



ΤΜΗΜΑ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ & ΜΗΧΑΝΙΚΩΝ Η/Υ

ECE415 - ΣΥΣΤΗΜΑΤΑ ΥΠΟΛΟΓΙΣΜΟΥ ΥΨΗΛΩΝ
ΕΠΙΔΟΣΕΩΝ

Lab 4 - Βελτιστοποιήσεις προγράμματος CUDA

Φοιτητές:

Καραγεώργος Νικόλαος 02528, nkarageorgos@uth.gr

Λαμπρινός Ισίδωρος 02551, ilamprinos@uth.gr

Ιανουάριος 2022

Περιεχόμενα

Εισαγωγή	1
Βήμα 0	1
Βήμα 1	3
Βήμα 2	6
Βήμα 3	7
Βήμα 4	8
Βήμα 5	8

Εισαγωγή

Στην παρούσα εργασία μας ζητήθηκε να πραγματοποιήσουμε βελτιστοποιήσεις στον κώδικα CUDA της προηγούμενης εργασίας ο οποίος εφαρμόζει, με χρήση συνέλιξης, ένα διαχωρίσιμο δισδιάστατο φίλτρο πάνω σε ένα δισδιάστατο πίνακα (ο οποίος υποθέσαμε ότι αντιστοιχεί σε μια εικόνα). Πιο συγκεκριμένα, σε αυτή την άσκηση κληθήκαμε να βελτιστοποιήσουμε τον προηγούμενο κώδικα, να τον εκτελέσουμε για πίνακες πολύ μεγαλύτερου μεγέθους και να προσπαθήσουμε να επικαλύψουμε μεταξύ τους διαφορετικούς υπολογισμούς καθώς και υπολογισμό με μεταφορές δεδομένων. Χρησιμοποιείται η έκδοση του κώδικα με padding που χρησιμοποιεί doubles.

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 (CUDART 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
17     = (4096, 4096, 4096)
18     Maximum Layered 1D Texture Size, (num) layers   1D=(16384), 2048 layers
19     Maximum Layered 2D Texture Size, (num) layers   2D=(16384, 16384), 2048 layers
20     Total amount of constant memory:                 65536 bytes
21     Total amount of shared memory per block:         49152 bytes
22     Total shared memory per multiprocessor:          114688 bytes
23     Total number of registers available per block:   65536
24     Warp size:                                       32
25     Maximum number of threads per multiprocessor:    2048
26     Maximum number of threads per block:            1024
27     Max dimension size of a thread block (x,y,z):   (1024, 1024, 64)
28     Max dimension size of a grid size (x,y,z):       (2147483647, 65535, 65535)
29     Maximum memory pitch:                           2147483647 bytes
30     Texture alignment:                              512 bytes
31     Concurrent copy and kernel execution:            Yes with 2 copy engine(s)
32     Run time limit on kernels:                       No
33     Integrated GPU sharing Host Memory:              No
34     Support host page-locked memory mapping:         Yes
35     Alignment requirement for Surfaces:              Yes
36     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

Βελτιστοποιήσεις του κώδικα σε GPU:

-Caching του πίνακα του φίλτρου d_Filter

Τοποθετήσαμε τον πίνακα του φίλτρου d_Filter σε κάποια cached memory ώστε να μειωθούν οι ακριβές προσπελάσεις στην global memory. Ειδικότερα, επιλέξαμε την constant memory η οποία υποστηρίζει μόνο αναγνώσεις καθώς ο συγκεκριμένος πίνακας παραμένει σταθερός και χρησιμοποιείται από όλα τα blocks για αναγνώσεις. Δεν επιλέξαμε την shared memory καθώς η συγκεκριμένη υποστηρίζει και εγγραφές οι οποίες δεν χρειάζονταν σε αυτό το σημείο και θα μπορούσαν να αξιοποιηθούν αργότερα.

