

Sistemi di Elaborazione Accelerata

Author: Bumma Giuseppe

Professors: Mattoccia Stefano, Tosi Fabio

Università degli studi di Bologna
Ingegneria Informatica Magistrale

2025/2026

The preface to your notes

Contents

1	Introduzione	1
1.1	Perché Scegliere la Piattaforma CUDA?	1
1.2	Cos'è il CUDA Toolkit?	1
1.2.1	Componenti chiave del CUDA Toolkit	1
1.2.1.1	Strumenti di Debugging e Profiling	2
1.2.1.2	Relazione tra CUDA Toolkit e CUDA Version	3
1.2.1.3	Retrocompatibilità del CUDA Toolkit	3
1.3	CUDA Compute Capability (CC)	3
1.3.1	Relazione tra Compute Capability (CC) e CUDA Version	3
1.4	Evoluzione delle Architetture GPU NVIDIA	4
1.5	Anatomia di un Programma CUDA	4
1.5.1	Struttura del Codice Sorgente	4
1.5.2	Componenti Principali	4
1.5.3	Flusso di Compilazione	5
1.6	Compilazione di un Programma CUDA	6
1.7	Hello World in CUDA C	6
1.7.1	Passo 1: Creare il File Sorgente	6
1.7.1.1	Passo 2: Scrivere il codice	6
1.8	Hello World in CUDA C - Analisi	7
1.9	Hello World in CUDA C	7
1.9.0.1	Passo 3: Compilazione	7
1.10	Hello World in CUDA C	8
1.10.0.1	Passo 4: Esecuzione	8
1.11	Ottenere Informazioni sulla GPU tramite API CUDA	8
1.11.1	Utilizzo delle API CUDA	8
1.12	Ottenere Informazioni sulla GPU tramite API CUDA	9
1.12.0.1	Esempio di Output	9
1.13	Cosa Significa Programmare in CUDA?	9
1.13.1	Pensare in Parallelo	9
1.13.2	Scrittura di codice in CUDA C	9
1.13.2.1	Testi Generali	10
1.13.2.2	NVIDIA Docs	10
1.13.2.3	Materiale di Approfondimento	10
2	Modello di Programmazione CUDA	11
2.1	La Struttura Stratificata dell'Ecosistema CUDA	11
2.2	Ruolo del Modello e del Programma	11
2.3	Livelli di Astrazione nella Programmazione Parallela CUDA	12
2.3.1	Livello Dominio	12
2.3.2	Livello Logico	12
2.3.2.1	Livello Hardware	12
2.4	Thread CUDA: L'Unità Fondamentale di Calcolo	12
2.4.1	Cos'è un Thread CUDA?	12
2.4.2	Cosa Fa un Thread CUDA?	13
2.5	Struttura di Programmazione CUDA	13
2.5.1	Caratteristiche Principali	13
2.6	Flusso Tipico di Elaborazione CUDA	13

2.7	Gestione della Memoria in CUDA	14
2.7.1	Modello di Memoria CUDA	14
2.7.2	Caratteristiche PCIe	14
2.8	Collegamento Fisico della GPU tramite PCIe	15
2.8.1	Connessione Fisica GPU	15
2.9	Modello di Memoria CUDA	15
2.10	Gerarchia di Memoria	15
2.11	Allocazione della Memoria sul Device	16
2.11.1	Esempio di Allocazione di Memoria sulla GPU	17
2.12	Trasferimento Dati	17
2.12.1	Spazi di Memoria Differenti	18
2.13	Deallocazione della Memoria sul Device	18
2.13.1	Esempio di Allocazione e Trasferimento Dati	19
2.14	Organizzazione dei Thread in CUDA	19
2.14.1	Struttura Gerarchica	19
2.14.2	Perché una Gerarchia di Thread?	20
2.15	Identificazione dei Thread in CUDA	21
2.15.1	Dimensione delle Griglie e dei Blocchi	21
2.16	Struttura <code>dim3</code>	22
2.17	Esecuzione di un Kernel CUDA	25
2.17.1	Cos'è un Kernel CUDA?	25
2.18	Qualificatori di Funzione in CUDA	26
2.18.1	Combinazione dei qualificatori host e device	26
2.19	Kernel CUDA: Regole e Comportamento	26
2.20	Configurazioni di un Kernel CUDA	27
2.20.1	Combinazioni di Griglia 1D (Esempi)	27
2.20.2	Combinazioni di Griglia 2D (Esempi)	27
2.20.3	Combinazioni di Griglia 3D (Esempi)	28
2.21	Numero di Thread per Blocco	28
2.22	Compute Capability (CC) - Limiti SM	29
2.23	Identificazione dei Thread in CUDA	29
2.23.1	Esempio Codice CUDA	29
2.24	Tecniche di Mapping e Dimensionamento	29
2.24.1	Somma di Array in CUDA	29
2.25	Confronto: Somma di Vettori in C vs CUDA C	30
2.26	Identificazione dei Thread e Mapping dei Dati in CUDA	32
2.26.1	Struttura dei Dati e Calcolo dell'Indice Globale	32
2.26.2	Calcolo dell'Indice Globale del Thread - Grid 1D, Block 1D	32
2.27	Calcolo dell'Indice Globale del Thread - Grid 1D, Block 2D	34
2.28	Calcolo dell'Indice Globale del Thread - Grid 2D, Block 2D	35
2.29	Metodo Basato su Coordinate per Indici Globali in CUDA	35
2.29.1	Calcolo degli Indici Coordinati	36
2.29.2	Esempio di Utilizzo (Caso 2D)	36
2.29.3	Come Calcolare la Dimensione della Griglia e del Blocco	36
2.29.4	Esempio 1 (Dati Residui): <code>dataSize = 1030, blockSize = 256</code>	37
2.29.5	Esempio 2 (Multiplo Perfetto): <code>dataSize = 1024, blockSize = 256</code>	37
2.29.6	Calcolo delle Dimensioni (Caso 2D)	37
2.29.7	Calcolo delle Dimensioni (Caso Generale 3D)	37

2.30	Analisi delle Prestazioni	38
2.30.1	Verifica del Kernel CUDA (Somma di Array)	38
2.30.2	Gestione degli Errori in CUDA	38
2.30.3	Profiling delle Prestazioni in CUDA	39
2.30.4	Metodi Principali	39
2.30.4.1	Timer CPU	39
APPENDICES	44

CHAPTER 1

Introduzione

1.1 Perché Scegliere la Piattaforma CUDA?

Dominio di Mercato e Standard Industriale

Standard de facto per calcolo parallelo e GPU, ampiamente supportata da software e librerie come TensorFlow e PyTorch. Le GPU NVIDIA sono prevalenti in HPC, AI e simulazione scientifica.

