



ΤΜΗΜΑ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ & ΜΗΧΑΝΙΚΩΝ Η/Υ

ECE415 - ΣΥΣΤΗΜΑΤΑ ΥΠΟΛΟΓΙΣΜΟΥ ΥΨΗΛΩΝ
ΕΠΙΔΟΣΕΩΝ

Lab 3 - Εισαγωγή στην CUDA

Φοιτητές:

Καραγεώργος Νικόλαος 02528, nkarageorgos@uth.gr

Λαμπρινός Ισίδωρος 02551, ilamprinos@uth.gr

Δεκέμβριος 2021

Περιεχόμενα

Εισαγωγή	1
Βήμα 0	1
Βήμα 1	2
Βήμα 2	3
Βήμα 3	3
Βήμα 4	5
Βήμα 5	6
Βήμα 6	9
Βήμα 7	11
Βήμα 8	13

Εισαγωγή

Η συνέλιξη (convolution) βρίσκει εφαρμογή σε πλήθος εφαρμογών της μηχανικής και των μαθηματικών. Πολλά φίλτρα μετασχηματισμού εικόνας στην ουσία υλοποιούνται ως συνελίξεις, όπως για παράδειγμα το φίλτρο Gaussian blur. Στο πεδίο της επεξεργασίας εικόνας, ένα φίλτρο συνέλιξης είναι απλά το σημείο-προς-σημείο γινόμενο των βαρών του φίλτρου με τα pixels της εικόνας εισόδου, εντός ενός παραθύρου που περιλαμβάνει κάθε pixel εξόδου.

Στην παρούσα εργασία θα ασχοληθούμε με ένα διαχωρίσιμο δισδιάστατο φίλτρο συνέλιξης που εφαρμόζεται σε μία $N \times N$ εικόνα. Πιο συγκεκριμένα, θα υλοποιήσουμε τον κώδικα του παραπάνω φίλτρου σε CUDA και θα πραγματοποιήσουμε διάφορα πειράματα πάνω στην υλοποίηση του σε GPU ώστε να απαντήσουμε τα ερωτήματα της εργασιάς όσον αφορά την απόδοση του και άλλες παραμέτρους.

Όλη η ανάπτυξη έγινε στο σύστημα inf-mars1 (10.64.82.31) το οποίο διαθέτει μία κάρτα GTX690, με 2 chips, ενώ τελικές μετρήσεις έγιναν στο πιο ισχυρό csl-artemis (10.64.82.65) το οποίο διαθέτει μία κάρτα Tesla K80 η οποία έχει 2 GK210 GPU chips

Βήμα 0

```
1 ./deviceQuery Starting...
2
3  CUDA Device Query (Runtime API) version (CUDA static linking)
4
5  Detected 2 CUDA Capable device(s)
6
7  Device 0: "Tesla K80"
8  CUDA Driver Version / Runtime Version      11.4 / 11.5
9  CUDA Capability Major/Minor version number: 3.7
10 Total amount of global memory:             11441 MBytes (11997020160 bytes)
11 (013) Multiprocessors, (192) CUDA Cores/MP: 2496 CUDA Cores
12 GPU Max Clock rate:                        824 MHz (0.82 GHz)
13 Memory Clock rate:                         2505 Mhz
14 Memory Bus Width:                          384-bit
15 L2 Cache Size:                             1572864 bytes
16 Maximum Texture Dimension Size (x,y,z)     1D=(65536), 2D=(65536, 65536), 3D
    =(4096, 4096, 4096)
17 Maximum Layered 1D Texture Size, (num) layers 1D=(16384), 2048 layers
18 Maximum Layered 2D Texture Size, (num) layers 2D=(16384, 16384), 2048 layers
19 Total amount of constant memory:            65536 bytes
20 Total amount of shared memory per block:    49152 bytes
21 Total shared memory per multiprocessor:     114688 bytes
22 Total number of registers available per block: 65536
23 Warp size:                                  32
24 Maximum number of threads per multiprocessor: 2048
25 Maximum number of threads per block:       1024
26 Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
27 Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
28 Maximum memory pitch:                       2147483647 bytes
29 Texture alignment:                          512 bytes
```

```

30 Concurrent copy and kernel execution: Yes with 2 copy engine(s)
31 Run time limit on kernels: No
32 Integrated GPU sharing Host Memory: No
33 Support host page-locked memory mapping: Yes
34 Alignment requirement for Surfaces: Yes
35 Device has ECC support: Enabled
36 Device supports Unified Addressing (UVA): Yes
37 Device supports Managed Memory: Yes
38 Device supports Compute Preemption: No
39 Supports Cooperative Kernel Launch: No
40 Supports MultiDevice Co-op Kernel Launch: No
41 Device PCI Domain ID / Bus ID / location ID: 0 / 6 / 0
42 Compute Mode:
43 < Default (multiple host threads can use ::cudaSetDevice() with device
    simultaneously) >
44
45 Device 1: "Tesla K80"
46 CUDA Driver Version / Runtime Version 11.4 / 11.5
47 CUDA Capability Major/Minor version number: 3.7
48 Total amount of global memory: 11441 MBytes (11997020160 bytes)
49 (013) Multiprocessors, (192) CUDA Cores/MP: 2496 CUDA Cores
50 GPU Max Clock rate: 824 MHz (0.82 GHz)
51 Memory Clock rate: 2505 Mhz
52 Memory Bus Width: 384-bit
53 L2 Cache Size: 1572864 bytes
54 Maximum Texture Dimension Size (x,y,z) 1D=(65536), 2D=(65536, 65536), 3D
    =(4096, 4096, 4096)
55 Maximum Layered 1D Texture Size, (num) layers 1D=(16384), 2048 layers
56 Maximum Layered 2D Texture Size, (num) layers 2D=(16384, 16384), 2048 layers
57 Total amount of constant memory: 65536 bytes
58 Total amount of shared memory per block: 49152 bytes
59 Total shared memory per multiprocessor: 114688 bytes
60 Total number of registers available per block: 65536
61 Warp size: 32
62 Maximum number of threads per multiprocessor: 2048
63 Maximum number of threads per block: 1024
64 Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
65 Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
66 Maximum memory pitch: 2147483647 bytes
67 Texture alignment: 512 bytes
68 Concurrent copy and kernel execution: Yes with 2 copy engine(s)
69 Run time limit on kernels: No
70 Integrated GPU sharing Host Memory: No
71 Support host page-locked memory mapping: Yes
72 Alignment requirement for Surfaces: Yes
73 Device has ECC support: Enabled
74 Device supports Unified Addressing (UVA): Yes
75 Device supports Managed Memory: Yes
76 Device supports Compute Preemption: No
77 Supports Cooperative Kernel Launch: No
78 Supports MultiDevice Co-op Kernel Launch: No
79 Device PCI Domain ID / Bus ID / location ID: 0 / 7 / 0
80 Compute Mode:
81 < Default (multiple host threads can use ::cudaSetDevice() with device
    simultaneously) >
82 > Peer access from Tesla K80 (GPU0) -> Tesla K80 (GPU1) : Yes
83 > Peer access from Tesla K80 (GPU1) -> Tesla K80 (GPU0) : Yes
84
85 deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 11.4, CUDA Runtime Version
    = 11.5, NumDevs = 2
86 Result = PASS

