών (*GPU*), σε συνδυασμό με την συνεχή ανάπτυξη του παράλληλου λογισμικού, εισήγαγε μια νέα αρχιτεκτονική που επιτρέπει την ανάπτυξη κοινών, μη γραφικών, εφαρμογών, τμήματα των οποίων ανατίθενται προς εκτέλεση στην *GPU*. Το ποια τμήματα της εφαρμογής θα ανατεθούν στην *GPU* έχει σχέση με την δυνατότητα παραλληλοποίησης τους και επιλέγονται από τον προγραμματιστή.

Η CUDA¹ είναι μια πλατφόρμα παράλληλου προγραμματισμού και ταυτόχρονα ένα προγραμματιστικό μοντέλο διεπαφής (*API*) που έχει δημιουργήσει η nVidia. Επιτρέπει στους μηχανικούς λογισμικού να χρησιμοποιήσουν τις δυνατότητες των *GPU* για να επιλύσουν προβλήματα γενικής χρήσεως (*GPGPU*²). Προκειμένου η CUDA να παραμείνει όσο το δυνατόν πιο προσιτή στους προγραμματιστές, δημιουργήθηκε ως μια επέκταση της πολύ δημοφιλούς γλώσσας *C* χρησιμοποιώντας σε γενικές γραμμές συντακτικό και έννοιες προγραμματισμού ίδιες με αυτήν. Σε αντίθεση με την *C*, όπου η εκτέλεση είναι συνήθως μονονηματική, η CUDA μας δίνει την δυνατότητα να ορίσουμε συναρτήσεις (*kernels*) οι οποίες εκτελούνται παράλληλα από έναν καθορισμένο αριθμό νημάτων (*multi-threaded execution*).

Στην παρούσα εργασία θα ασχοληθούμε με ένα πρόβλημα γραμμικής άλγεβρας, με τον υπολογισμό του γινομένου του ανάστροφου ενός μητρώου με το αρχικό μητρώο. Δηλαδή, το γινόμενο $A^T \cdot A$. Φυσικά ο υπολογισμός αυτός συνήθως αποτελεί τμήμα μόνο του συνολικού προγράμματος, ωστόσο, μπορεί να απαιτεί αρκετά σημαντικό ποσοστό από τον συνολικό χρόνο εκτέλεσης της εφαρμογής. Κατά συνέπεια μια αποδοτική υλοποίηση είναι απαραίτητη. Θα υλοποιήσουμε και θα μετρήσουμε την απόδοση της προαναφερθείσας πράξης στο προγραμματιστικό μοντέλο *CUDA*.

Παρακάτω παρουσιάζουμε κάποια βασικά στοιχεία της *GPU* που χρησιμοποιήθηκε για να υλοποιήσουμε την εργασία.



```
team13@tesla:~/workspace$ ./clevermul2
Device Number: 0
Device name: Tesla C2075
Memory Clock Rate (KHz): 1566000
Memory Bus Width (bits): 384
Peak Memory Bandwidth (GB/s): 150.336000
```

