

1) Heterogeno računarstvo (GPU, CPU)

$$\text{index} = \text{redniBroj_niti_u_bloku} + \text{redniBroj_bloka} * \text{dimenzija_bloka}$$

GPU nema virtuelnu memoriju, prekide, i saradnju sa uređajima poput tastature i misa. Veoma je neefikasna u slučaju da posao koji joj je dodeljen nije tipa SIMD

Heterogeno računarstvo : Multicore CPU, Manycore GPU i rekonfigurabilne jedinice (FPGA)

GPU - CPU

- CPU -- Optimizovan za malo kasnjenje zbog cache memorija, poseduje kontrolnu logiku za izvršenje preko reda i predikcije, desetine niti, dobar za sve vrste obrada

- Cilj CPU - minimizovati kasnjenje (cache i kompleksna logika)

- GPU -- Optimizovan za paralelnu obradu i propusnost, tolerise mem. kasnjenja, Hiljade niti (više tranzistora), dobar za paralelnu obradu

- Cilj GPU - maksimizirati propusnost (multithreading može da skrije kasnjenje, kontrolna logika je ista za više niti)

GPU se razlikuje od CPU-a po samom dizajnu, njihova kompleksnost je manja. Optimizovana je za compute-intensive i highly parallel proračune.

2) Konkurentno vs Paralelno

Konkurentno (kvazi paralelno):

- 1 CPU - 1 Core

Cilj:

- bolje iskoristiti procesor
- poboljšati odzivnost sistema
- podržati sto više korisnika

Problem:

- Prioriteti
- Raspoređivanje

Paralelno (više niti se izvršava na više jezgara):

- više CPU - više Core-ova

Cilj:

- Paralelno izvršavanje
- Speedup

Problemi:

- Paralelni algoritmi i strukture podataka
- Mapiranje i balansiranje opterećenja

Zajednički problemi: sinhronizacija, komunikacija, race-condition, deadlock, konzistentnost podataka

3) Nivoi paralelizma

- Na nivou instrukcije
 - izvršenje preko reda, predikcije grananja
- Na nivou podataka
 - vektorski procesori, SIMD izvršavanje - GPU
- Na nivou niti
 - povećanje broja jezgara, niti, procesa

4) GPU - memorijska hierarhija (tipovi, vidljivost)

- Registri - Najbrzi pristup - dostupni tokom života niti
- Lokalna - sporiji od registara i deljene memorije, dostupna tokom života niti - nalazi se u globalnoj memoriji
- Deljena - Može biti brza kao i registarska ako nema konflikta. Dostupna svim nitima iz blok (tokom života bloka)
- Globalna - sporija od registara i deljene memorije, dostupna i device-u i host-u (životni vek app)
- L1 - kesira pristup globalnoj memoriji - brzina ista kao i kod deljene memorije
- L2 - upravlja atomske operacije - sporiji od L1, ali i dalje brzi pristup nego kod globalne memorije.

5) GPU - hierarhija niti

- Nit se izvršava na jednom jezgri (1 Warp = 32 niti) - jedinstven ID
- Blok se izvršava na jednom SM-u (blokovi ne migriraju) - jedinstven ID
- Grid se izvršava na jednom Device-u - ~~do 16 kernela (gridova) po uređaju - sekvencijalno se izvršavaju.~~

Blok se dodeljuje SM. U SM-u blok se deli na Warp-ove (32/64 niti). Svih 32 niti u Warpu moraju izvršavati isti set instrukcija. Warp-ovi se izvršavaju konkurentno na SM-u. - Nit se izvršava na jednom jezgri, grupa od 32 niti se zove warp više warp-ova čine blok koji se izvršava na SM-u. Blokovi se ne prključuju između SM-a. Više blokova se može izvršavati na istom SM-u - konkurentno. Grupa blokova se zove grid i izvršava se na device-u (sekvencijalni kerneli).

6) Paralelni programi - odnos GPU-CPU

- Imamo master proces koji je izvršava na CPU i izvršava sledeće korake:
 - Inicijalizacija device-a
 - Alokacije memorije na hostu i device-u
 - Kopiranje podataka sa hosta na device
 - Pokretanje više instanci izvršavanja na device-u
 - Kopiranje podataka sa device-a na host
 - 3-5 ponavljanja prethodnog postupka
 - Dealokacija memorije i kraj programa!
- Na GPU:
 - Svaki blok se izvršava na Streaming Multiprocessor-u (SM)
 - Ako je broj blokova veći od broja SM-a, onda će se samo jedan blok izvršavati u jednom trenutku (queue)
 - Sve niti jednog bloka mogu da pristupe deljenoj memoriji, ali ne mogu videti šta rade druge niti
 - Ne možemo znati redosled izvršavanja niti i blokova.