-Tiling της εικόνας εισόδου στην shared memory

Η εικόνα εισόδου διασπάται σε μικρότερα τετράγωνα κομμάτια τα οποία αναθέτονται το καθένα σε ένα block για υπολογισμό. Τα threads του ίδιου block έχουν πρόσβαση στην ίδια shared memory και χρησιμοποιούν συνολικά για τους υπολογισμούς τους $(\text{blockDim.x} + 2\text{filterR}) * \text{blockDim.y}$ στοιχεία του κομματιού που τους ανατέθηκε. Λαμβάνοντας υπόψιν τα παραπάνω αποφασίσαμε να χρησιμοποιήσουμε την shared memory ως cache. Έτσι κάθε φορά όταν τα threads αρχίζουν την εκτέλεση τους φέρνουν κάποια διαφορετικά στοιχεία του κομματιού τους της αρχικής εικόνας από την global memory στην shared. Στο τέλος όλα τα threads έχουν συνεργαστεί και έχουν φέρει μαζί όλα τα στοιχεία του κομματιού που έχει ανατέθει στο block τους στην shared memory, τα οποία θα χρειαστούν το καθένα αργότερα για τους υπολογισμούς τους και θα τα κάνουν caching (βελτιστοποίηση coalescing). Αυτή η βελτιστοποίηση ωστόσο μας περιορίζει στο μέγεθος της ακτίνας του φίλτρου (filterR) που μπορούμε να χρησιμοποιήσουμε αφού το μέγεθος της shared memory κάθε block είναι αρκετά περιορισμένο, όπως βλέπουμε από το deviceQuery 49152 bytes. Πρέπει να ισχύει:

$(\text{blockDim.x} + 2\text{filterR}) * \text{blockDim.y} * \text{sizeof(double)} \leq \text{Size of shared memory}$

Δηλαδή πρέπει $\text{filterR} \leq 80$.

Επίσης με τον τρόπο με τον οποίο έχουμε υλοποίησει την μεταφορά των κομματιών στην shared memory το μικρότερο μέγεθος ακτίνας φίλτρου που μπορεί να χρησιμοποιηθεί είναι 16.

-Μείωση των ορισμάτων των συναρτήσεων του kernel

Αφαιρούμε από όρισμα μία απο τις δύο διαστάσεις της εικόνας αφού η εικόνα είναι τετράγωνη και είναι ίδιες. Αφαιρούμε επίσης την ακτίνα του φίλτρου από όρισμα και την ορίζουμε με `#define` κάπως την χρειαζομαστε και για τον ορισμό του πίνακα `d_Filter` στην constant memory ο οποίος επίσης αφαιρείται απο όρισμα της συνάρτησης.

GPU activities	Αρχική έκδοση	Βελτιστοποιημένη έκδοση
CUDA memcpy DtoH	407.20	408.77
CUDA memcpy HtoD	215.90	201.83
convolutionRowGPU	197.00	43.520
convolutionColumnGPU	167.29	32.210