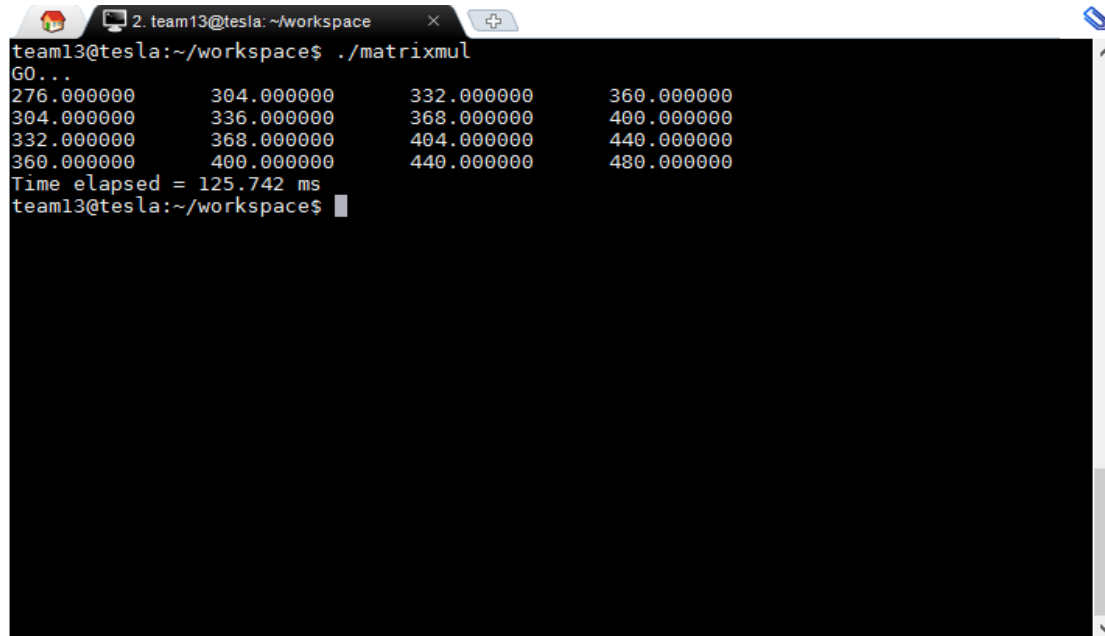
Εικόνα 1: Χαρακτηριστικά *GPU*

¹ Compute Unified Device Architecture

² General-Purpose computation on graphics processing units

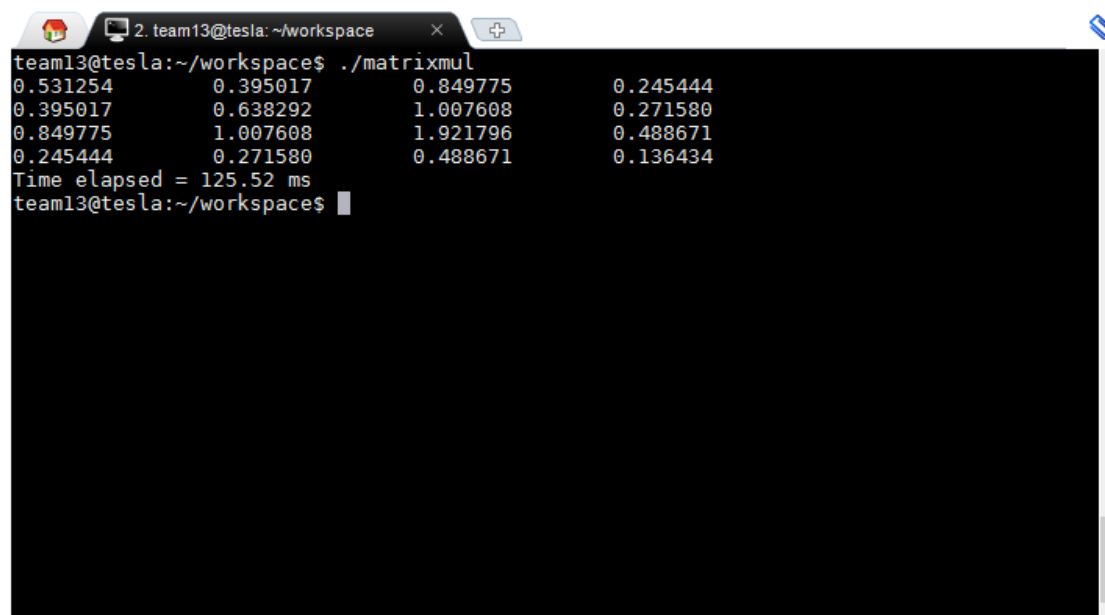
Ερώτημα 1

Σε αυτό το ερώτημα χρησιμοποιήσαμε την συνάρτηση της *CUDA cublasDgemm* για διάφορα μεγέθη μητρώων και μετρήσαμε τον χρόνο εκτέλεσης της πράξης. Παρακάτω χρησιμοποιήσαμε ένα γνωστό 4x4 μητρώο [1, 2, 3, 4, 5 κ.λ.π] και ένα τυχαίο μη τετραγωνικό μητρώο για να επιβεβαιώσουμε ότι λειτουργεί σωστά ο πολλαπλασιασμός.



```
team13@tesla:~/workspace$ ./matrixmul
G0...
276.000000      304.000000      332.000000      360.000000
304.000000      336.000000      368.000000      400.000000
332.000000      368.000000      404.000000      440.000000
360.000000      400.000000      440.000000      480.000000
Time elapsed = 125.742 ms
team13@tesla:~/workspace$
```

