

Ιόνιο Πανεπιστήμιο – Τμήμα Πληροφορικής
Παράλληλος Προγραμματισμός
2023-24

Προγραμματισμός σε GPUs

(Το προγραμματιστικό μοντέλο CUDA)

<https://mixstef.github.io/courses/parprog/>

*Το εργαστήριο του μαθήματος χρησιμοποιεί υπολογιστικούς πόρους
AWS Cloud χρηματοδοτούμενους από το ΕΔΥΤΕ*

Μ.Στεφανιδάκης

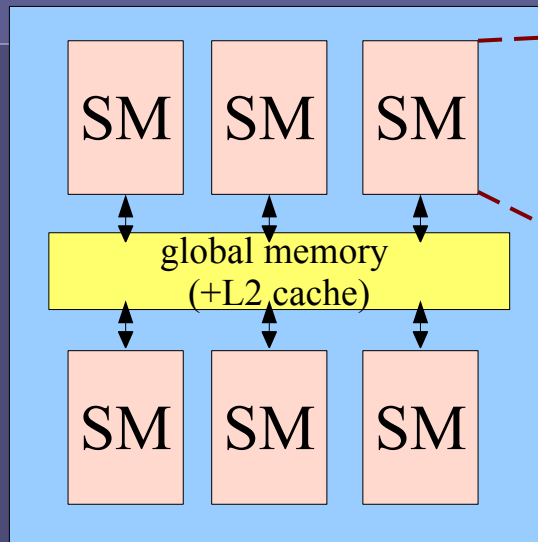


Προγραμματισμός σε GPUs: ιδιαιτερότητες

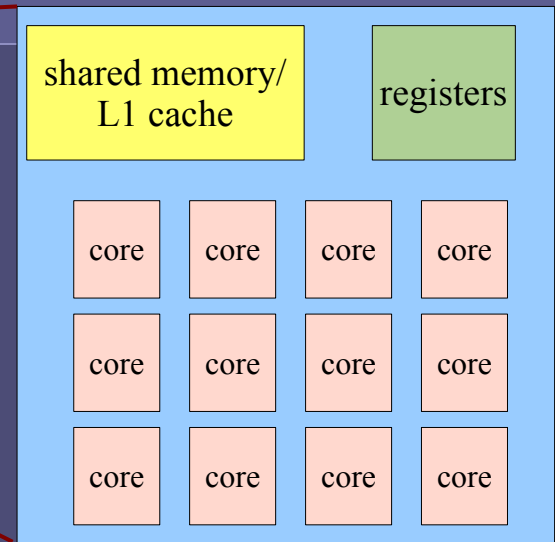
- Πού εκτελείται το πρόγραμμα;
 - Ένα μέρος στη CPU (**host**)
 - και ένα μέρος σε κάποια GPU/accelerator (**device**)
- Που βρίσκονται τα δεδομένα;
 - Στην κύρια μνήμη του υπολογιστή
 - και σε κάποια μνήμη της GPU
 - Μια τυπική GPU έχει διάφορα είδη μνήμης
 - Ο προγραμματιστής πρέπει να φροντίζει για την επιλογή της πιο κατάλληλης
 - Συνεπώς προκύπτει η ανάγκη μεταφοράς δεδομένων μεταξύ διαφορετικών μνημών
 - Πρέπει να συνυπολογίζεται το κόστος μεταφοράς στη συνολική απόδοση

Οργάνωση τυπικής GPU

GPU



SM

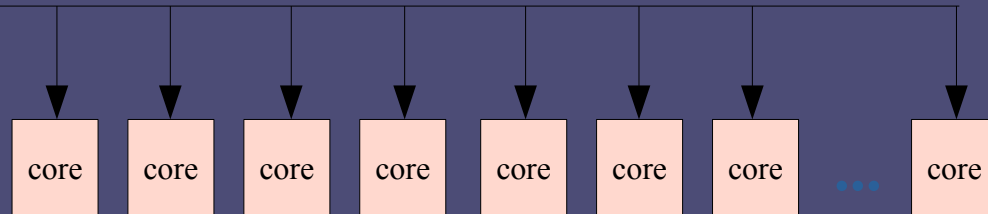


- Μια GPU αποτελείται από Streaming Multiprocessors (SMs)
 - Πρόσβαση όλων των SM σε global memory
- Κάθε SM διαθέτει
 - Έναν αριθμό υπολογιστικών cores
 - Απλή λογική λειτουργίας, σχεδιασμένα να εκτελούν μαζικά πράξεις (ίδια πράξη σε μεγάλες ομάδες cores)
 - Ένα ξεχωριστό set καταχωρητών
 - Μια γρήγορη κοινή μνήμη (shared memory)