```

Βήμα 1

Μελετάμε τις παραμέτρους εκτέλεσης του compiler nvcc, γράφοντας nvcc -help σε ένα τερματικό. Για τη μετταγλώττιση του κώδικα .cu χρησι-

μπορούμε τον nvcc compiler με το μέγιστο βαθμό βελτιστοποιήσεων (-O4).

Βήμα 2

Σε αυτό το βήμα ζητείται να γράψουμε κώδικα που εκτελεί την συνέλιξη δύο διαστάσεων σε μια εικόνα διαστάσεων NxN στη GPU με βάση τον κώδικα που μας δίνεται για τη CPU.

Σε αυτόν τον κώδικα παρατηρούμε πως μπορεί να πραγματοποιηθεί μια σημαντική βελτιστοποίηση loop invariant code motion.

Η έκφραση:

$$h_Dst[y * imageW + x] = sum$$

στην συνάρτηση για την συνέλιξη κατά γραμμές αλλά και στην συνάρτηση για την συνέλιξη κατά γραμμές βρίσκεται μέσα σε μία for με μετρητή k. Ωστόσο όπως παρατηρούμε δεν εξαρτάται κανένας όρος της άπο το k οπότε την τοποθετούμε εκτός της for αυτής.

Όσον αφορά την υλοποίηση για την GPU σε αυτό το βήμα ζητείται να χρησιμοποιηθούν τα threads ενός μόνο block ώστε κάθε thread να υπολογίζει ένα στοιχείο της εικόνας-αποτελέσματος. Κάθε thread της GPU δουλεύει για μία διαφορετική θέση στον πίνακα των αποτελεσμάτων οπότε τα threads δεν χρειάζεται να συγχρονιστούν μεταξύ τους. Για να εντοπίσουν την θέση που δουλεύουν στο block υπάρχουν οι δείκτες threadIdx.x και threadIdx.y. Ο πρώτος δείχνει την οριζόντια συντεταγμένη του thread στο block ενώ ο δεύτερος την κατακόρυφη. Έτσι δεν χρειάζεται να χρειάζεται να χρησιμοποιηθεί οι δύο εξωτερικές for με τους μετρητές x και y και αυτοί αντικαθιστούνται από τους threadIdx.x και threadIdx.y αντίστοιχα.

Βήμα 3

3α. Μέχρι ποιο μέγεθος εικόνας (πειραματιστείτε πάντα με μεγέθη εικόνας που είναι δυνάμεις του 2 και μεγαλύτερες από το μήκος του φίλτρου) μπορείτε να υποστηρίξετε χωρίς να συμβεί κάποιο σφάλμα χρόνου εκτέλεσης; Γιατί συμβαίνει αυτό το σφάλμα;

Το μέγιστο μέγεθος εικόνας που μπορεί να υποστηριχθεί χωρίς να συμβεί κάποιο σφάλμα χρόνου εκτέλεσης είναι 32x32. Όταν δοκιμάζουμε σαν είσοδο μέγεθος εικόνας 64x64 εμφανίζεται το παρακάτω:

```
Enter filter radius : 16
Enter image size. Should be a power of two and greater than 33 : 64
Image Width x Height = 64 x 64

Allocating and initializing host arrays...
Allocating and initializing device arrays...
CPU computation...
GPU computation...
Cuda error Convolution2D.cu: 247: 'invalid configuration argument'.
```

Figure 1: Σφάλμα χρόνου εκτέλεσης για είσοδο εικόνα μεγέθους 64x64

Αυτό μας φαίνεται απόλυτα λογικό καθώς όπως παρατηρήσαμε απο το deviceQuery νωρίτερα κάθε block διαθέτει 1024 threads. Στο συγκεκριμένο ερώτημα χρησιμοποιούμε μόνο ένα block άρα έχουμε διαθέσιμα 1024 threads για μία εικόνα με imageW*imageH pixels όπου imageW=imageH οπότε και μαθηματικά μπορούμε να φτάσουμε μέχρι $32 \times 32 = 1024$.