Prestazioni Elevate

CUDA consente di sfruttare la potenza delle GPU NVIDIA per eseguire milioni di thread simultaneamente, migliorando significativamente le prestazioni rispetto ai processori CPU tradizionali.

Ampia Documentazione e Risorse

NVIDIA fornisce una documentazione dettagliata e risorse pratiche, mentre la comunità attiva consente il supporto e la condivisione di conoscenze tra sviluppatori.

Facilità d'Uso

CUDA estende i linguaggi di programmazione C, C++, e Fortran, permettendo agli sviluppatori di utilizzare sintassi e concetti già noti.

Versatilità

È utilizzato in vari campi, dalla grafica 3D alla simulazione scientifica, dall'elaborazione video al deep learning, rendendolo una scelta flessibile per molti progetti.

Ecosistema Ricco

CUDA offre un ampio set di librerie ottimizzate (cuBLAS, cuDNN, Thrust, etc.) e strumenti di sviluppo per facilitare l'ottimizzazione e il debugging.

1.2 Cos'è il CUDA Toolkit?

- Il CUDA Toolkit è un insieme completo di strumenti di sviluppo fornito da NVIDIA per creare applicazioni accelerate tramite GPU.
- È essenziale per lo sviluppo di applicazioni CUDA, poiché fornisce tutti gli strumenti necessari per scrivere, compilare e ottimizzare codice CUDA.

1.2.1 Componenti chiave del CUDA Toolkit

Driver NVIDIA

- Fondamento Invisibile: Essenziali per CUDA, ma gli sviluppatori raramente interagiscono direttamente con essi.
- Ruolo: Funzionano da ponte tra il sistema operativo e la GPU, gestendo l'hardware, il caricamento del codice e il trasferimento dati tra CPU e GPU.
- Installazione: Necessari per CUDA e di solito installati separatamente.
- Compatibilità: Devono essere compatibili sia con la versione del CUDA Toolkit utilizzata che con la GPU in uso.

CUDA Runtime / CUDA Driver API (Application Programming Interface)

- Gli sviluppatori possono scegliere tra due interfacce per interagire con la GPU:
 - CUDA Runtime API: Livello di astrazione più alto, più semplice da utilizzare.
 - CUDA Driver API: Livello di astrazione più basso, offre un controllo più granulare sulle operazioni.

Compilatore CUDA (NVIDIA CUDA Compiler - `nvcc`)

- Traduzione del codice: Compila il codice CUDA, scritto in linguaggi come C, C++ e Fortran, in un formato eseguibile dalle GPU NVIDIA.
- Fasi:
 - Separazione: Separa il codice destinato alla CPU da quello per la GPU.
 - Compilazione: Compila il codice GPU in PTX (linguaggio intermedio) o direttamente in codice macchina per l'architettura GPU target (tenendo conto della Compute Capability e della CUDA Version utilizzata).
 - Linking: Combina il codice CPU e GPU con le librerie CUDA per creare l'applicazione finale.

Librerie CUDA

- Le librerie forniscono implementazioni ottimizzate e parallele di operazioni comuni, così gli sviluppatori non devono reinventare algoritmi complessi da zero.
- Esempi:
 - cuBLAS: Algebra lineare (operazioni su matrici e vettori).
 - cuFFT: Fast Fourier Transform (analisi di segnali, elaborazione di immagini).
 - cuDNN: Primitive per deep neural networks (convoluzione, pooling, attivazione).
 - cuRAND: Generazione di numeri casuali (simulazioni Monte Carlo, crittografia).
 - cuSPARSE: Operazioni su matrici sparse (risoluzione di sistemi lineari sparsi).
 - Thrust: Algoritmi paralleli generici (ordinamento, ricerca, trasformazioni) su GPU.

Esempi di Codice (CUDA Samples)

- Utilità: Forniscono implementazioni concrete di algoritmi e applicazioni comuni che utilizzano CUDA.
- Scopo: Aiutano gli sviluppatori a imparare le best practice di programmazione CUDA e a iniziare rapidamente nuovi progetti.

1.2.1.1 Strumenti di Debugging e Profiling

- Nsight Systems:
 - Profiling a livello di sistema. Offre una visione d'insieme del comportamento dell'applicazione su CPU e GPU, evidenziando eventuali colli di bottiglia.
- Nsight Compute:
 - Profiling approfondito della GPU. Permette di analizzare le prestazioni dei kernel CUDA in dettaglio, identificando aree di ottimizzazione per la memoria e l'utilizzo dei core.
- CUDA-GDB:
 - Debugging a riga di comando. Permette di eseguire il debug del codice CUDA a livello di sorgente, sia sulla CPU che sulla GPU. Supporta breakpoint, ispezione di variabili e stack trace su thread GPU.
- NVIDIA Visual Profiler (NVVP) [non più supportato]:
 - Offre una rappresentazione grafica timeline delle attività della CPU e della GPU, facilitando l'identificazione dei colli di bottiglia. Oggi le sue funzionalità sono state integrate in Nsight.