7) CUDA - osnovni koncepti

- Proizvedena od strane Nvidia-e
- Programski model: Prenos podataka na GPU, Izvršavanje, Pruzimanje rezultata
- Racunarska arhitektura opšte namene.
- Izvorni kod host-a i device-a se nalazi u istom fajlu (u OpenCL-u nije tako)
- CUDA je racunarska arhitektura opšte namene - uključuje CUDA set instrukcija i paralelni engine na GPU.
- Koristi C/C++ uz minimalne izmene - minimalni nivo apstrakcije
- Relativno low-level pristup: Apstrakcija CUDA-e je usko vezana za mogućnosti/performance GPU-a
- CUDA na Nvidia karticama, OpenCL na bilo kojim GPU-ovima i CPU-ovima.
- Apstrakcija: Hierarhija niti, "Jednostavne" metode sinhronizacije, Deljena memorija za saradju između niti
- CUDA virtualizuje HW tako što niti virtualizuje skalarni procesor a blok virtualizuje multiprocesor. Sve se raspoređuje na fizički HW GPU-a. Niti i Blokovi se pokreću i izvršavaju do kraja, blokovi su nezavisni.

8) GPU - virtuelizacija

GPU virtualizuje HW tako što niti virtualizuje skalarni procesor a blok virtualizuje multiprocesor. Sve se raspoređuje na fizički HW GPU-a. Niti i Blokovi se pokreću i izvršavaju do kraja, blokovi su nezavisni.

9) Tipovi promenljivih u CUDA-i i njihova vidljivost

Variable declaration	Memory
<code>int var;</code>	register
<code>int array_var[10];</code>	local
<code>__shared__ int shared_var;</code>	shared
<code>__device__ int global_var;</code>	global
<code>__constant__ int constant_var;</code>	constant

10) Kreiranje kernela

- definisanje `__global__` kernel funkcija (programa) koje će se izvršavati na GPU.
- alokacija memorije na GPU i oslobađanje memorije na GPU nakon završenog posla.
- kreiranje (startovanje) kernela sa definisanim brojem niti, blokova, i argumentima
`myKernel<<<nBlocks, nThreads>>>(par1, par2)`
- sinhronizacija CPU i GPU (čekamo da se izvrši kod na GPU)

11) GPU - atomic i sinhronizacija

Na GPU ne možemo koristiti klasične load i store operacije zbog Race-Condition-a, umesto toga koriste se atomske instrukcije. Atomske operacije nam garantuju da samo jedna nit ima pristup memoriji sve dok se operacija ne izvrši do kraja (neprekidno). Definicija atomskih instrukcija se postiže ključnom reči `atomic`. CUDA obezbeđuje atomske operacije nad globalnom memorijom.

Sinhronizacija:

- Sve niti u warp-u se moraju eksplicitno sinhronizovati.

- `__syncthreads()` se koristi unutar kernela kako bi se zaustavile sve niti i sačekale da sve budu na istom mestu. Veoma je bitno kada se piše kod, da se obezbedi da se sve niti mogu izvršiti do sinhronizacije u protivnom nastaje deadlock. Između blokova nema sinhronizacije.

Ako hocemo da sinhronizujemo rad hosta i device-a tako da se nalaze na istoj tacki izvršenja koristimo: `cudaDeviceSynchronize();`

- osigurava da se svi asinhroni zadaci završe pre nastavka rada!

`cudaStreamSynchronize();` - za sinhronizaciju stream-ova kernela (više kernela).

Blokovi moraju biti nezavisni. Sva moguća preplitanja blokova bi trebala da budu validna:

- izvršenje do kraja bez prekida, bilo koji redosled izvršavanja, mogu se izvršavati konkurentno ili sekvencijalno.

Blokovi nisu sinhronizovani ali mogu biti koordinisani: deljeni pokazivac u redu čekanja.

12) Strategija obrade podataka na GPU

- Podela podataka na delove koji se mogu smestiti u deljenu memoriju
- Svakim delom podatak rukuje jedna blok
- Iz globalne memorije se učitavaju podaci u deljene memorije blokova
- Izvode se proračuni na blokovima, sa datim delovima podataka
- Rezultati se kopiraju nazad - iz deljene u globalnu memoriju.

13) Tipovi adresnih prostora na GPU

Ugrađeni primitiv **memcpy(cudaMemcpy)** nam omogućuje kopiranje podataka između adresnih prostora na host-u i device-u.