3β) Για το μέγιστο μέγεθος εικόνας που μπορείτε να υποστηρίξετε, ποια είναι η μέγιστη ακρίβεια, σε σχέση με το μέγεθος του φίλτρου, που μπορεί να υποστηριχθεί χωρίς να παρουσιάζονται σφάλματα σύγκρισης;

Ως ακρίβεια ορίζουμε την μέγιστη απόλυτη διαφορά μεταξύ των αποτελεσμάτων της συνέλιξης στη CPU και τη GPU (σφάλμα).

Για το μέγιστο μέγεθος εικόνας που μπορεί να υποστηριχθεί δηλαδή 32x32 το μέγιστο μέγεθος της ακτινής του φίλτρου που μπορεί να χρησιμοποιηθεί είναι 15.

Ακτίνα φίλτρου	Μέγιστο απόλυτο σφάλμα
1	0.007812
2	0.015625
3	0.0625
4	0.1875
5	0.25
6	0.5
7	0.375
8	0.75
9	1
10	1
11	1.5
12	1
13	1.5
14	1
15	1

Παρατηρούμε ότι η μέγιστη ακρίβεια δηλαδή το μικρότερο απόλυτο σφάλμα(0.007812) παρατηρείται με την μικρότερη ακτίνα φίλτρου που είναι 1 και αρα για συνολικό μήκος φίλτρου 3.

Βήμα 4

Σε αυτό το ερώτημα χρησιμοποιούμε πολλαπλά τετράγωνα blocks που είναι οργανωμένο σε ένα τετράγωνο grid. Τα block αυτά όπως είδαμε και

προηγούμενως διαθέτουν 1024 threads άρα έχουν διαστάσεις 32x32. Τα threads πλέον για να βρουν το στοιχείο της εικόνας το οποίο τους έχει ανατεθεί για υπολογισμό χρησιμοποιούν και τα blockDim.x, blockIdx.x, blockDim.y και blockIdx.y.

Το μέγιστο μέγεθος εικόνας που μπορεί υποστηριχθεί πλέον στο csl-artemis (10.64.82.65) είναι 16384x16384. Μετά από αυτό το μέγεθος εικόνας εμφανίζεται:

```
Enter filter radius : 4
Enter image size. Should be a power of two and greater than 9 : 32768
Image Width x Height = 32768 x 32768

Allocating and initializing host arrays...
Allocating and initializing device arrays...
Error while allocating device memory.
```

Figure 2: Σφάλμα χρόνου εκτέλεσης για είσοδο εικόνα μεγέθους 32768x32678

Αυτό συμβαίνει γιατί πλέον δεν υπάρχει αρκετή μνήμη στη GPU για τόσο μεγάλα μεγέθη.

Βήμα 5

5a)

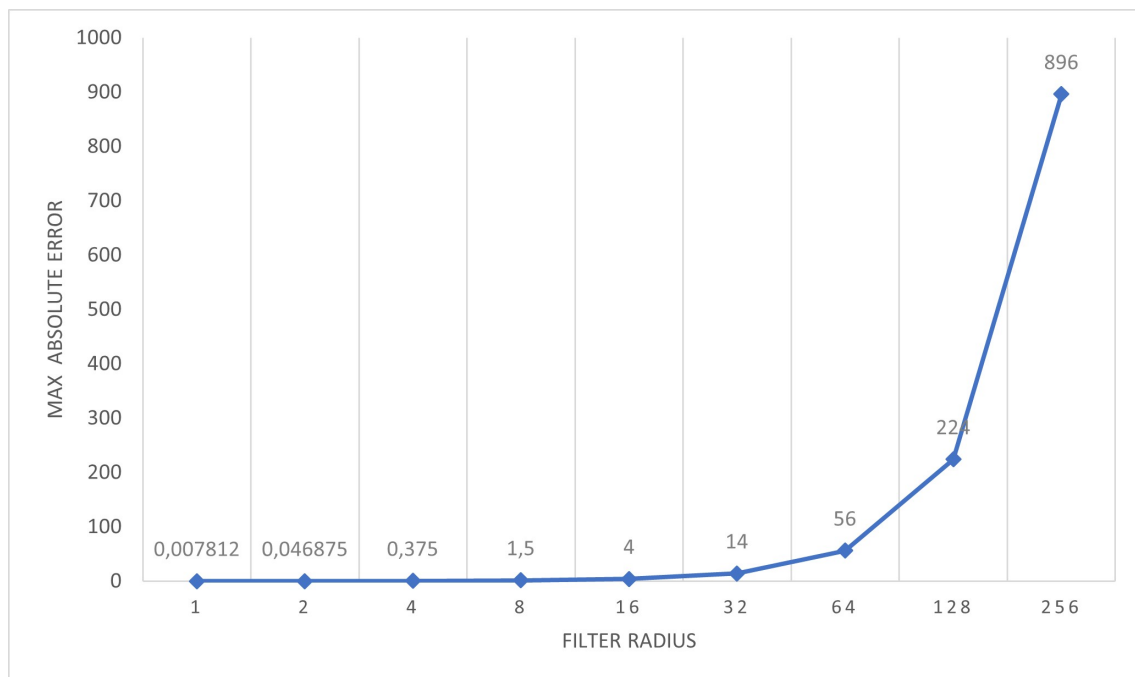


Figure 3: Ακρίβεια σε σχέση με την ακτίνα του φίλτρου για εικόνα μεγέθους 16384x16384