GPUs και threads

program counter
instruction decoder

ίδια εντολή (διαφορετικά δεδομένα)



- Στις GPUs η έννοια του thread είναι διαφορετική απ' ό,τι σε μια CPU
 - Διαθέσιμος (μαζικά) μεγάλος αριθμός **hardware threads** (cores)
 - Παράλληλη εκτέλεση της **ίδιας εντολής** από **πολλά threads** την **ίδια στιγμή** (μοντέλο SIMT)
 - Οι πόροι (state) των εκτελούμενων threads βρίσκονται συνεχώς στους καταχωρητές του SM
 - Κάθε SM χρονοδρομολογεί διαδοχικά ομάδες threads με τον **ίδιο program counter** (PC) στα cores

Ιδιαιτερότητες εκτέλεσης SIMT

- Τι συμβαίνει στις διακλαδώσεις;

```
if cond {  
    A;  
}  
else {  
    B;  
}
```

- Η ομάδα των threads με ίδιο PC πρέπει να εκτελεστεί **δύο φορές**
 - Την πρώτη φορά με ενεργοποιημένα μόνο τα threads που εκτελούν τον κώδικα A
 - και τη δεύτερη φορά με ενεργοποιημένα τα threads που εκτελούν το B

Ιδιαιτερότητες εκτέλεσης SIMT (2)

- Αν ένα thread της ομάδας δεν μπορεί να προχωρήσει (π.χ. αναμονή για πόρους);
 - Όλη η ομάδα threads πρέπει να περιμένει
 - Το SM επιλέγει άλλη ομάδα threads που είναι έτοιμη να εκτελεστεί
 - Προσοχή στα barriers: πρέπει να εκτελούνται εκτός διακλαδώσεων
 - Δηλαδή, από όλα τα threads της ομάδας
- Τα παραπάνω προσπαθούν να διορθώσουν νεώτερες αρχιτεκτονικές GPU
 - Επιτρέποντας ανεξάρτητο PC ανά thread
 - Η μέγιστη απόδοση όμως επιτυγχάνεται όταν όλα τα threads της ομάδας εκτελούν τον ίδιο κώδικα

Το προγραμματιστικό μοντέλο CUDA

- Ένας τρόπος «αφαιρετικής» περιγραφής του hardware (threads, cores, SMs) σε ένα προγραμματιστικό interface
 - Πλαισιώνεται από τις βιβλιοθήκες και τα εργαλεία που το υλοποιούν
- Ο στόχος: να καλυφθούν με ενιαίο τρόπο GPUs με διαφορετικά χαρακτηριστικά και δυνατότητες
 - Η γνώση όμως των ιδιοτεροτήτων του hardware εκτέλεσης είναι σε μεγάλο βαθμό αναγκαία
 - Ο προγραμματισμός σε CUDA δεν είναι εύκολος εάν θέλουμε να επιτύχουμε τη μέγιστη απόδοση...

Ορολογία CUDA

- **Kernel**
 - Μια συνάρτηση που περιγράφει τι θα εκτελεστεί στη GPU από κάθε ένα thread
- **Thread**
 - Η μικρότερη μονάδα εκτέλεσης του κώδικα του kernel
- **Block (Cooperative Thread Array – CTA)**
 - Μια ομάδα threads που εκτελούν τον ίδιο kernel στο ίδιο SM
 - Αν υπάρχουν πόροι (hardware), το SM μπορεί να εκτελεί και άλλο block παράλληλα
- **Grid**
 - Το σύνολο των blocks που εκτελούν τον kernel, όχι κατ' ανάγκη ταυτόχρονα
 - Σε πρόσφατες GPU υπάρχει και το **cluster**, μεταξύ block και grid

CUDA block

- Μια ομάδα από threads που εκτελούνται ταυτόχρονα σε ένα SM
 - Τα threads αυτά μοιράζονται τους καταχωρητές (registers) και την κοινή μνήμη (shared memory) του SM
 - Κάθε thread έχει δεσμευμένο το δικό του μέρος από τους παραπάνω πόρους όσο διαρκεί η εκτέλεση του block
- Τα threads ενός block μπορούν να συγχρονιστούν μεταξύ τους
 - Αντιθέτως, δεν υπάρχει τρόπος να συγχρονιστούν διαφορετικά blocks μεταξύ τους
 - Ούτε μπορούμε να υποθέσουμε με ποια σειρά θα εκτελεστούν τα blocks