1.2.1.2 Relazione tra CUDA Toolkit e CUDA Version

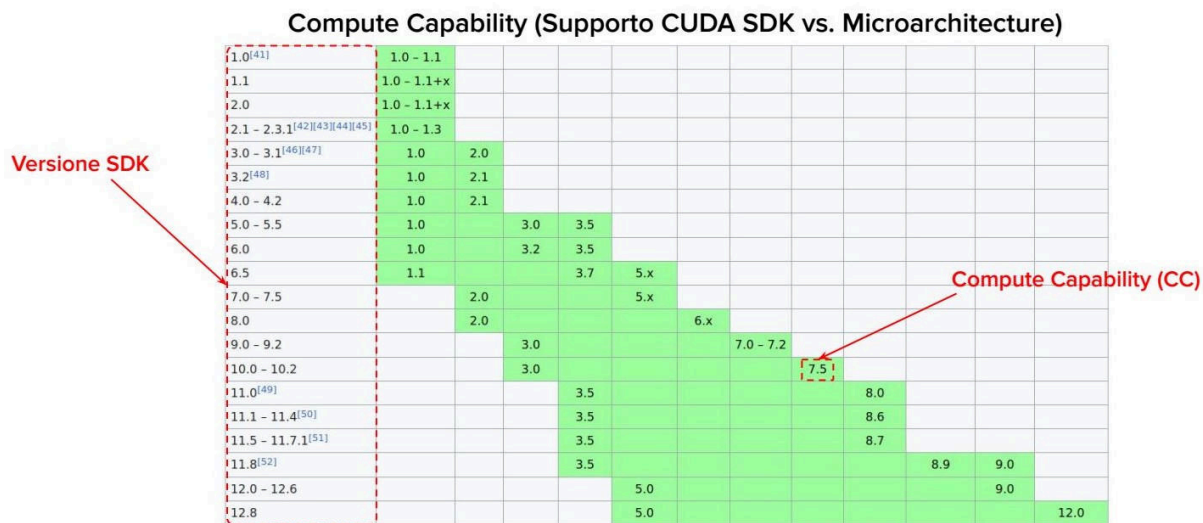
- Il CUDA Toolkit viene rilasciato in versioni numerate (es. CUDA 12.6, CUDA 11.0). La versione installata determina la CUDA Version in uso sul sistema.
- Ogni nuova CUDA Version include aggiornamenti software e bugfix, nuovi strumenti di sviluppo e librerie, supporto per le più recenti architetture GPU NVIDIA.
- Aggiornare il CUDA Toolkit alla versione più recente garantisce compatibilità con le GPU più recenti e l'accesso alle ultime ottimizzazioni, migliorando performance e stabilità del codice.

1.2.1.3 Retrocompatibilità del CUDA Toolkit

- **Supporto per GPU Precedenti:** Le nuove versioni del CUDA Toolkit mantengono la compatibilità con GPU più vecchie, anche se non tutte le funzionalità più recenti sono disponibili su queste architetture.
- **Limitazioni:** Alcune funzionalità avanzate introdotte nelle nuove versioni potrebbero non essere supportate su GPU datate, e il supporto per architetture molto vecchie può essere gradualmente ridotto o deprecato.
- **Compatibilità del Codice:** In generale, il codice scritto con versioni precedenti del Toolkit può essere eseguito su versioni più recenti, ma può richiedere piccoli adattamenti.

1.3 CUDA Compute Capability (CC)

- La Compute Capability (CC) è un numero in formato X.Y che identifica le caratteristiche e le capacità di una GPU NVIDIA in termini di funzionalità supportate e limiti hardware
- Il numero X (principale) indica la generazione dell'architettura (ad es. Turing, Ampere), mentre Y (secondario) indica una revisione della stessa architettura, con piccoli miglioramenti o varianti hardware.



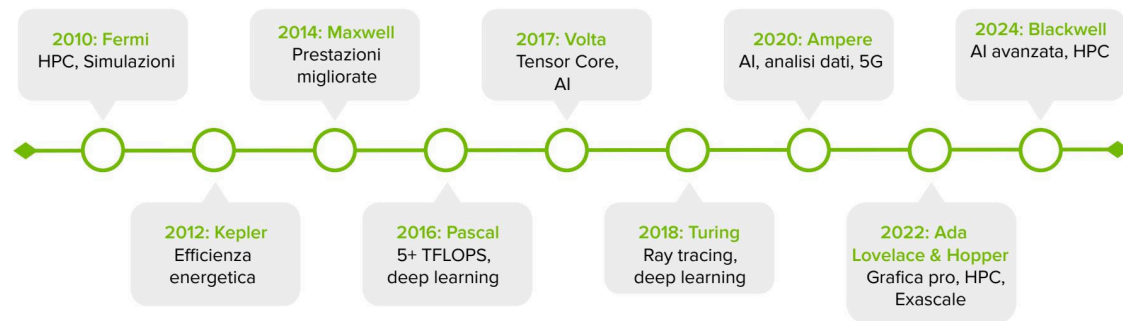
1.3.1 Relazione tra Compute Capability (CC) e CUDA Version

- **Compute Capability (CC)**
 - Indica le caratteristiche e i limiti hardware di una GPU.
 - È indipendente dalla CUDA Version, ma la CUDA Version deve supportare la Compute Capability della GPU per sfruttare appieno le sue capacità.
- **CUDA Version**
 - Determina quali funzionalità software, API e librerie sono disponibili per lo sviluppo.

- Può supportare più Compute Capabilities: una singola versione di CUDA Toolkit può essere compatibile con diverse generazioni di GPU (es. CUDA 11.x supporta Volta, Turing, Ampere).

1.4 Evoluzione delle Architetture GPU NVIDIA

- Progressione Tecnologica: Da Fermi a Blackwell, ogni generazione ha portato significativi avanzamenti nelle capacità di calcolo e nell'efficienza energetica.
- Adattamento al Mercato: L'evoluzione riflette il passaggio da un focus su grafica e HPC a un'enfasi crescente su AI, deep learning e calcolo ad alte prestazioni.