Device ima 3 tipa adresnih prostora: **private** - za svaku nit, **shared** - za svaki blok, **global** - za svaki program.

14) Masivni paralelizam (fina i gruba podela posla)

Masivni paralelizam:

- Upotreba par stotina/hiljada blokova

Blok se izvršava na jednom SM-u. Potreban nam je veliki broj blokova kako bi upotrebili 10-ak SM-a jer jedan SM izvršava efikasno 2-8 blokova konkurentno. Ogroman broj blokova nam treba kako bi posao skalirali na više GPU-ova.

Ovu logiku koristimo kada želimo ugrubo da podelimo posao (task-ove)

- Upotreba stotine niti po bloku

Svaka nit se izvršava na jednom jezgri. Treba nam do 512 niti po SM-u kako bi koristili sva jezgra istovremeno. Koristi se umnožak od 32 niti (warp) po bloku. Finija podela posla - paralelizam na nivou podatak, niti, instrukcija.

Zaključak: Više blokova - grublja podela posla, Više niti po bloku - finija podela posla.

15) SIMT - koncept

SIMT - niti se izvršavaju u grupi od 32 niti - warp. Sve niti u warp-u izvršavaju iste instrukcije. HW multithreading - obuhvata alokaciju resursa na HW-u i raspoređivanje niti. HW se oslanja na niti kako bi sakrio kašnjenja. Context-switch između warp-ova je bez overheada.

SIMT:

- Warp - set od 32 paralelne niti koje izvršavaju jednu instrukciju. SIMT - dodeljuje instrukciju warp-u

SM sva jezgra deli na dva warp-a, zato ima i dva raspoređivaca. SIMT warp izvršava svaku instrukciju na 32 niti. Predikati omogućuju ili onemogućuju individualno izvršavanje niti. Stek upravlja grananjima za svaku nit. Redundantna izračunavanja su brza od neispravnog grananja (prilikom grananja izračunava se obe strane jer je to brže nego da izaberemo pogrešno grananje pa da se ponovo vratimo i računamo).

16) Occupancy

Occupancy: odnos broja aktivnih warp-ova po SM-u u odnosu na maksimalan broj warp-ova po SM-u.

- Treba da je sto veća kako bi se sakrilo kašnjenje
- Limitirana je upotrebom resursa
- Zavisi od arhitekture ali i od same aplikacije.

Occupancy optimization:

- Trenutnu zauzetost saznajemo preko profajlera.

- Prilagođavanje upotrebe resursa kako bi se zauzetost povećala (promena veličine bloka, limitiranje upotrebe registara, dinamička alokacija deljene memorije)

17) Deljena memorija i optimizacija

- Za komunikaciju izmedju niti u bloku
- Koristi sekao kes kako bi se zaobisao pristup globalnoj memoriji (kod vise pristupa istom podatku)
- Izbegavanje pristupa koji nije po redosledu
- Kasnjenje je od 4-8 puta manje nego u slucaju L2 I DRAM memorija.
- Koristi sinhronizaciju kako bi se izvrila komunikacija - sinhronizacija je jako brza!
- Niti u warp-u pristupaju deljenoj memoriji preko crossbar-mreze. Konflikti mogu da postoje ali oni neznatno uticu na performanse. Favorizuje se upotreba deljene memorije u odnosu na globalnu.

18) Restrikcije CUDA-e u odnosu na C/C++

- moze se pristupiti samo memoriji GPU-a, nema promenljivog broja argumenata, nema statickih promenljivih, nema rekurzije, nema polimorfizama.
- sve funkcije moraju imati kvalifikator:
 - `__global__` - pokrenute od CPU-a, ne mogu biti pozvane sa GPU-a, moraju da vracaju void
 - `__device__` - mogu da se pozovu samo sa GPU-a
 - `__host__` - samo sa CPU-a
 - `__host__` i `__device__` se mogu kombinovati

19) Preporučen način za pristup deljenoj memoriji od strane GPU-a

Niti bi trebale deljenoj memoriji da pristupaju redosledno, po indexu bez preplitanja

Preporuke:

- Koristiti sukcesivne mem. adrese za niti u warp-u
- Koristiti pristup bez presecanja kada god je to moguće (kao kod niza - sekvencijalan pristup po indeksima)
- Load i Store operacije mogu narusiti konzistentnost - atomic nam serializuje pristup mem.

20) Memorijski prostori i pointeri

Najcesce se koriste pointeri, na osnovu vrednosti pointera ne mozemo znati da li je to na GPU-u ili na CPU-u. Dereferenciranje pointera koji pokazuje na memoriju na CPU na GPU-u moze dovesti do kraha celog procesa i obrnuto.