Εικόνα 2: Μητρώο 4x4

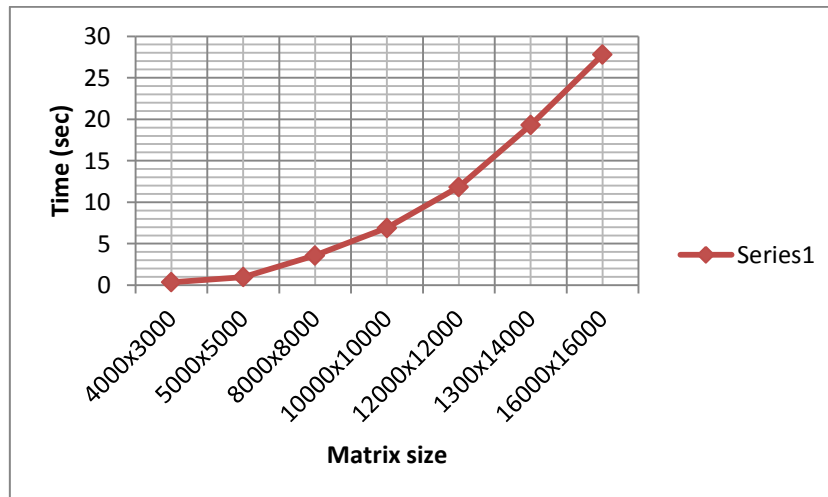


```
team13@tesla:~/workspace$ ./matrixmul
0.531254      0.395017      0.849775      0.245444
0.395017      0.638292      1.007608      0.271580
0.849775      1.007608      1.921796      0.488671
0.245444      0.271580      0.488671      0.136434
Time elapsed = 125.52 ms
team13@tesla:~/workspace$
```

Εικόνα 3: Μητρώο 3x4

Εδώ βρίσκονται οι χρόνοι εκτέλεσης των πράξεων για τα παρακάτω μεγέθη μητρώων.

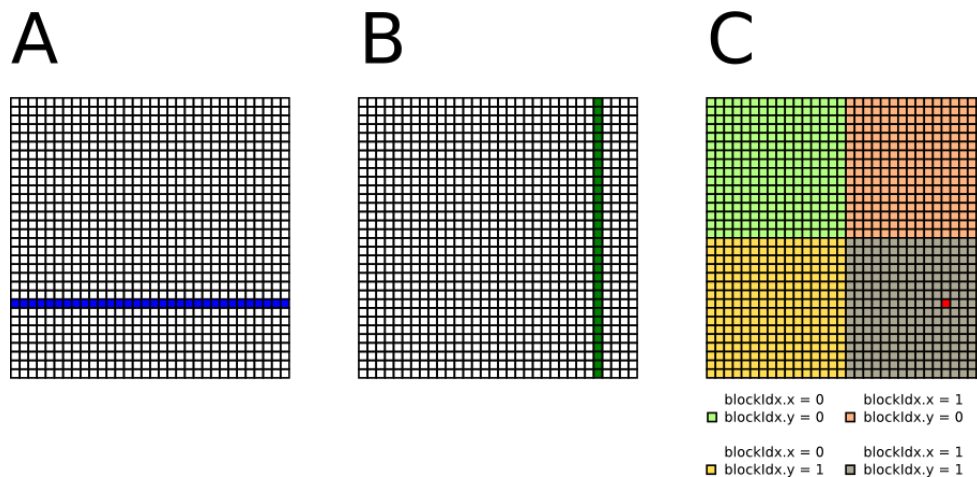
Size	4000x3000	5000x5000	8000x8000	10000x10000	12000x12000	13000x14000	16000x16000
Time(sec)	0,3	0,9	3,5	6,9	11,8	18,2	27,7



Εικόνα 4: Γραφική παράσταση του χρόνου συναρτήσει της εισόδου για την *cublasDgemm*

Ερώτημα 2

Σε αυτό το ερώτημα χρειάστηκε να υλοποιήσουμε την δική μας συνάρτηση υπολογισμού του γινομένου $C = A^T \cdot A$. Χρησιμοποιήσαμε τον πιο απλό τρόπο, δημιουργήσαμε μια συνάρτηση kernel (ονόματι *matrix_mul*) και αναθέσαμε σε κάθε νήμα να διαβάσει από μια γραμμή και μία στήλη του πίνακα A και να υπολογίζει ένα στοιχείο του πίνακα C. Έχουμε τρέξει τον κώδικα για τρία διαφορετικά μεγέθη block.



Εικόνα 5: Παρουσιάζει τον τρόπο που υπολογίζεται κάθε στοιχείο του πίνακα C

Εδώ παρουσιάζουμε δύο στιγμιότυπα εκτέλεσης

```
team13@tesla:~/workspace$ ./naivemul
Array size : ( 16000 X 16000 )
GPU will create : 1000.000000 blocks
GPU will create : 256 threads per block
Time elapsed = 240575 ms
team13@tesla:~/workspace$
```