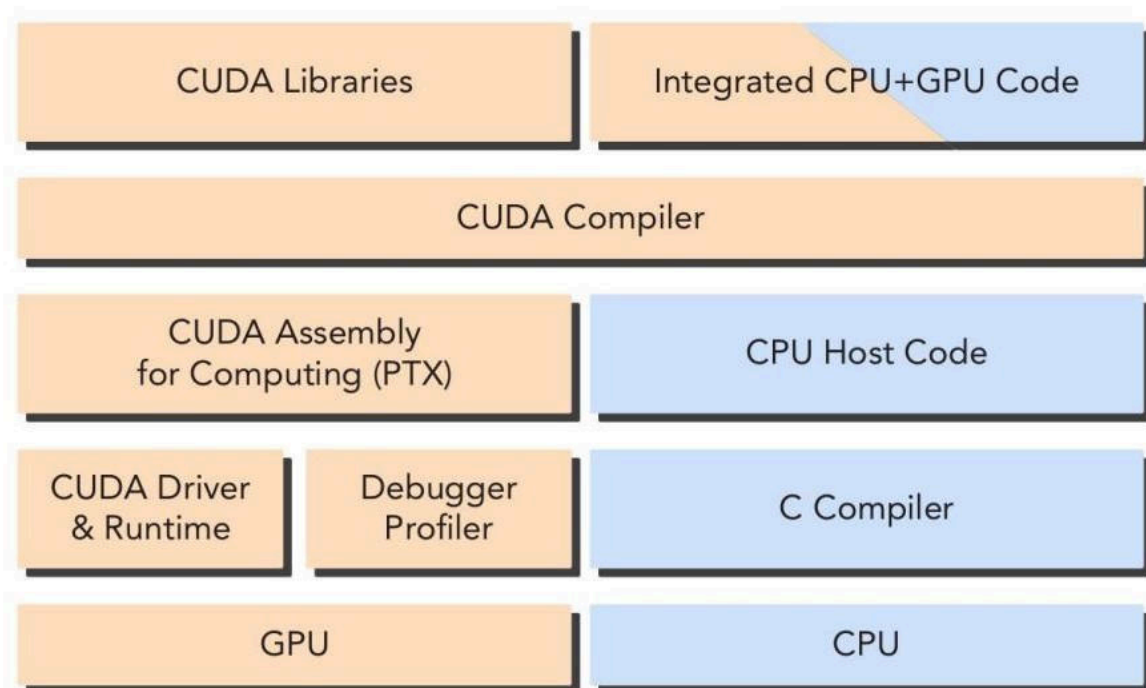
1.5 Anatomia di un Programma CUDA

1.5.1 Struttura del Codice Sorgente

- File Sorgente: Estensione `.cu`
- Codice host + Codice device

1.5.2 Componenti Principali

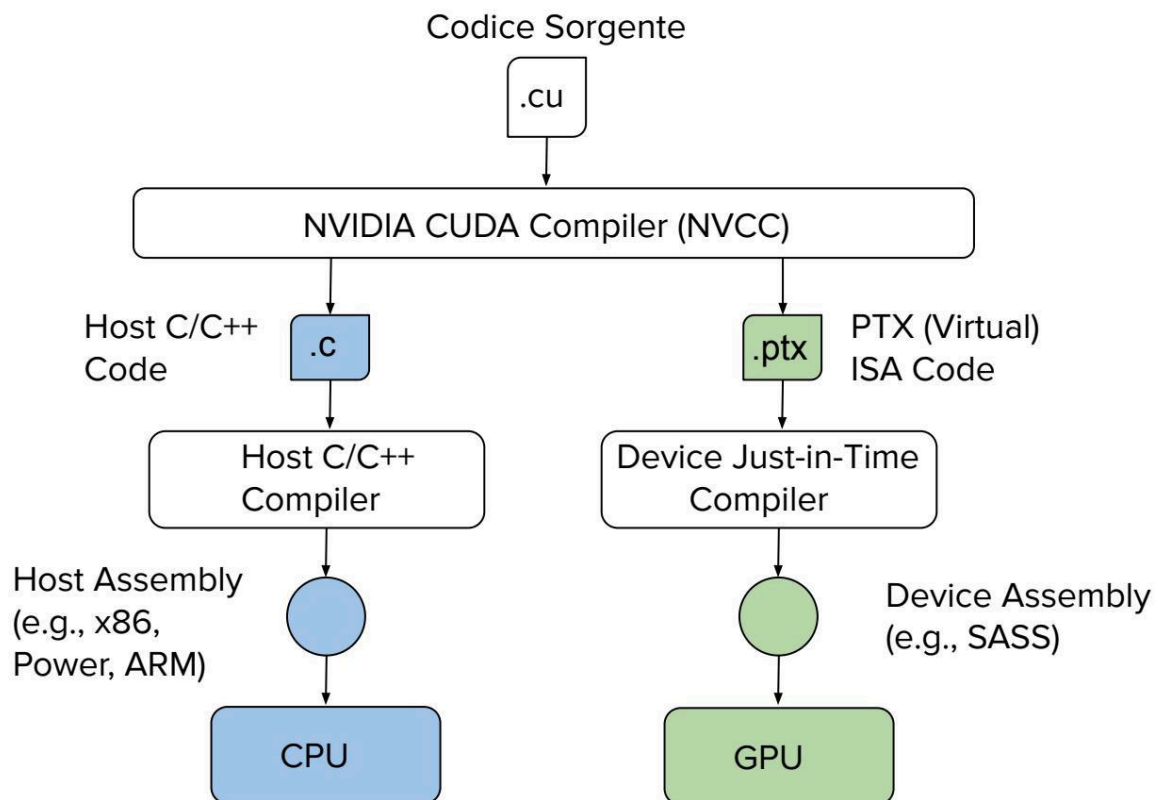
- Codice Host
 - Codice C/C++ eseguito sulla CPU
 - Gestisce la logica dell'applicazione
 - Alloca memoria sulla GPU
 - Trasferisce dati tra CPU e GPU
 - Lancia i kernel GPU
 - Gestisce la sincronizzazione
- Codice Device:
 - Codice CUDA C eseguito sulla GPU
 - Contiene i kernel (funzioni parallele)
 - Esegue operazioni computazionali intensive in parallelo