Table 1: GPU Activities Execution Time

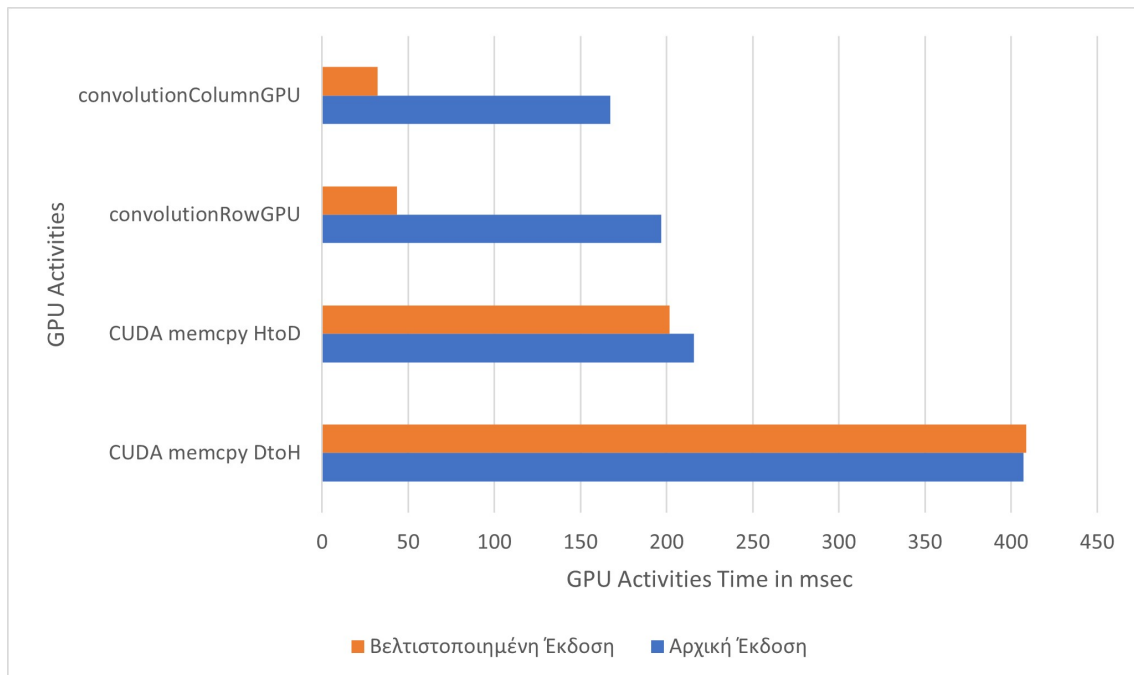


Figure 1: GPU Activities Time (msec)

GPU activities	Αρχική έκδοση	Βελτιστοποιημένη έκδοση
CUDA memcpy DtoH	41.00%	57.47%
CUDA memcpy HtoD	21.62%	30.17%
convolutionRowGPU	20.02%	6.84%
convolutionColumnGPU	17.01%	5.04%

Table 2: GPU Activities Execution Time (percentage)

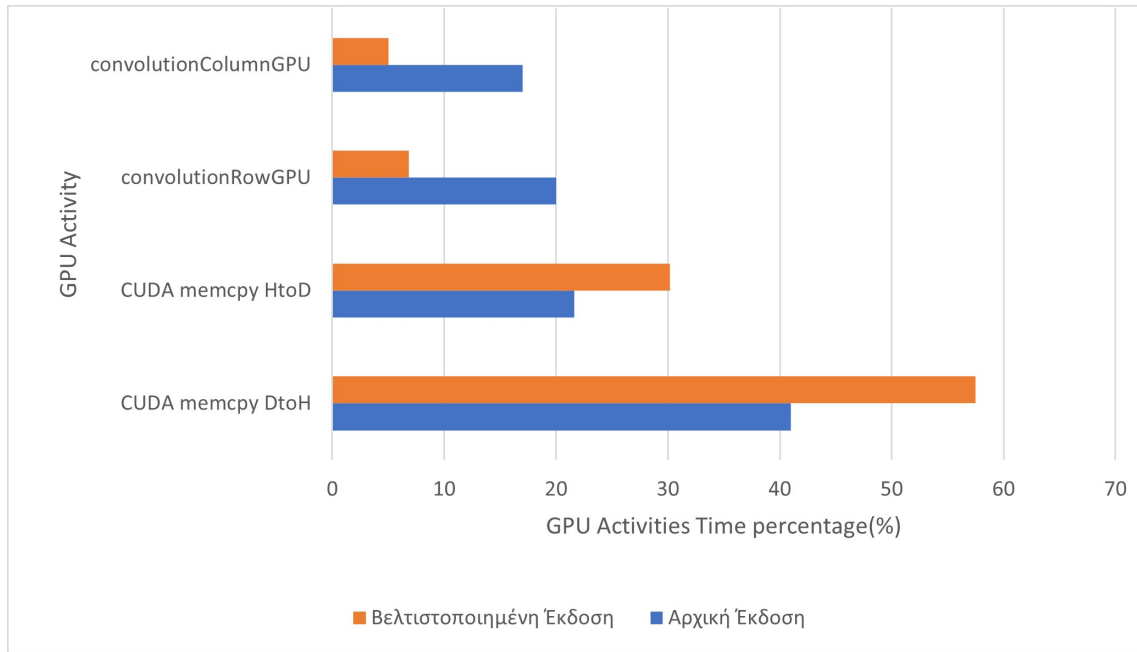


Figure 2: GPU Activities Time Percentage(%)

Βήμα 2

Κάνοντας profiling για διαδορετικές τιμές φίλτρου παρατηρούμε ότι όσο αυξάνεται το μέγεθος του φίλτρου τόσο η εκτέλεση των kernels ακριβ-
αίνει σε χρόνο. Το ίδιο ωστόσο δεν ισχύει για τις μεταφορές δεδομένων,
οι οποίες κυμαίνονται στις ίδιες τιμές ανεξαρτήτως ακτίνας φίλτρου και
δε φαίνεται να επηρεάζονται απ' αυτό. Για το πείραμα χρησιμοποιήσαμε
εικόνα μεγέθους 8192x8192.