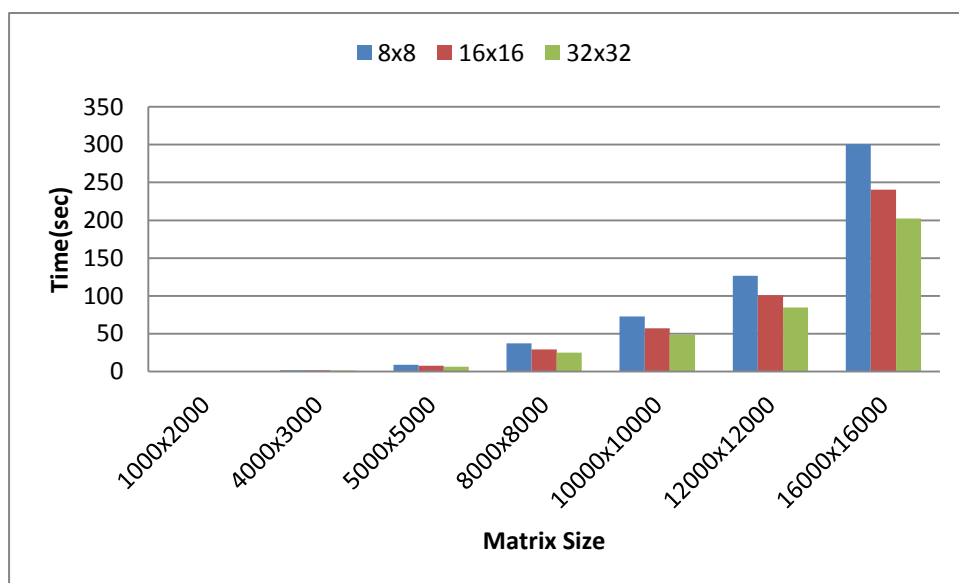
Εικόνα 6: Block size 16x16, Matrix size 16K x 16K

```
team13@tesla:~/workspace$ ./naivemul
Array size : ( 13000 X 14000 )
GPU will create : 438.000000 blocks
GPU will create : 1024 threads per block
Time elapsed = 130882 ms
team13@tesla:~/workspace$
```

Εικόνα 7: Block size 32x32, Matrix Size 13K x 14K

Παρακάτω βρίσκονται οι μετρήσεις που λάβαμε για τα επόμενα μεγέθη μητρώων. Παρατηρούμε ότι όσο μεγαλύτερο είναι το *block size* τόσο καλύτερη είναι και η απόδοση του αλγορίθμου μας.

Size	4000x3000	5000x5000	8000x8000	10000x10000	12000x12000	13000x14000	16000x16000
Block 8x8	1,92	8,9	37	72,9	126,6	200	300,6
Block 16x16	1,62	7,7	29	57,3	100,6	160	240,4
Block 32x32	1,33	6,1	25	48,8	84,8	134	202,2