Δήλωση kernel

- Χρήση του keyword `__global__`

```
// a kernel function - must return void
__global__ void add(int a,int b,int *c) {
    *c = a+b;
}
```

- CUDA C/C++
 - Προσθήκη επεκτάσεων (extensions) της CUDA
 - Ο compiler της CUDA (`nvcc`), στέλνει την κλασσική C/C++ στον gcc (ή αντίστοιχο) ενώ χειρίζεται διαφορετικά τις επεκτάσεις της CUDA

Εκτέλεση kernel

- Ο προγραμματιστής ζητά την εκτέλεση ενός kernel με συγκεκριμένο αριθμό **blocks ανά grid** και **threads ανά block**
 - Τα blocks που αποτελούν το grid του kernel κατανέμονται στα SM της GPU
- Τα SM εκτελούν τα blocks που τους αντιστοιχούν, το ένα μετά το άλλο ή παράλληλα (ανάλογα με τους διαθέσιμους πόρους)
 - Κάθε SM χρονοδρομολογεί ομάδες από threads του ίδιου block για εκτέλεση από τα cores του
 - Σε ομάδες με τον ίδιο PC (“warps” στην ορολογία CUDA, 32 threads με την τρέχουσα τεχνολογία)

Εκκίνηση kernel από host (CPU)

- `kernel_name<<blocks-per-grid,threads-per-block>>(arg, ...)`

```
// this call is asynchronous - host continues execution  
add<<<1,1>>>(2,7,dev_c);
```

- Η εκκίνηση («κλήση») του kernel γίνεται από το κυρίως πρόγραμμα που εκτελείται στη CPU (host)
- Η κλήση είναι **ασύγχρονη** – η CPU συνεχίζει την εκτέλεση των επόμενων εντολών χωρίς να περιμένει την ολοκλήρωση της εκτέλεσης του kernel στην GPU

Τυπική ροή εκτέλεσης

- Δέσμευση μνήμης στη συσκευή GPU (device)
 - Για την υποδοχή των δεδομένων εισόδου από host
- Μεταφορά δεδομένων στη GPU
 - Στην βασική μορφή CUDA, ευθύνη του προγραμματιστή
- Εκτέλεση ενός kernel
 - Ανεξάρτητα και ασύγχρονα από CPU (host)
- Μεταφορά των αποτελεσμάτων πίσω στη μνήμη του υπολογιστή (host)
 - Συγχρονισμός με την ολοκλήρωση του kernel
- Αποδέσμευση μνήμης στη συσκευή GPU

Παράδειγμα

```
int c; // host's (CPU) 'c' variable
int *dev_c; // ptr to device's (GPU) 'c' variable

// allocate space for 'c' on device's memory
cudaMalloc((void **)&dev_c,sizeof(int))

// call the kernel on device, 1 block/1 thread
// this call is asynchronous - host continues execution
add<<<1,1>>>(2,7,dev_c);

// transfer device's 'c' into host's 'c' - synchronous call,
// waits until kernel is done
cudaMemcpy(&c,dev_c,sizeof(int),cudaMemcpyDeviceToHost);

// free memory of device's c
cudaFree(dev_c);
```

- **Προσοχή:** Οι συναρτήσεις της CUDA επιστρέφουν ένδειξη επιτυχίας η όχι – σε κανονικό πρόγραμμα θα πρέπει να ελέγχουμε εάν κάθε κλήση ήταν επιτυχής!

Σφάλματα του kernel

- Δεν υπάρχει τρόπος να ειδοποιηθεί η CPU (host) για σφάλματα κατά την εκκίνηση και εκτέλεση ενός kernel
 - Εκτελείται ασύγχρονα, σε διαφορετικό device (GPU)
- Πώς θα ελέγξουμε αν υπήρξαν σφάλματα στον kernel;
 - Σε αντίθεση με τις συναρτήσεις των CUDA APIs, οι οποίες επιστρέφουν ένα error code
- Μπορούμε να ελεγχουμε αν υπήρξαν σφάλματα μετά την ολοκλήρωση του kernel

```
// a catchall msg here, will catch kernel launch failures, too!  
printf("Last error msg is: %s\n",  
      cudaGetErrorString( cudaGetLastError() ));
```