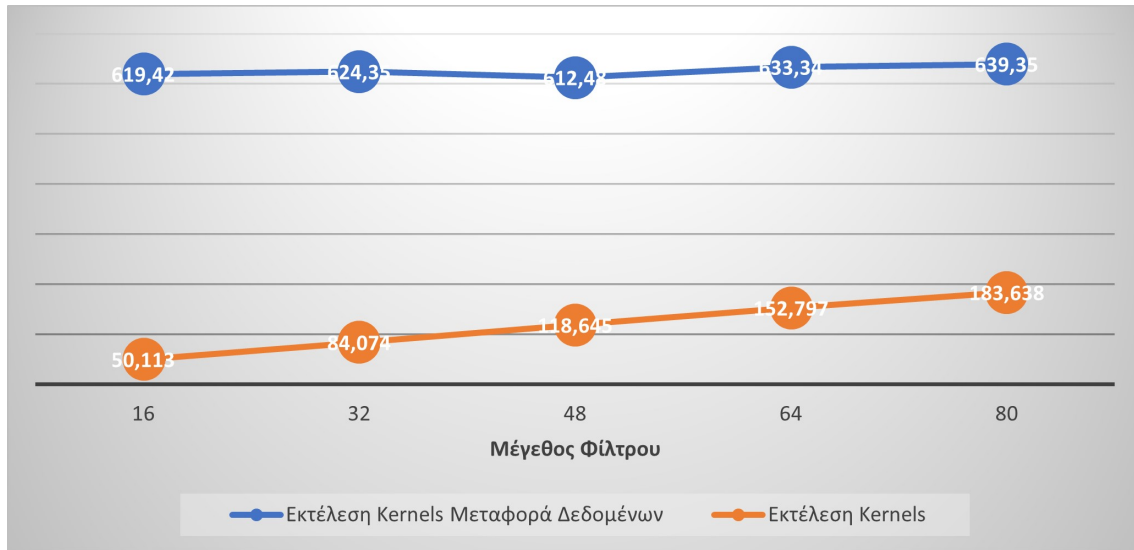


Figure 3: Kernel and Data Transfer diagram (msec)

Filter Radius	Εκτέλεση Kernel	Μεταφορά Δεδομένων
16	50,113	619,42
32	84,074	624,35
48	118,645	612,48
64	152,797	633,34
80	183,638	639,35

Table 3: Kernel and Data Transfer Table

Βήμα 3

Για να καταφέρουμε να υποστηρίξουμε πολύ μεγαλύτερες εικόνες θα πρέπει η εκτέλεση των kernels να είναι blocked, με τέτοιο τρόπο ώστε όλα τα δεδομένα που απαιτούνται για την εκτέλεση κάθε βήματος να χωράνε στη GPU. Για αυτό το λόγο η αρχική εικόνα σπάει σε μικρότερα κομμάτια και σε καθένα απο αυτά εφαρμόζεται ξεχωριστά η διασάσταση συνέλιξη. Το μικρότερο κομμάτι της εικόνας που δίνεται για υπολογισμό στους kernels αποθηκεύεται σε ένα πίνακα σε μνήμη του host και μετά με την cudaMemcpy μεταφέρεται σε μνήμη του device. Αφου ολοκληρωθεί ο υπολογισμός της δισδιάστατης συνέλιξης αυτού του κομματιού τότε περνάει στη θέση του στη μνήμη του device το επόμενο κομμάτι. Έτσι τώρα δεν μας περιορίζει πια η global memory του device αλλά η μνήμη του host καθώς έχει αποθηκεύεται η συνολική εικόνα.

Το μέγεθος του κομματιού (παράμετρος BLOCKING) μπορεί να μεταβληθεί, αλλά πρέπει να παίρνει τιμές που είναι δύναμη του 2 και να ναι μικρότερο από το μέγεθος της εικόνας. Στα πειράματά μας χρησιμοποιήσαμε μέγεθος BLOCKING 512.

Αυτό το βήμα είχε αρνητική επίδραση στον χρόνο εκτέλεσης, κάτι το οποίο είναι λογικό καθώς πλέον γίνονται πολύ περισσότερες κλήσεις των kernels και μεταφορές δεδομένων.