Εικόνα 8: Γραφική παράσταση απόδοσης *naive_mul*

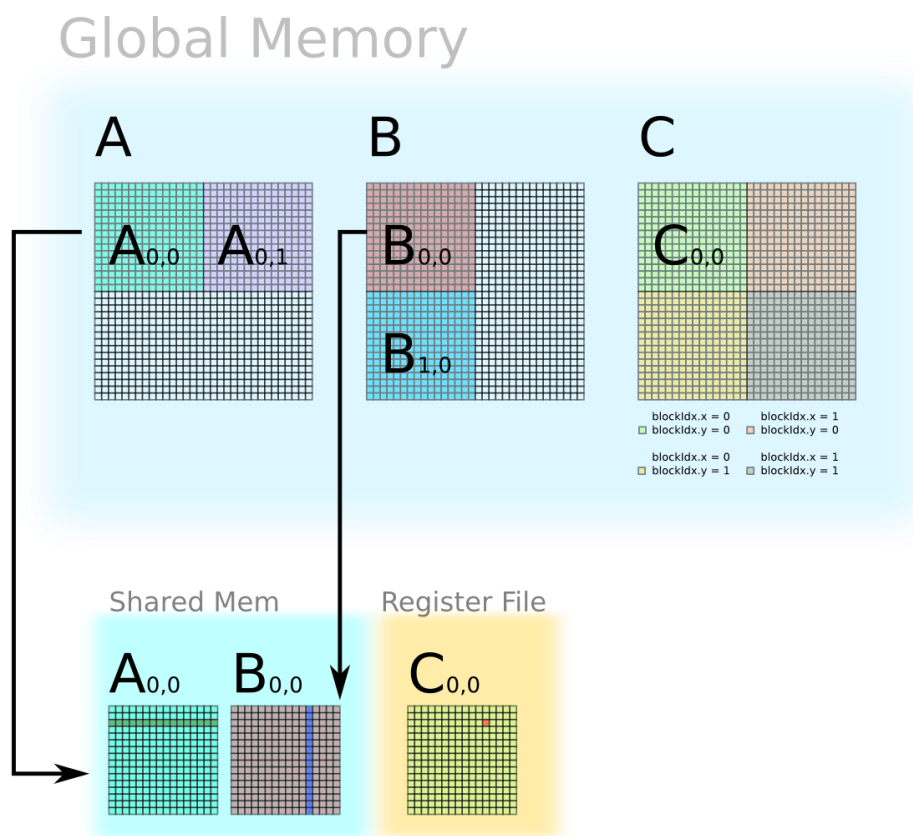
Ερώτημα 3

Σε αυτό το ερώτημα μας ζητήθηκε να βελτιστοποιήσουμε τον κώδικα μας χρησιμοποιώντας διάφορες γνωστές τεχνικές για τον πολλαπλασιασμό $A^T \cdot A$. Μερικές από τις πιο γνωστές τεχνικές είναι το tiling, coalesced memory access, prefetching, loop unrolling κ.λ.π. Έχουμε δημιουργήσει δύο εκδόσεις αυτού του ερωτήματος. Το κάναμε αυτό για να δείξουμε ότι σε κάθε επόμενη έκδοση χρησιμοποιούμε και έναν λίγο πιο βελτιστοποιημένο αλγόριθμο και έτσι πετυχαίνουμε καλύτερα αποτελέσματα σε ότι αφορά τους χρόνους εκτέλεσης.

➤ 1^η Έκδοση: **clevermulSM.cu** | 2^η Έκδοση: **clevermulSym.cu**

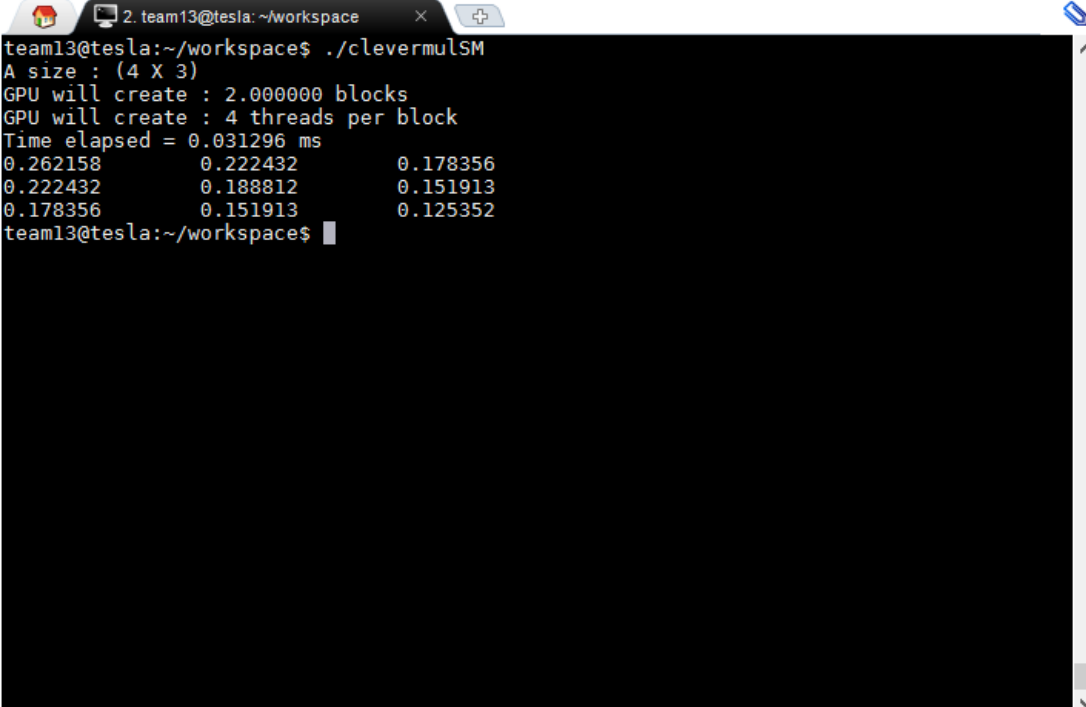
1^η Έκδοση

Η τεχνική που εφαρμόστηκε εδώ είναι η διαίρεση του αλγορίθμου σε πλακίδια, η οποία χρησιμοποιεί την κοινόχρηστη (shared memory) μνήμη για να μειώσει συνολικά τις προσπελάσεις προς την καθολική μνήμη. Χρησιμοποιήσαμε δύο πίνακες *Asub*, *Bsub*, οι οποίοι ορίστηκαν στην κοινή μνήμη και έχουν μέγεθος $[TILE_WIDTH][TILE_WIDTH]$ ο καθένας. Έπειτα διαβάσαμε τον πίνακα *A* που βρίσκεται στην global memory της GPU και αποθηκεύσαμε $TILE_WIDTH \times TILE_WIDTH$ στοιχεία στους πίνακες *Asub*, *Bsub*. Κάθε νήμα ουσιαστικά είναι υπεύθυνο σε κάθε φάση του αλγορίθμου να μεταφέρει και ένα στοιχείο από τον πίνακα *A* που βρίσκεται στην global memory στην κοινή μνήμη. Σε αυτήν την έκδοση του προγράμματος χρησιμοποιήθηκε η τεχνική *coalesced memory access* κατά την διάρκεια που προσπελαύναμε τον πίνακα *A* για να αποθηκεύσουμε τα στοιχεία του στον *Bsub* που βρίσκονται στην κοινή μνήμη.