1.5.3 Flusso di Compilazione

- Separazione del codice
 - Il compilatore NVIDIA CUDA Compiler (nvcc) separa il codice device dal codice host
- Compilazione del codice host
 - È codice C standard o C++
 - Compilato con compilatori C tradizionali (gcc)
- Compilazione del codice device
 - Compilato da nvcc in formato intermedio PTX (Parallel Thread Execution) .
 - Il driver NVIDIA poi traduce il PTX in codice macchina specifico per la GPU (SASS - Streaming Assembly) al momento dell'esecuzione, usando un compilatore Just-In-Time (JIT).
- Linking
 - Aggiunta delle librerie runtime CUDA
 - Supporto per chiamate ai kernel e manipolazione esplicita della GPU
- Eseguibile finale
 - File unico con codice per CPU e GPU

1.6 Compilazione di un Programma CUDA



1.7 Hello World in CUDA C

1.7.1 Passo 1: Creare il File Sorgente

Nome file: hello.cu

1.7.1.1 Passo 2: Scrivere il codice

Codice C e CUDA C a confronto

Linguaggio C C

```

1 include <stdio.h>
2 int main(void) {
3     printf("Hello World from CPU!\n");
4     return 0;
5 }
  
```

Linguaggio CUDA C

CUDA

```

1  include <stdio.h>
2  __global__ void helloFromGPU()
3  {
4      printf("Hello World from GPU thread
5      %d!\n", threadIdx.x);
6  }
7  int main()
8  {
9      // Lancio del kernel
10     helloFromGPU<<<1, 10>>>();
11     // Attendere che la GPU finisca
12     cudaDeviceSynchronize();
13     return 0;
14 }

```

Linguaggio CUDA C

1.8 Hello World in CUDA C - Analisi

- Definizione kernel GPU:
 - `__global__`: Qualificatore CUDA che indica una funzione eseguita sulla GPU, ma chiamata dalla CPU. (In C standard non esiste).
 - `threadIdx.x`: Variabile built-in CUDA che fornisce l'ID univoco del thread all'interno del blocco.

```

1  __global__ void helloFromGPU(){
2      printf("Hello World from GPU thread %d!\n", threadIdx.x);}

```

CUDA

- Funzione main: Punto di ingresso del programma, eseguito sulla CPU come in C standard.
- Lancio del kernel
 - `<<<1, 10>>>`: Configurazione di esecuzione (1 blocco, 10 thread). Avvia 10 istanze parallele del kernel sulla GPU.
- Sincronizzazione GPU-CPU
 - `cudaDeviceSynchronize()`: la CPU attende che la GPU completi tutte le operazioni prima di proseguire/terminare.

```

1  int main(){
2      helloFromGPU<<<1, 10>>>();
3      cudaDeviceSynchronize();
4      return 0;
5  }

```

CUDA

1.9 Hello World in CUDA C

1.9.0.1 Passo 3: Compilazione

- Salvare il codice nel file `hello.cu`
 - Un file `.cu` può contenere sia codice C/C++ standard che codice CUDA C
 - Questo permette di mescolare codice per CPU e GPU nello stesso file
- Compilare il programma usando il compilatore CUDA `nvcc`

- nvcc può gestire sia il codice C/C++ standard che le estensioni CUDA
- nvcc separa internamente il codice host e device, compilando ciascuno in modo appropriato

```
1 $ nvcc hello.cu -o hello
```

\$ Shell

1.10 Hello World in CUDA C

1.10.0.1 Passo 4: Esecuzione

Eseguire il file eseguibile:

```
1 $ ./hello
```

Output:

Variante Linguaggio C

```
1 Hello World from CPU!
```

Variante Linguaggio CUDA C

```
1 Hello World from GPU thread 0!
2 Hello World from GPU thread 1!
3 Hello World from GPU thread 2!
4 Hello World from GPU thread 3!
5 Hello World from GPU thread 4!
6 Hello World from GPU thread 5!
7 Hello World from GPU thread 6!
8 Hello World from GPU thread 7!
9 Hello World from GPU thread 8!
10 Hello World from GPU thread 9!
```

1.11 Ottenere Informazioni sulla GPU tramite API CUDA

1.11.1 Utilizzo delle API CUDA

CUDA fornisce API per ottenere informazioni dettagliate sulle GPU direttamente dal codice

```
1 int main() {
2     cudaDeviceProp prop;
3     cudaGetDeviceProperties(&prop, 0); // Ottieni proprietà del dispositivo 0
4     printf("Nome Dispositivo: %s\n", prop.name);
5     printf("Memoria Globale Totale: %.0f MB\n", prop.totalGlobalMem / 1024.0 /
6         1024.0);
7     printf("Clock Core: %d MHz\n", prop.clockRate / 1000);
8     printf("Compute Capability: %d.%d\n", prop.major, prop.minor);
9     // Stampa di altre proprietà
10    return 0;
11 }
```

CUDA

- Utilizzo della struttura `cudaDeviceProp` per memorizzare le proprietà del device
- Utilizzo della funzione `cudaGetDeviceProperties` per ottenere le proprietà del device specificato

- Accesso alle proprietà della GPU, come nome, memoria totale, clock core e compute capability
- Per una lista completa delle proprietà disponibili, consulta la documentazione ufficiale di CUDA(<https://docs.nvidia.com/cuda/cuda-runtime-api/structcudaDeviceProp.html>)