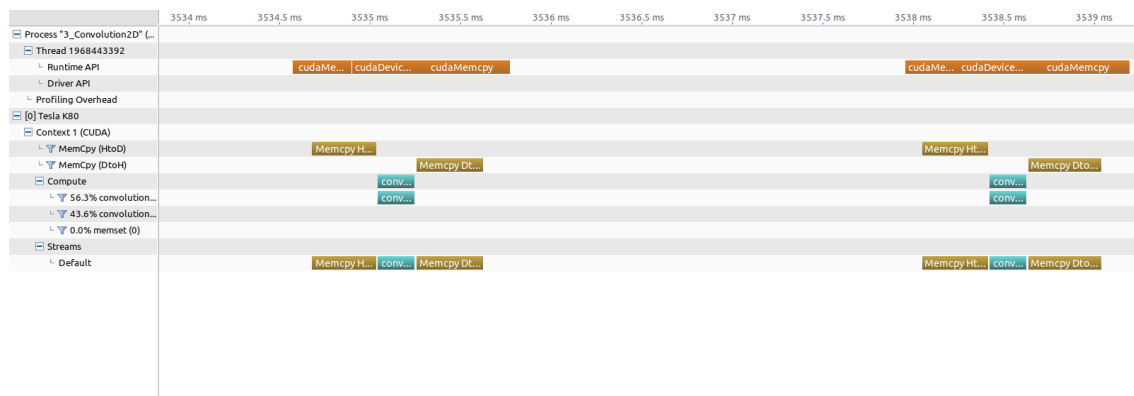


Figure 4: Profiling Row Convolution Step 3

Βήμα 4

Όπως παρατηρήσαμε απο το profiling στον κώδικα του προηγούμενου βήματος δεν υπάρχει καμία επικάλυψη ανάμεσα στον κώδικα που εκτελείται από διαφορετικούς kernels ή σε εκτέλεση κώδικα και μεταφορά δεδομένων. Σε αυτό το βήμα για να πετύχουμε επικάλυψη χρησιμοποιούμε 2 streams. Κάθε stream μεταφέρει ασύγχρονα ένα κομμάτι της αρχικής εικόνας στο device με τη χρήσης της cudaMemcpyAsync πλέον, καλείται ο kernel για υπολογισμός της συνέλιξης και το αποτέλεσμα έπειτα μεταφέρεται πάλι ασύγχρονα απο το device στον host.

Όπως είναι αναμενόμενο ο χρόνος εκτέλεσης μειώνεται απο το προηγούμενο βήμα αρκετά (για εικόνα 8192x8192 και ακτίνα φίλτρου 31 ο χρόνος πέφτει από 0.141 s σε 0.123 s) και όπως παρατηρούμε παρακάτω απο το στιγμιότυπο του profiling πλέον υπάρχει επικάλυψη μεταξύ των εκτελέσεων των kernels και των μεταφορών δεδομένων.

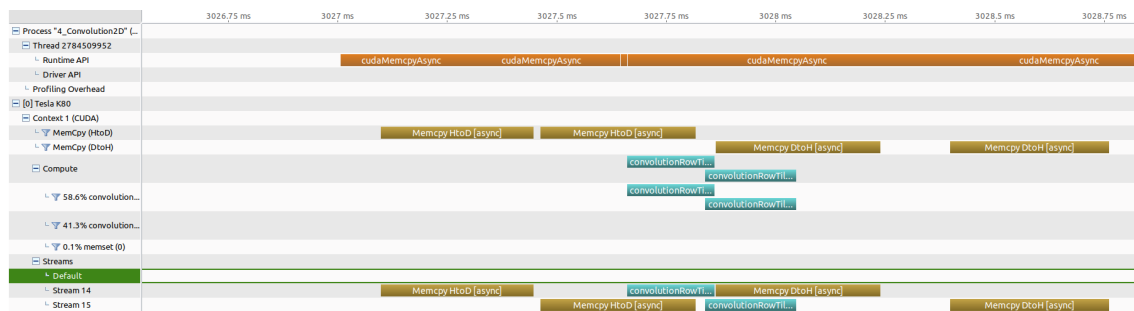


Figure 5: Profiling Row Convolution Step 4 with Streams

Βήμα 5

Το μέγιστο μέγεθος εικόνας που μπορεί να υποστηριχθεί πλέον είναι 65536x65536. Για εικόνα μεγέθους 131072x131072 εμφανίζεται το μήνυμα: Error while allocating host memory .Ο πόρος που μας περιορίζει πλέον δεν είναι η μνήμη του device αλλά η μνήμη του host η οποία με την εντολή cat /proc/meminfo βλέπουμε ότι είναι 128 GB.