Όπως παρατηρούμε από τα παραπάνω υπάρχει μία εκθετική αύξηση της μέγιστης ακρίβειας των αποτελεσμάτων της GPU σε σχέση με τη CPU όσο ανεβαίνει το μέγεθος της ακτίνας του φίλτρου. Αυτό συμβαίνει καθώς οι GPU πετυχαίνουν πολύ καλύτερες επιδόσεις στον χρόνο εκτέλεσης καθώς κάνουν fused multiply-add (FMAD) υπολογισμούς, με συνέπεια όμως να θυσιάζουν κάποια ακρίβεια. Ειδικά στην συγκεκριμένη περίπτωση χρησιμοποιούνται αριθμοί κινητής υποδιαστολής απλής ακρίβειας (float). Αυτοί οι αριθμοί κάνουν σημαντικά προσεγγιστικά λάθη με αποτέλεσμα να μην πετυχαίνεται η ακρίβεια που θα είχαμε αν χρησιμοποιούνταν αριθμοί κινητής υποδιαστολής διπλής ακρίβειας (doubles).

5b)

Διαστάσεις Εικόνas	Μέσος Χρόνος Εκτέλεσης CPU για floats	Μέσος Χρόνος Εκτέλεσης GPU για floats
64x64	0,000364805	0,000121654
128x128	0,001507745	0,000237222
256x256	0,006107532	0,000673197
512x512	0,024851352	0,002491741
1024x1024	0,101809057	0,008731715
2048x2048	0,485297299	0,033883184
4096x4096	2,326409672	0,149977916
8192x8192	9,183927318	0,542511845
16384x16384	37,31040972	2,116709781

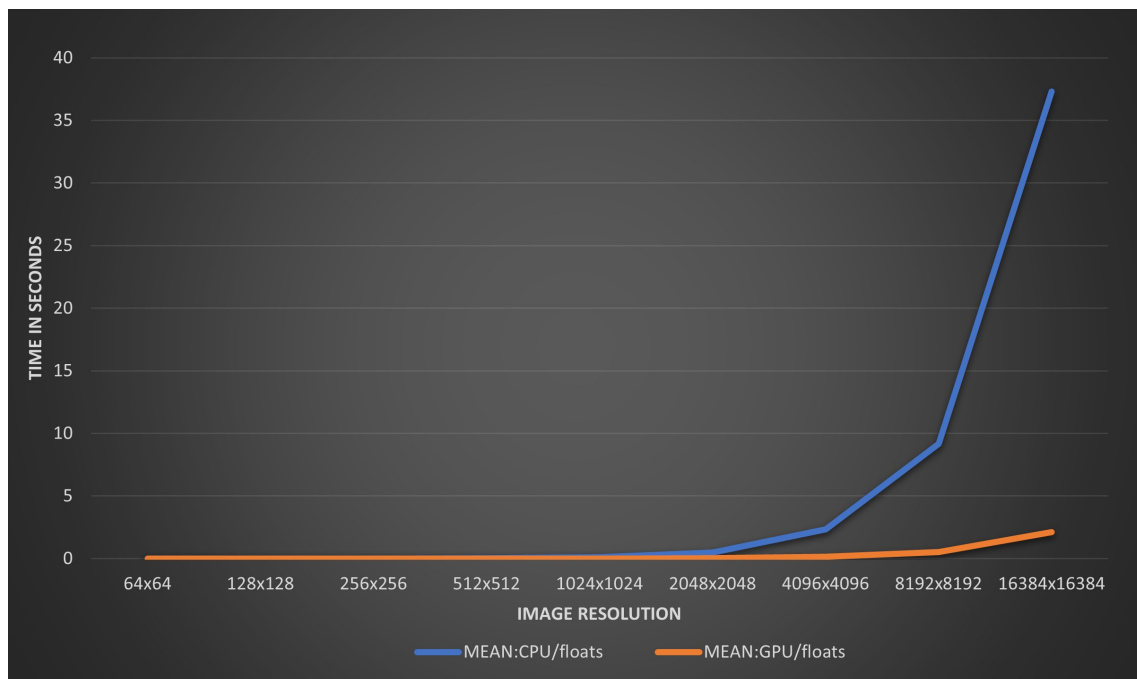


Figure 4: Χρόνος εκτέλεσης για τύπο floats σε CPU και GPU

Παρατηρούμε πως οι επιδόσεις του χρόνου εκτέλεσης στην GPU είναι πολύ καλύτερες σε σχέση με αυτές στην CPU. Αυτό γίνεται ακόμη πιο ξεκάθαρο όσο μεγαλώνει το μέγεθος της εικόνας-εισόδου.