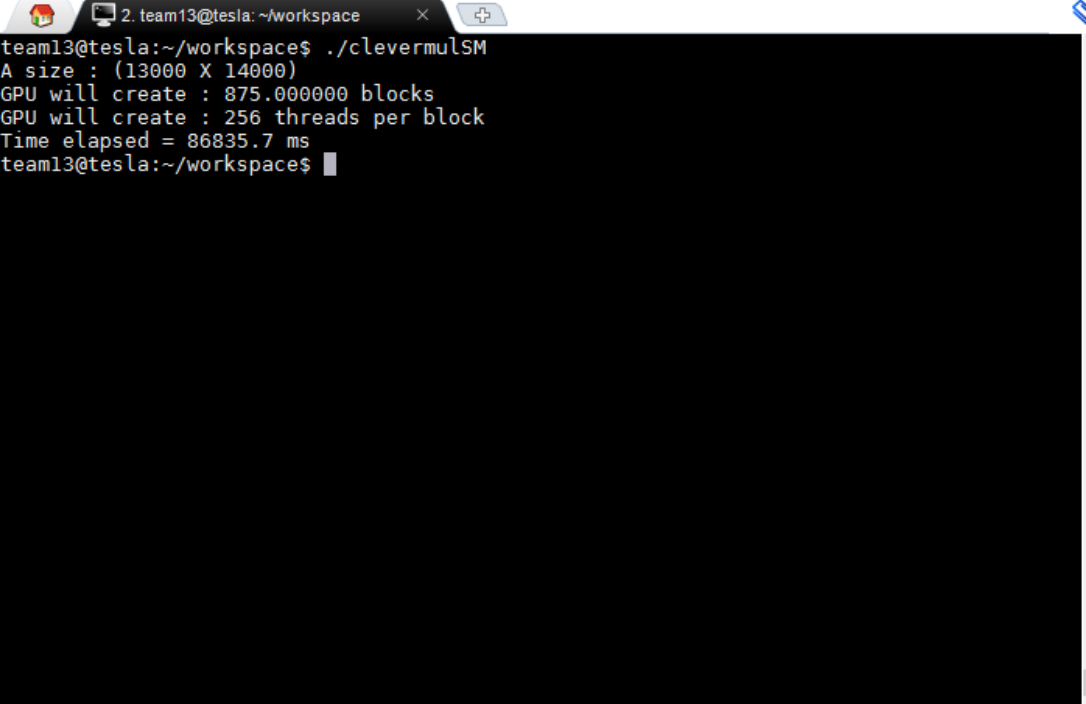
Εικόνα 9: Τρόπος υλοποίησης tiled αλγορίθμου

Εδώ παρουσιάζουμε δύο στιγμιότυπα εκτέλεσης για να αποδείξουμε την ορθότητα του αλγορίθμου. Στην αρχή τρέχουμε ένα μη τετραγωνικό μητρω και έπειτα ένα πολύ μεγάλο μητρώο(16K x 16K)



```
team13@tesla:~/workspace$ ./clevermulSM
A size : (4 X 3)
GPU will create : 2.000000 blocks
GPU will create : 4 threads per block
Time elapsed = 0.031296 ms
0.262158      0.222432      0.178356
0.222432      0.188812      0.151913
0.178356      0.151913      0.125352
team13@tesla:~/workspace$
```