1.12 Ottenere Informazioni sulla GPU tramite API CUDA

1.12.0.1 Esempio di Output

```

1  CUDA System Information:
2  CUDA Driver Version: 12.2
3  CUDA Runtime Version: 11.3
4  Numero di dispositivi CUDA: 2
5  Dispositivo 0: NVIDIA GeForce RTX 3090
6  1. Compute Capability: 8.6
7  2. Memoria Globale Totale: 23.69 GB
8  3. Numero di Multiprocessori: 82
9  4. Clock Core: 1695 MHz
10 5. Clock Memoria: 9751 MHz
11 6. Larghezza Bus Memoria: 384 bit
12 7. Dimensione Cache L2: 6144 KB
13 8. Memoria Condivisa per Blocco: 48 KB
14 9. Numero Massimo di Thread per Blocco: 1024
15 10. Dimensioni Massime Griglia: (2147483647, 65535, 65535)
16 11. Dimensioni Massime Blocco: (1024, 1024, 64)
17 12. Warp Size: 32
18 13. Memoria Costante Totale: 65536 bytes
19 14. Texture Alignment: 512 bytes

```

1.13 Cosa Significa Programmare in CUDA?

1.13.1 Pensare in Parallelo

- **Decomposizione del Problema:** Identificare le parti del problema che possono essere eseguite in parallelo per sfruttare al meglio le risorse della GPU.
- **Architettura della GPU:** Le GPU sono composte da migliaia di core in grado di eseguire thread in parallelo. CUDA fornisce gli strumenti per organizzare e gestire questi thread.
- **Scalabilità:** Progettare algoritmi che si adattano a diversi numeri di thread (e GPU).
- **Gerarchia di Thread:** Organizzare il lavoro in blocchi e griglie per massimizzare l'efficienza.
- **Gerarchia di Memoria:** Utilizzare strategicamente memoria globale, condivisa, locale e registri per ridurre i tempi di accesso.
- **Sincronizzazione:** Gestire la coordinazione tra thread e il trasferimento dati tra CPU e GPU senza conflitti.
- **Bilanciamento del Carico:** Distribuire il lavoro in modo uniforme fra thread per evitare colli di bottiglia.

1.13.2 Scrittura di codice in CUDA C

- CUDA estende C/C++ con costrutti specifici come `__global__`, `__shared__` e la sintassi `<<<...>>>` per lanciare kernel sulla GPU.

- Ogni kernel viene scritto come codice sequenziale, ma viene eseguito in parallelo da migliaia di thread, permettendo di pensare in modo semplice ma scalare.

1.13.2.1 Testi Generali

- Cheng, J., Grossman, M., McKercher, T. (2014). Professional CUDA C Programming. Wrox Pr Inc. (1^a edizione)
- Kirk, D. B., Hwu, W. W. (2022). Programming Massively Parallel Processors. Morgan Kaufmann (4^a edizione)

1.13.2.2 NVIDIA Docs

- CUDA Programming:
 - <http://docs.nvidia.com/cuda/cuda-c-programming-guide/>
- CUDA C Best Practices Guide
 - <http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/>
- CUDA University Courses
 - <https://developer.nvidia.com/educators/existing-courses=2>
- An Even Easier Introduction to CUDA
 - <https://developer.nvidia.com/blog/even-easier-introduction-cuda/>

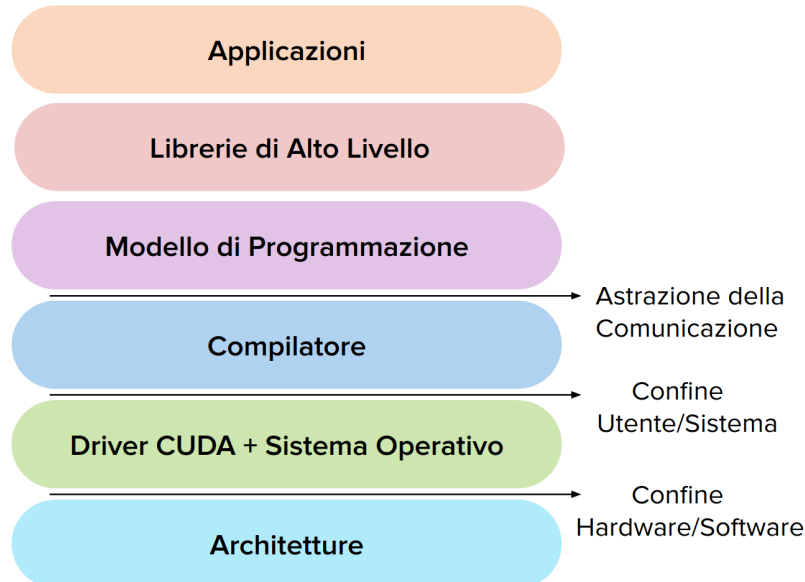
1.13.2.3 Materiale di Approfondimento

- Branch Education (canale YouTube con spiegazioni su GPU, ray tracing, hardware, ecc.)
 - <https://www.youtube.com/@BranchEducation>

CHAPTER 2

Modello di Programmazione CUDA

2.1 La Struttura Stratificata dell'Ecosistema CUDA



Ecosistema stratificato per algoritmi paralleli su GPU, con semplicità e controllo hardware ottimizzati.

Applicazioni: Programmi scritti dagli sviluppatori per risolvere problemi specifici utilizzando CUDA.

Librerie: Raccolte di funzioni ottimizzate (es. cuBLAS, cuDNN) che semplificano lo sviluppo.

Modello di Programmazione: CUDA fornisce un'astrazione per la programmazione GPU, offrendo concetti come thread, blocchi e griglie.

Compilatore: Strumenti (nvcc) che traducono il codice in istruzioni GPU eseguibili.

Driver CUDA + Sistema Operativo: Il sistema operativo gestisce le risorse; il driver CUDA traduce le chiamate CUDA in comandi per la GPU.

Architetture: Le specifiche GPU NVIDIA su cui il codice CUDA viene eseguito, con diverse capacità e caratteristiche.