Βήμα 6

Διαστάσεις Εικόνας	Μέσος Χρόνος Εκτέλεσης CPU doubles	Μέσος Χρόνος Εκτέλεσης GPU doubles
64x64	0,000408155	0,000149894
128x128	0,001703403	0,00032137
256x256	0,00669445	0,000950256
512x512	0,027434598	0,00323968
1024x1024	0,111931962	0,011528599
2048x2048	0,567745006	0,044662288
4096x4096	2,467726286	0,200996171
8192x8192	9,756948556	0,790095544
16384x16384	39,51873891	3,072551799

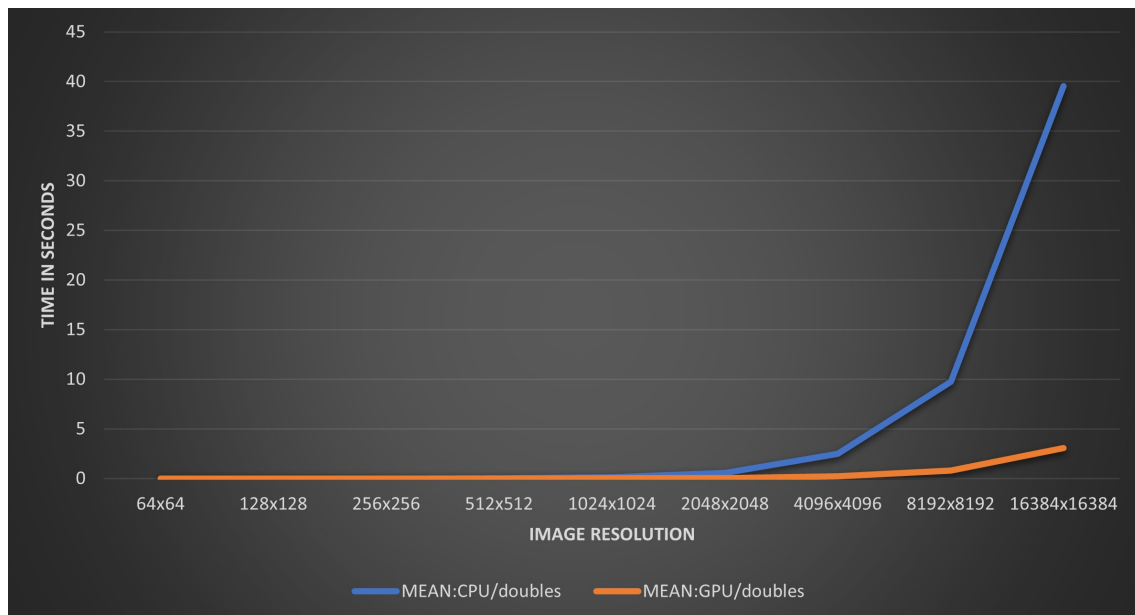


Figure 5: Χρόνος εκτέλεσης για τύπο doubles σε CPU και GPU

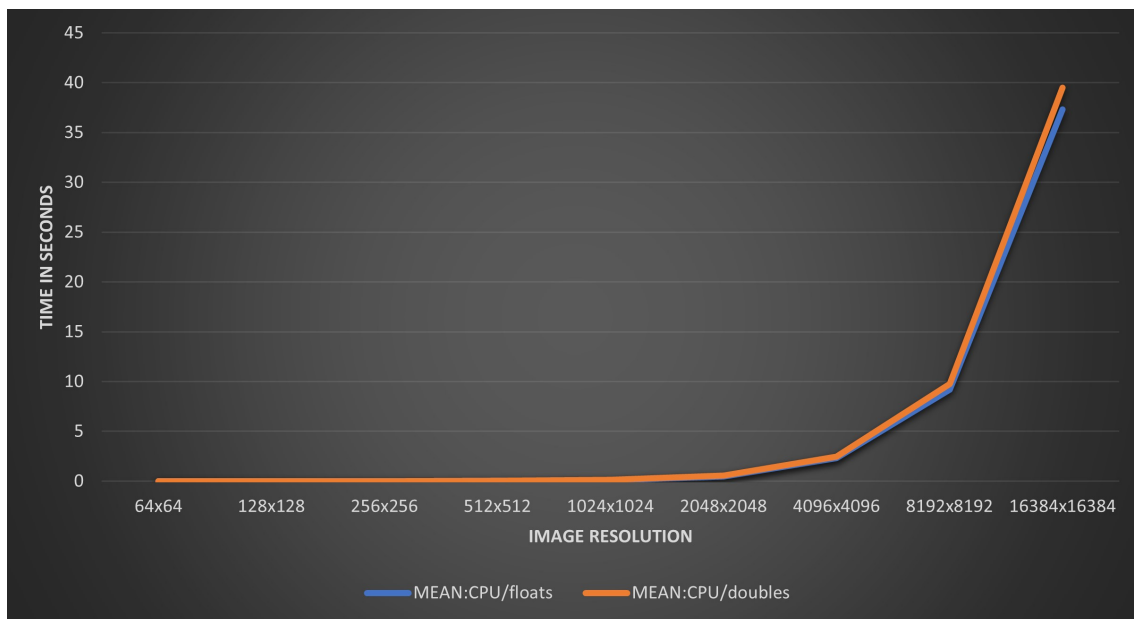


Figure 6: Χρόνος εκτέλεσης για τύπους floats και doubles σε CPU

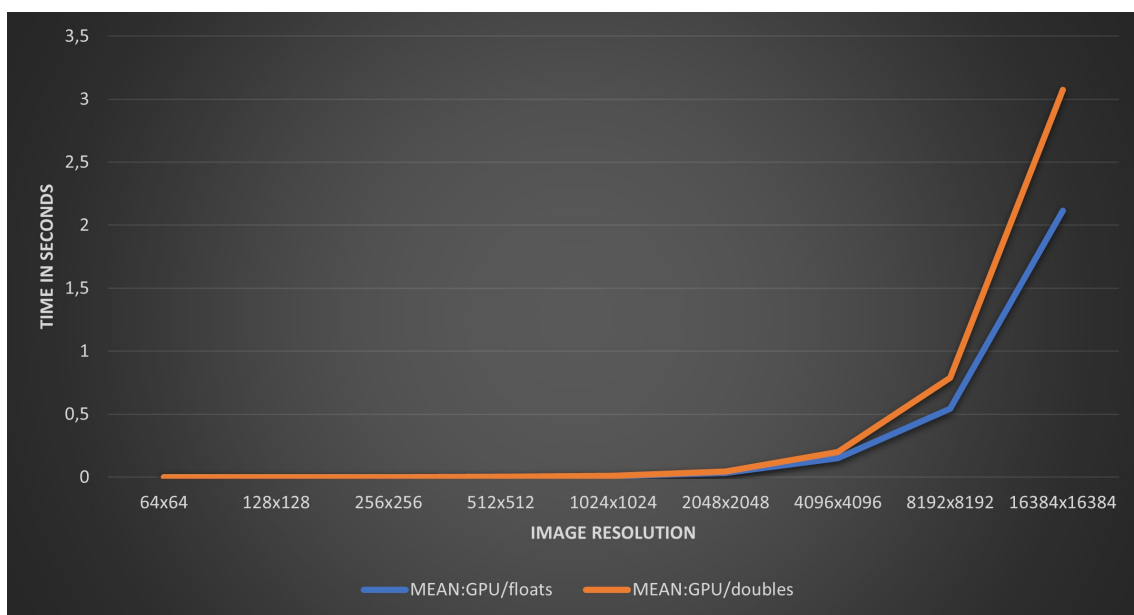


Figure 7: Χρόνος εκτέλεσης για τύπους floats και doubles σε GPU