Εικόνα 10: 4x3 matrix size & 2x2 block size

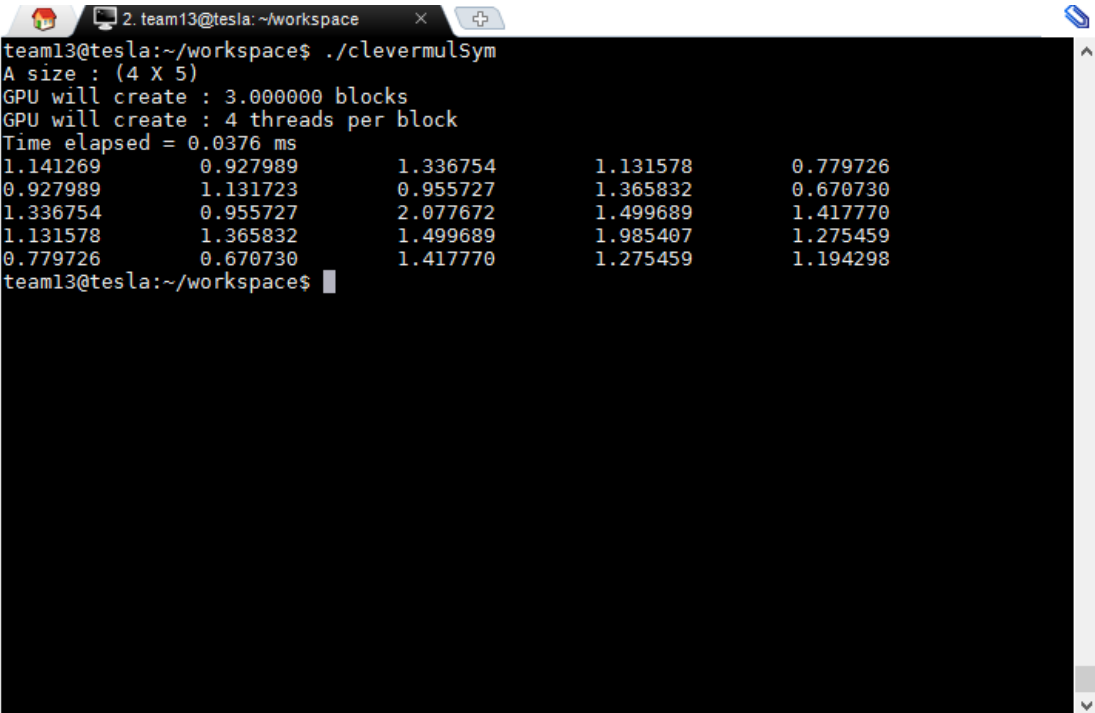


```
team13@tesla:~/workspace$ ./clevermulSM
A size : (13000 X 14000)
GPU will create : 875.000000 blocks
GPU will create : 256 threads per block
Time elapsed = 86835.7 ms
team13@tesla:~/workspace$
```

Εικόνα 11: 13K x 14K matrix size & 16x16 block size

Σε αυτήν την έκδοση του προγράμματος εκμεταλλευτήκαμε το γεγονός ότι έχουμε να πολλαπλασιάσουμε έναν πίνακα με τον ανάστροφο του ($A^T \cdot A$). Οπότε προκειμένου να βελτιώσουμε την απόδοση του αλγορίθμου μας υπολογίσαμε μόνο τα στοιχεία που βρίσκονται στην κύρια διαγώνιο και πάνω απο αυτήν και μετά τα αντιγράψαμε και κάτω απο την κύρια διαγώνιο. Το γεγονός ότι ο πίνακας που προκύπτει είναι συμμετρικός μας επιτρέπει να το κάνουμε. Τέλος χρησιμοποιήσαμε και την τεχνική του *loop unrolling*. Με αυτόν τον τρόπο βελτιώνουμε λίγο ακόμα την απόδοση του αλγορίθμου μας καθώς έχουμε αφαιρέσει τον δεύτερο βρόγχο for. Κερδίζουμε επειδή έχει αφαιρεθεί η εντολή διακλάδωσης και η ενημέρωση του μετρητή που υπάρχει μέσα στο for loop.

Ακολουθούν ορισμένα screenshot από διάφορες εκτελέσεις του αλγορίθμου για να διαπιστωθεί η ορθή επαλήθευση των αποτελεσμάτων. Παρακάτω βλέπουμε ότι ο αλγόριθμος μας δουλεύει και όταν τα μητρώα δεν είναι πολλαπλάσια του block size.



```

team13@tesla:~/workspace$ ./clevermulSym
A size : (4 X 5)
GPU will create : 3.000000 blocks
GPU will create : 4 threads per block
Time elapsed = 0.0376 ms
1.141269      0.927989      1.336754      1.131578      0.779726
0.927989      1.131723      0.955727      1.365832      0.670730
1.336754      0.955727      2.077672      1.499689      1.417770
1.131578      1.365832      1.499689      1.985407      1.275459
0.779726      0.670730      1.417770      1.275459      1.194298
team13@tesla:~/workspace$