2.2 Ruolo del Modello e del Programma

Il Modello di Programmazione:

Definisce la struttura e le regole per sviluppare applicazioni parallele su GPU. Elementi fondamentali:

- **Gerarchia di Thread:** Organizza l'esecuzione parallela in thread, blocchi e griglie, ottimizzando la scalabilità su diverse GPU.
- **Gerarchia di Memoria:** Offre tipi di memoria (globale, condivisa, locale, costante, texture) con diverse prestazioni e scopi, per ottimizzare l'accesso ai dati.
- **API:** Fornisce funzioni e librerie per gestire l'esecuzione del kernel, il trasferimento dei dati e altre operazioni essenziali.

Il Programma:

Rappresenta l'implementazione concreta (il codice) che specifica come i thread condividono dati e coordinano le loro attività. Nel programma CUDA, si definisce:

- Come i dati verranno suddivisi e elaborati tra i vari thread.
- Come i thread accederanno alla memoria e condivideranno dati.
- Quali operazioni verranno eseguite in parallelo.
- Quando e come i thread si sincronizzeranno per completare un compito.

2.3 Livelli di Astrazione nella Programmazione Parallela CUDA

Il calcolo parallelo si articola in tre livelli di astrazione: dominio, logico e hardware, guidando l'approccio del programmatore.

2.3.1 Livello Dominio

- Focus sulla decomposizione del problema.
- Definizione della struttura parallela di alto livello.

Chiave: Ottimizza la strategia di parallelizzazione.

2.3.2 Livello Logico

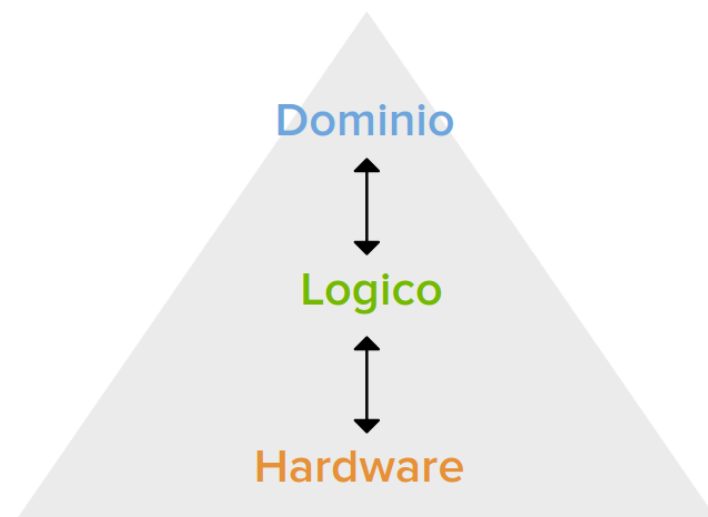
- Organizzazione e gestione dei thread.
- Implementazione della strategia di parallelizzazione.

Chiave: Massimizza l'efficienza del parallelismo.

2.3.2.1 Livello Hardware

- Mappatura dell'esecuzione sull'architettura GPU.
- Ottimizzazione delle prestazioni hardware.

Chiave: Sfrutta al meglio le risorse GPU.



2.4 Thread CUDA: L'Unità Fondamentale di Calcolo

2.4.1 Cos'è un Thread CUDA?

- Un thread CUDA rappresenta un'unità di esecuzione elementare nella GPU.
- Ogni thread CUDA esegue una porzione di un programma parallelo, chiamato kernel.

- Sebbene migliaia di thread vengano eseguiti concorrentemente sulla GPU, ogni singolo thread segue un percorso di esecuzione sequenziale all'interno del suo contesto.

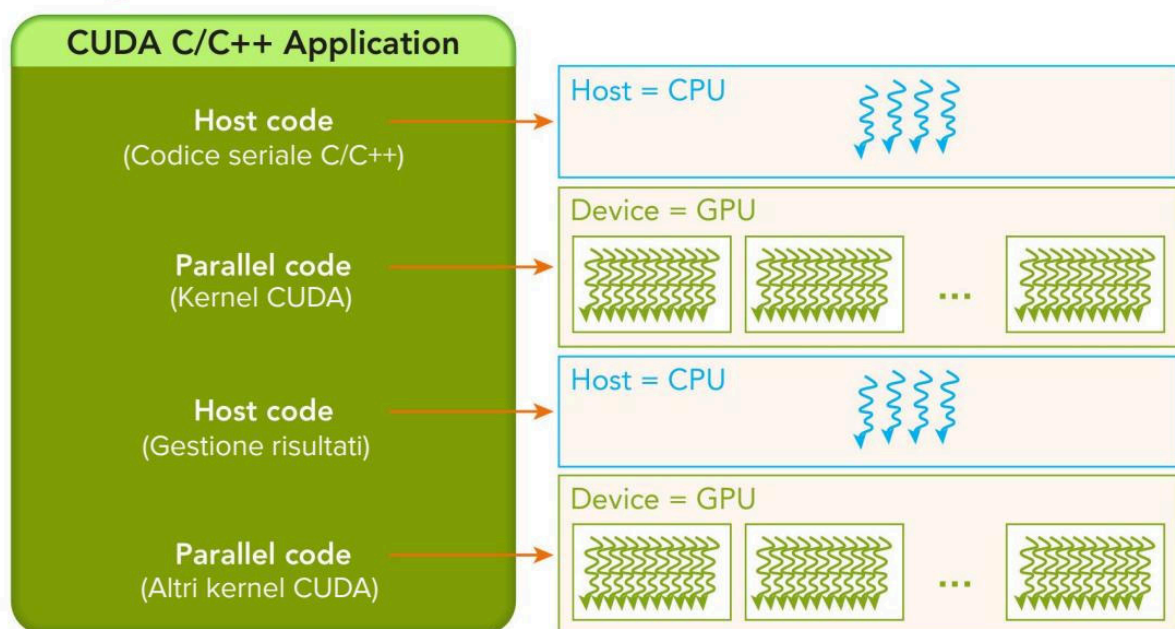
2.4.2 Cosa Fa un Thread CUDA?