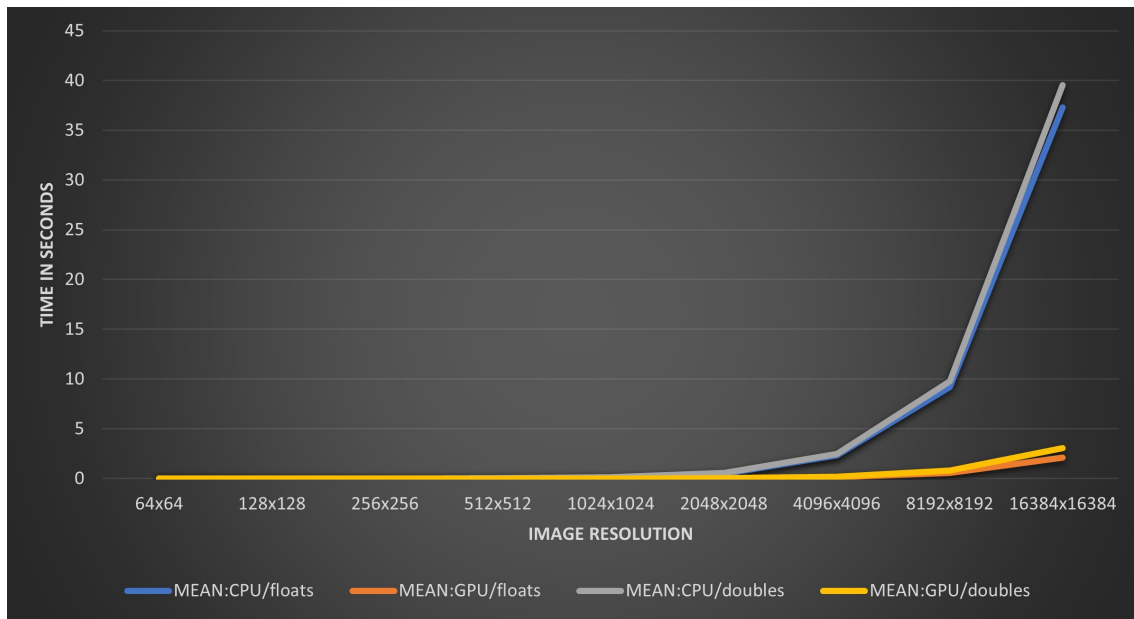


Figure 8: Χρόνος εκτέλεσης για τύπους floats και doubles σε CPU και GPU

Χρησιμοποιώντας στοιχεία κινητής υποδιαστολής διπλής ακρίβειας (double) παρατηρούμε ότι πετυχαίνεται πάρα πολύ καλή ακρίβεια αφού δοκιμάζοντας μέχρι και μέγεθος ακτίνας φίτρου ίσο με 128 το μέγιστο απόλυτο σφάλμα παραμένει 0. Αυτό συμβαίνει καθώς πλέον χρησιμοποιούνται 64 bits για την αναπαράσταση των αριθμών κινητής υποδιαστολής ενώ προηγουμένως με floats χρησιμοποιούνταν 32 bits οπότε γίνονται πολύ λιγότερα προσεγγιστικά λάθη. Ωστόσο, παρατήρουμε μια μικρή αύξηση στο χρόνο εκτέλεσης, το οποίο είναι λογικό αφού πλέον έχουμε να κάνουμε υπολογισμούς με αριθμούς μεγαλύτερου μεγέθους.