```

Εικόνα 12: 4x5 Matrix size με 2x2 block size

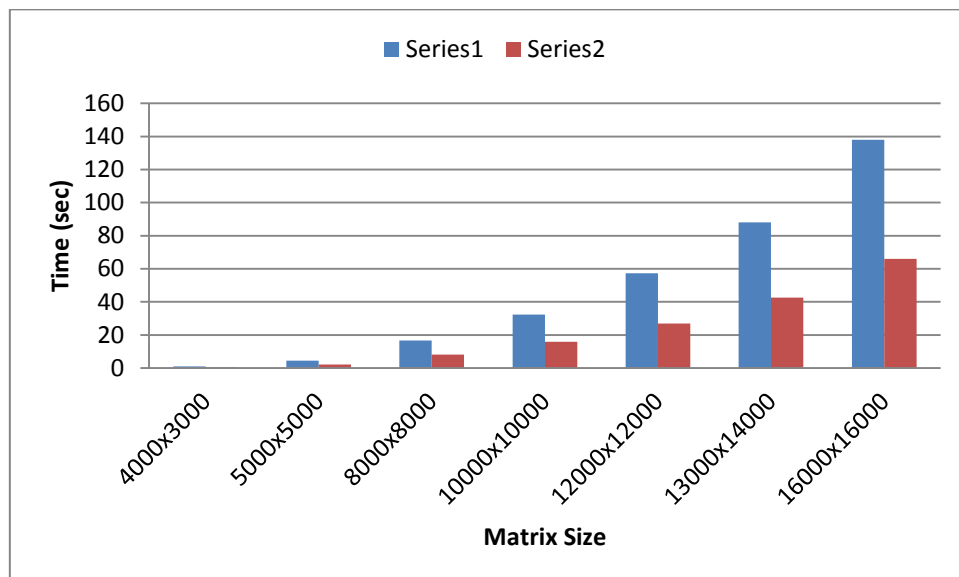
```

team13@tesla:~/workspace$ ./clevermulSym
A size : (16000 X 16000)
GPU will create : 1000.000000 blocks
GPU will create : 256 threads per block
Time elapsed = 67719.2 ms
team13@tesla:~/workspace$ █

```

Εικόνα 13: 16K x 16K

Size	4000x3000	5000x5000	8000x8000	10000x10000	12000x12000	13000x14000	16000x16000
Shared Memory	0.9	4.4	16.7	32.4	57.2	88.1	138.2
Symmetric/LU	0.4	2.1	8.2	15.8	27.3	42.6	66.8

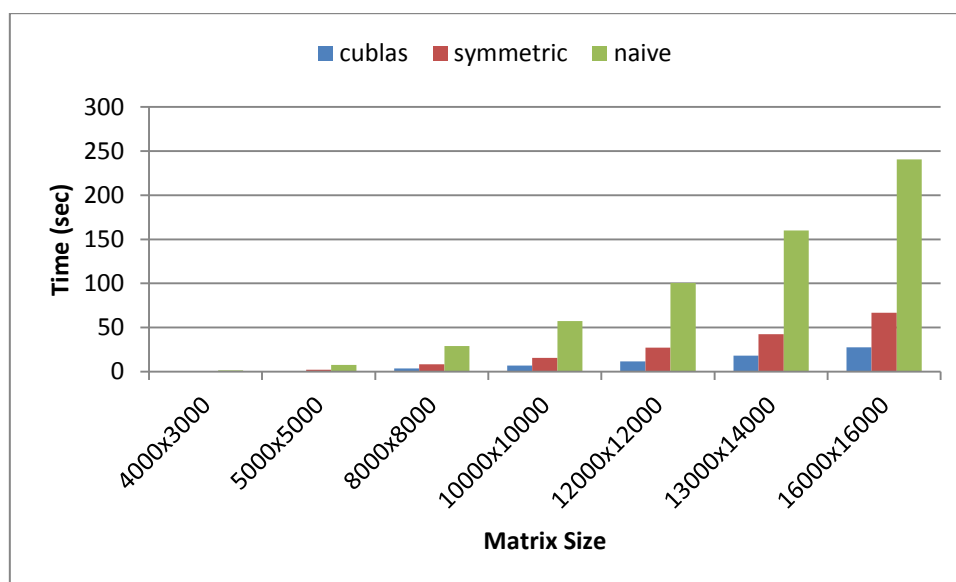


Εικόνα 14: Series 1 → Shared Memory, Series 2 → Symmetric & Loop Unrolling

Παρατήρηση:

- Σε αυτό το σημείο πρέπει να επισημάνουμε την αισθητή διαφορά που βλέπουμε στον χρόνο εκτέλεσης από την στιγμή που εκμεταλλευόμαστε την συμμετρικότητα του αποτελέσματος. Για την ακρίβεια κάνουμε λίγες περισσότερες από τις μισές πράξεις σε σχέση με την πρώτη έκδοση του προγράμματος που χρησιμοποιεί μόνο shared memory & coalesced memory access και γι αυτό τον λόγο ο χρόνος εκτέλεσης μειώνεται στην μέση.

Τέλος ακολουθεί μια τελευταία γραφική παράσταση που παρουσιάζει την απόδοση του `naivemul` (ερωτήματος 2) σε σχέση με τον τελικό αλγόριθμο που υλοποιήσαμε στο ερώτημα 3 (`clevermulSym`) ο οποίος χρησιμοποιεί shared memory, coalesced memory access, loop unrolling και εκμεταλλεύεται και την συμμετρικότητα του τελικού πίνακα και σε σχέση με την συνάρτηση της `cuda cublasDgemm`.



Εικόνα 15: Απόδοση cublas VS symmetric VS naive