Παραμετρική εκτέλεση

- Πώς γνωρίζει κάθε thread ποιο μέρος της συνολικής εργασίας θα εκτελέσει
 - Κάθε thread έχει διαθέσιμες κατά την εκτέλεση του kernel μια σειρά από **μεταβλητές index** (read-only)
 - Οι μεταβλητές αυτές μπορούν να οργανωθούν σε 1, 2 ή 3 διαστάσεις
 - Προγραμματιστική ευκολία, για να ταιριάζουν με το είδος (και την τοπικότητα των δεδομένων) της εφαρμογής
 - Δεν αλλάζει η υποκείμενη οργάνωση του hardware σε threads/blocks/grid

Οργάνωση σε μία διάσταση

- **threadIdx.x**
 - «Ποιο thread του block είμαι»
- **blockIdx.x**
 - «Σε ποιο block ανήκω»
- **blockDim.x**
 - «Πόσα threads υπάρχουν σε κάθε block»
- **gridDim.x**
 - «Πόσα blocks υπάρχουν στο grid»
- Σε οργανώσεις με 2 ή 3 διαστάσεις υπάρχουν επιπλέον τα **.y** και **.z** των παραπάνω

Περιορισμοί διαστάσεων

- `/usr/local/cuda-12.2/extras/demo_suite/deviceQuery`

```
Maximum number of threads per multiprocessor: 2048  
Maximum number of threads per block: 1024  
Max dimension size of a thread block (x,y,z): (1024, 1024, 64)  
Max dimension size of a grid size (x,y,z): (2147483647,  
65535, 65535)
```

- Κάθε μοντέλο GPU έχει περιορισμούς στον μέγιστο αριθμό threads και blocks ανά διάσταση, όπως και στο πόσα threads/block ή blocks/grid μπορούν να εκτελεστούν

Παραδείγματα κατανομής εργασίας

- Έστω ότι εκτελούμε απλό μετασχηματισμό (map) σε N στοιχεία εισόδου
 - Πώς κατανέμω την εργασία;
- Λύση #1: Ένα και μοναδικό thread εκτελεί τα πάντα
 - Το γνωστό loop `for (i=0; i<N; i++)`
 - Δεν είναι λύση στην πραγματικότητα (αδικαιολόγητη σπατάλη των διαθέσιμων πόρων της GPU...)

Παραδείγματα κατανομής εργασίας (2)

- **Λύση #2: Ένα block με K threads**
 - Πόσο θα είναι το K ;
 - Συνήθως ένα (μικρό) πολλαπλάσιο του 32 (warp size)
 - Π.χ. με $K = 256$ θα έχουμε 8 warps προς εκτέλεση, ικανά να κρύψουν καθυστερήσεις σε κάποιο/α από αυτά
 - Πώς καλύπτουμε τα N στοιχεία εισόδου;
 - Striding
 - Το thread0 χειρίζεται τα στοιχεία 0, 256, 512, ...
 - Το thread1 χειρίζεται τα στοιχεία 1, 257, 513, ...
 - Το thread2 χειρίζεται τα στοιχεία 2, 258, 514, ...
 - Η ομοιομορφία στην προσπέλαση είναι απαραίτητη για την υψηλή απόδοση
 - Όμως με ένα ενεργό block, χρησιμοποιούμε μόνο ένα SM...

Παραδείγματα κατανομής εργασίας (3)

- **Λύση #3: M blocks με K threads το καθένα**
 - Το M πρέπει να είναι τέτοιο ώστε $M \times K \geq N$
 - Κάθε thread υπολογίζει ένα (ή κανένα) στοιχείο μόνο!
 - Υπολογισμός του M
 - Το γνωστό $(N + K - 1) / K$
 - Κρατώντας το $K = 256$
 - Ποιο στοιχείο χειρίζεται κάθε thread;
 - $i = \text{αριθμός-block} * \text{μέγεθος-block} + \text{αριθμός thread}$

Παραδείγματα κατανομής εργασίας (4)

- **Λύση #4: M blocks με K threads το καθένα**
 - Το M είναι σταθερό, π.χ. 32 φορές τον αριθμό των SM
 - Κάθε thread υπολογίζει πολλαπλά στοιχεία

```
int threads = 256;
int devId;
cudaGetDevice(&devId);
int numSM;
cudaDeviceGetAttribute(&numSM, cudaDevAttrMultiProcessorCount,
devId);
int blocks = numSM*32; // as a multiple of SMs in GPU
```

- Ποια στοιχεία χειρίζεται κάθε thread;
 - $start = \text{αριθμός-block} * \text{μέγεθος-block} + \text{αριθμός thread}$
- Επανάληψη με $stride = \text{συνολικό αριθμό threads στο grid}$
 - $\text{μέγεθος block (σε threads)} * \text{μέγεθος grid (σε blocks)}$