21) Definirati blokove i objasniti njihove osobine

Blok je sacinjen od niti. Niti se grupisu u warpove od po 32 niti, vise warpova cini blok. Blok se izvrsava na SM-u, blokovi ne migriraju izmedju SM-a. Rasporedjivanje blokova na SM-u: HW raspoređuje blokove na dostupne SM-e, nema garancije redosleda, blok ce biti rasporedjen cim neki drugi blok završi! Svaki blok se mapira na jedan ili vise warp-ova. SM moze efikasno da izvrsava 2-8 blokova konkurentno (nalaze se u queue-u). Blokovi moraju biti nezavisni. Blokovi nisu sinhronizovani, ali mogu biti koordinisani. Svaki blok ima deljenu memoriju kojoj mogu da pristupe sve niti iz tog bloka (moze biti brza kao i registrarska, brza od globalne)

22) Preklapanje CPU i GPU computinga i CPU-GPU (GPU-CPU) transfera podataka

- `cudaMemcpy()` izaziva prenos podataka na relaciji GPU - CPU i obrnuto. Ovo se moze preklapati sa periodom procesuiranja GPU ili CPU -> Resenje: u razlicitim kernel-ima kada se na jednom vrsi procesiranje na drugom se vrsi kopiranje podataka (u razlicitim taktovima) - da nije ovako procesuiranje ne bi bilo moguće sve dok se podaci ne kopiraju (Single vs Dual Stream, Overlap)

23) Kako se može optimizovati GPU algoritam množenja matrica

$$A \times B = C$$

1. Transponovaceмо matricu B, kako bi dobili redosledni pristup po indeksima (umesto a_{11}, a_{21}, a_{31} , bice $a_{11}, a_{12}, a_{13}...$)
2. Podelom matrica A i B na blokove, svakom bloku cemo dati npr 32-vrste A i 32-kolone B

24) Ako funkciju definisanu u C, C++ programskom jeziku želimo da koristimo kao kernel funkciju u CUDA programiranju šta teba da uradimo?

Da joj prilikom definicije stavimo prefix `__global__`, alociramo memoriju na GPU, iskopiramo potrebne podatka na GPU, pozovemo je preko def. broja blokova i niti `myKernel<<<nBlocks,nThreads>>>>(par1,par2,par3)`

25) Opisati programsku strategiju rada globalne i deljene memorije grafičkog procesora

- Podela podataka na delove koji se mogu smestiti u globalnu memoriju
- Svakim delom podatak rukuje jedna blok
- Iz globalne memorije se učitavaju podaci u deljene memorije blokova
- Izvode se proračuni na blokovima, sa datim delovima podataka
- Rezultati se kopiraju nazad - iz deljene u globalnu memoriju.

26) Objasniti u rutini cudaMemcpy(void *dst, void *src, size_t nbytes, enum cudaMemcpyKind direction) poslednji argument: enum cudaMemcpyKind direction

Oznacava smer prenosa podataka:

- cudaMemcpyHostToDevice - sa hosta na device
- cudaMemcpyDeviceToHost - sa device na host
- cudaMemcpyDeviceToDevice - sa device-a na device

27) Definirati nit, blok i rešetku (thread, block, grid) i gde se izvršavaju

- Nit se izvršava na jednom jezgri (1 Warp = 32 niti) - jedinstven ID
- Blok se izvršava na jednom SM-u (blokovi ne migriraju, sastoje se iz više warp-ova) - jedinstven ID
- Grid se izvršava na jednom Device-u - do 16 karnela(gridova, više blokova) po uređaju - sekvencijalno se izvršavaju.

28) Objasniti cache data locality

Koristi lokalnost podataka

- vremenska - podatak kome smo sada pristupili verovatno će se ponovo koristiti
- prostorna - podaci u okolini trenutnog podataka će se verovatno koristiti zato ih treba prebaciti u cache.

29) Uporediti brzine rada registra, lokalne memorije, deljene memorije i globalne memorije

- Registrarska - najbrza memorija
- Lokalna - nalazi se u globalnoj 150x sporija nego registrarska
- Deljena - može biti iste brzine kao i registrarska
- Globalna - 150x sporija od registrarske i deljene

30) Kako se Parallel Thread Execution (PTX) pseudo-assembly jezik (kod) koristi u Nvidia CUDA programskom okruženju (kodu)

PTX - asemblerski jezik za CUDA-u. CUDA se prevodi u PTX a drajver GPU-a prevodi PTX u binarni kod. Inline PTX izrazi se mogu koristiti u CUDA-I. PTX definiše viruelnu masinu i set instrukcija opšte namene za paralelno izvršavanje niti.