Βήμα 7

α) Πόσες φορές διαβάζεται κάθε στοιχείο της εικόνας εισόδου και του φίλτρου κατά την εκτέλεση του kernel;

Όσον αφορά την εικόνα εισόδου είναι ξεκάθαρο ότι τα στοιχεία-pixels που βρίσκονται εσωτερικά δηλαδή πιο κοντά στο κέντρο θα διαβάσουν περισσότερες φορές στην πράξη της συνέλιξης από αυτά που βρίσκονται στις άκρες της εικόνας. Πιο συγκεκριμένα, το κάθε thread αναλαμβάνει να υπολογίσει το αποτέλεσμα-pixel διαβάζοντας $2 \times \text{filterR} + 1$ στοιχεία γύρω από το αρχικό pixel. Αυτό συνεπάγεται ότι ένα στοιχείο μπορεί να δια-

βαστεί το πολύ $2\text{filterR}+1$ κατά την πράξη μίας συνέλιξης.

Μετά απο πειραματισμούς με μεγέθη φίτρου και εικόνας καταλήγουμε στους γενικούς τύπους όπου κάθε $\text{pixel}(x,y)$ με $x \in [0,\text{imageW}-1]$ και $y \in [0,\text{imageH}-1]$ διαβάζεται:

Στη συνέλιξη κατά γραμμές:

$$\text{Row_Reads}_{(x,y)} = \begin{cases} \text{filterR} + x + 1, & x \leq \text{filterR} \\ 2\text{filterR} + 1 & \text{filterR}, < x < \text{imageW} - \text{filterR} \\ \text{filterR} + \text{imageW} - x & x, \geq \text{imageW} - \text{filterR} \end{cases}$$

Στη συνέλιξη κατά στήλες:

$$\text{Col_Reads}_{(x,y)} = \begin{cases} \text{filterR} + y + 1, & y \leq \text{filterR} \\ 2\text{filterR} + 1 & \text{filterR}, < y < \text{imageW} - \text{filterR} \\ \text{filterR} + \text{imageW} - y, & y \geq \text{imageW} - \text{filterR} \end{cases}$$

Συνολικά το κάθε στοιχείο διαβάζεται:

$$\text{Reads}_{(x,y)} = \text{Col_Reads}_{(x,y)} + \text{Row_Reads}_{(x,y)}$$

Όσον αφορά το πόσες φορές διαβάζονται τα στοιχεία του φίλτρου είναι πάλι ευδιάκριτο πως τα κεντρικό στοιχείο του φίλτρου, που βρίσκεται στη θέση filterR από τα $2*\text{filterR}+1$ στοιχεία, θα διαβαστεί όσες φορές όσες είναι και τα στοιχεία της εικόνας εισόδου δηλαδή $\text{imageH}*\text{imageW}$. Ενώ όσο πιο εξωτερικά είναι τα στοιχεία στο φίλτρο θα διαβάζονται λιγότερες φορές. Πιο συγκεκριμένα, ένα στοιχείο που βρίσκεται στη θέση $\text{filterR}-n$ όπου $n \in [0,R]$ του φίλτρου θα διαβαστεί $(\text{imageW}-n)*\text{imageH}$ φορές στη συνέλιξη κατα γραμμές και $(\text{imageH}-n)*\text{image}$ φορές στη συνέλιξη κατά στήλες.

Συνολικά:

$$\text{Reads}_{(n)} = (\text{imageW} - |\text{filterR} - n|) * \text{imageH} + (\text{imageW} - |\text{filterR} - n|) * \text{imageH}$$

Γνωρίζοντας ότι $\text{imageW}=\text{imageH}$:

$$\text{Reads}_{(n)} = 2(\text{imageW} - |\text{filterR} - n|) * \text{imageW}$$

b) Ποιός είναι ο λόγος προσπελάσεων μνήμης προς πράξεις κινητής υποδιαστολής; Θεωρήστε τους πολλαπλασιασμούς και τις προσθέσεις ως ξεχωριστές πράξεις και αγνοείστε την αποθήκευση του αποτελέσματος. Μετρήστε μόνο τις αναγνώσεις από την global memory της GPU ως προσπελάσεις μνήμης.

Παρατηρούμε ότι για κάθε πράξη πολλαπλασιασμού και πρόσθεσης στη συνέλιξη γίνονται δύο προσπελάσεις στην μνήμη, μία για την ανάγνωση του στοιχείου του φίλτρου και μία για την ανάγνωση του στοιχείου της εικόνας εισόδου που είναι για υπολογισμό. Επομένως, συμπεραίνουμε πως ο λόγος προσπελάσεων μνήμης προς πράξεις κινητής υποδιαστολής είναι 1.

Βήμα 8

Διαστάσεις Εικόνας	Μέσος Χρόνος Εκτέλεσης CPU χωρίς padding	Μέσος Χρόνος Εκτέλεσης CPU με padding
64x64	0,000364805	0,000225059
128x128	0,001507745	0,000861735
256x256	0,006107532	0,003696144
512x512	0,024851352	0,014604245
1024x1024	0,101809057	0,058377305
2048x2048	0,485297299	0,240614454
4096x4096	2,326409672	1,054468082
8192x8192	9,183927318	4,283031374
16384x16384	37,31040972	17,3529493

Διαστάσεις Εικόνας	Μέσος Χρόνος Εκτέλεσης GPU χωρίς padding	Μέσος Χρόνος Εκτέλεσης GPU με padding
64x64	0,000121654	0,000161725
128x128	0,000237222	0,000290666
256x256 2	0,000673197	0,000769075
512x512	0,002491741	0,002616608
1024x1024	0,008731715	0,009078586
2048x2048	0,033883184	0,034716791
4096x4096	0,149977916	0,14242692
8192x8192	0,542511845	0,535792637
16384x16384	2,116709781	2,117604208

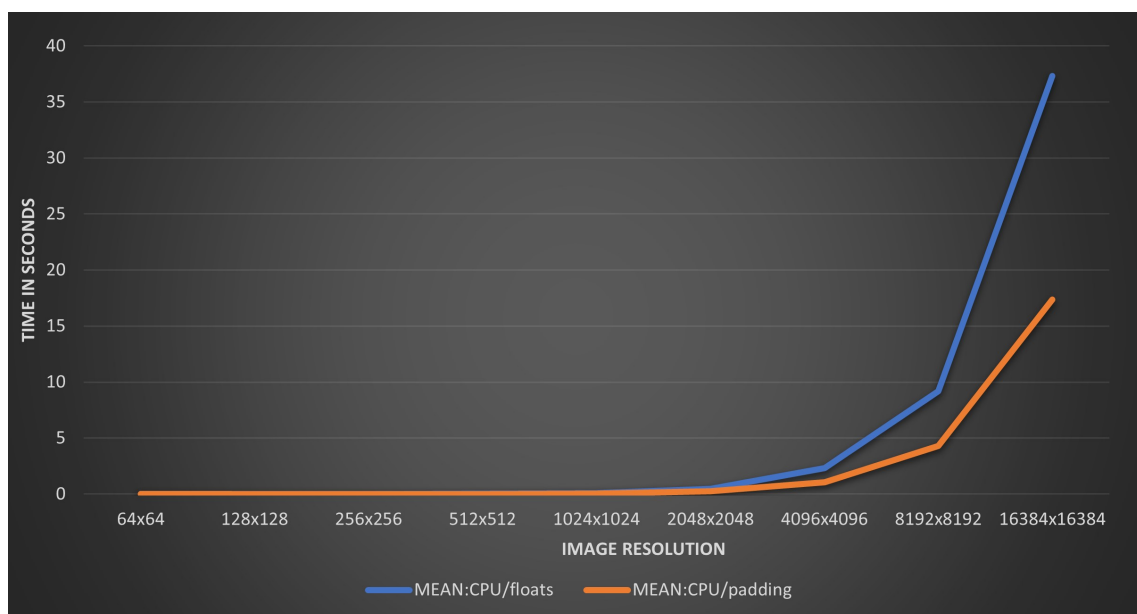


Figure 9: Χρόνος εκτέλεσης σε CPU με και χωρίς padding

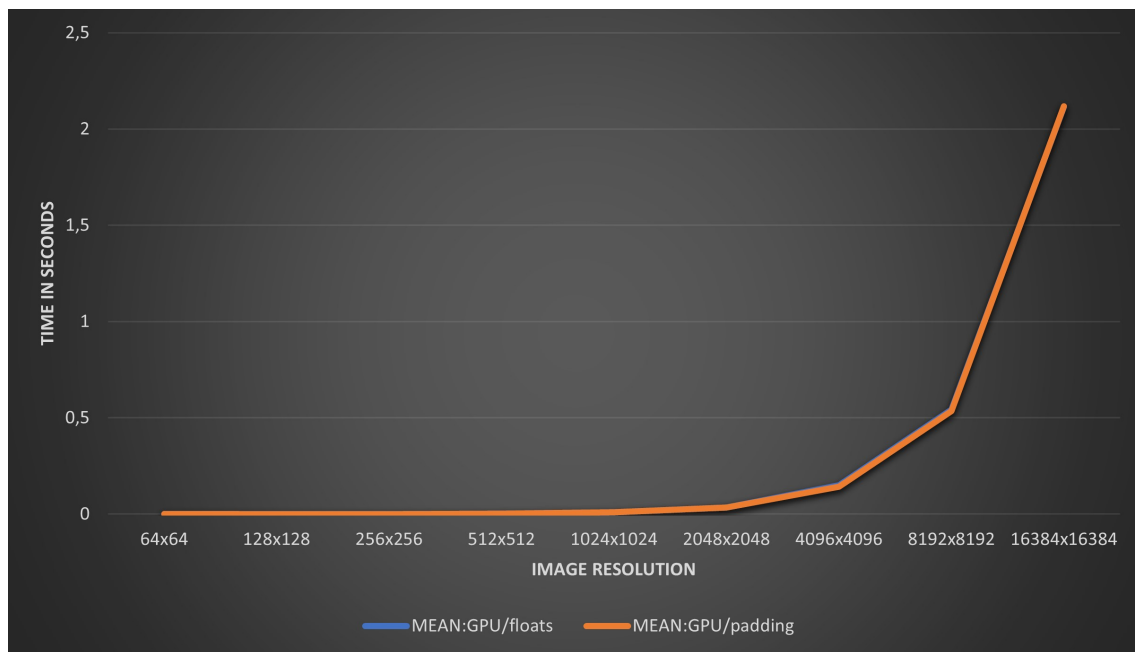


Figure 10: Χρόνος εκτέλεσης σε GPU με και χωρίς padding

Χρησιμοποιώντας padding καταφέρνουμε να απαλλαχθούμε εντελώς από τις if που ελέγχουν αν το φίλτρο του πίνακα βρίσκεται εκτός πίνακα. Αυτό όπως παρατηρείται παραπάνω είχε ως αποτέλεσμα την σημαντική πτώση του χρόνου εκτέλεσης για τη CPU ειδικά όσο ανεβαίνει το μέγεθος της εικόνας. Για τη GPU ωστόσο δεν παρατηρείται καποιά σημαντική διαφορά στους χρόνους εκτέλεσης.