- **Elaborazione di Dati:** Ogni thread CUDA si occupa di un piccolo pezzo del problema complessivo, eseguendo calcoli su un sottoinsieme di dati.
- **Esecuzione di Kernel:** Ogni thread esegue lo stesso codice del kernel ma opera su dati diversi, determinati dai suoi identificatori univoci (threadIdx, blockIdx).
- **Stato del Thread:** Ogni thread ha il proprio stato, che include il program counter, i registri, la memoria locale e altre risorse specifiche del thread.

Thread CUDA vs Thread CPU

- **GPU:** parallelismo massivo (migliaia di core leggeri), basso overhead di gestione.
- **CPU:** parallelismo limitato (pochi core complessi), overhead più elevato.

2.5 Struttura di Programmazione CUDA



2.5.1 Caratteristiche Principali

- **Codice Seriale e Parallelo:** Alternanza tra sezioni di codice seriale e parallelo (stesso file).
- **Struttura Ibrida Host-Device:** Alternanza tra codice eseguito sulla CPU (host) e sulla GPU (device).
- **Esecuzione Asincrona:** Il codice host può continuare l'esecuzione mentre i kernel GPU sono in esecuzione.
- **Kernel CUDA Multipli:** Possibilità di lanciare più kernel nella stessa applicazione, anche in overlapping temporale.
- **Gestione dei Risultati sull'Host:** Fase dedicata all'elaborazione dei risultati sulla CPU dopo l'esecuzione dei kernel.

2.6 Flusso Tipico di Elaborazione CUDA

1. Inizializzazione e Allocazione Memoria (Host)

Preparazione dati e allocazione di memoria su CPU (host) e GPU (device).

2. Trasferimento Dati (Host \rightarrow Device)

Copia degli input dalla memoria host alla memoria device.

3. Esecuzione del Kernel (Device)

La GPU esegue calcoli paralleli secondo la configurazione di griglia e blocchi.

4. Recupero Risultati (Device \rightarrow Host)

Copia dell'output dalla memoria device alla memoria host.

5. Post-elaborazione (Host)

Analisi o elaborazione aggiuntiva dei risultati sulla CPU.

6. Liberazione Risorse

Rilascio della memoria allocata su host e device.

Nota: i passi 2-5 possono essere ripetuti più volte o eseguiti in pipeline tramite stream per massimizzare l'overlap tra calcolo e trasferimento dati.

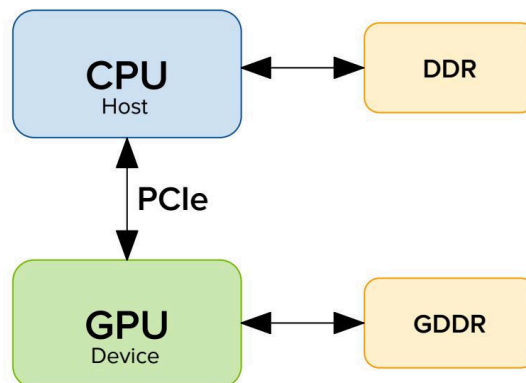
2.7 Gestione della Memoria in CUDA

2.7.1 Modello di Memoria CUDA

- Il modello prevede un sistema con host (CPU) e device (GPU), ciascuno con la propria memoria.
- La comunicazione tra memoria host e device avviene tramite PCIe (Peripheral Component Interconnect Express), interfaccia seriale point-to-point che sfrutta più lane indipendenti in parallelo per aumentare la banda.

2.7.2 Caratteristiche PCIe

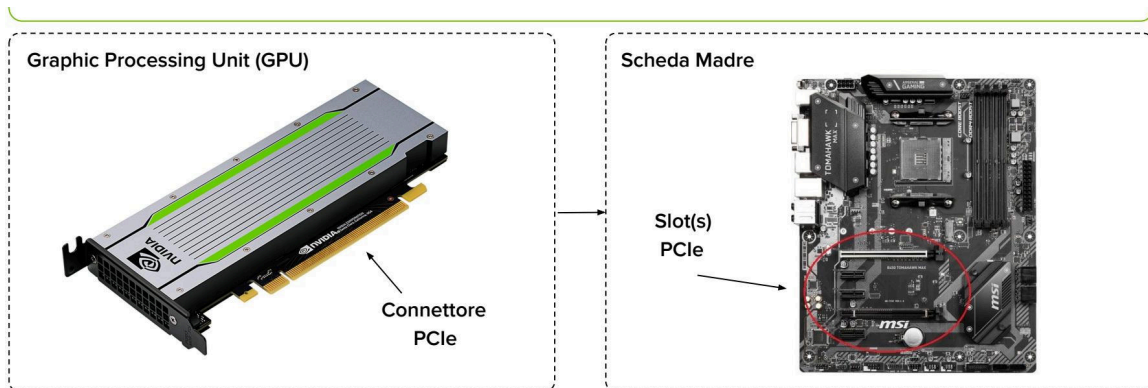
- **Lane:** Ogni lane (canale di trasmissione) è costituito da due coppie di segnali differenziali (quattro fili), una per ricevere (RX) e una per trasmettere (TX) dati.
- **Full-Duplex:** Trasmette e riceve dati simultaneamente in entrambe le direzioni.
- **Scalabilità:** La larghezza di banda varia a seconda del numero di lane (x1, x2, x4, x8, x16).
- **Bassa Latenza:** Garantisce trasferimenti rapidi, adatti a scambi frequenti.
- **Collo di Bottiglia:** Può diventare un collo di bottiglia in trasferimenti di grandi volumi tra CPU e GPU.



2.8 Collegamento Fisico della GPU tramite PCIe

2.8.1 Connessione Fisica GPU

- La GPU si collega alla scheda madre attraverso uno slot PCI Express (PCIe).
- Il connettore, costituito da contatti metallici dorati sul bordo della scheda, si inserisce nello slot PCIe corrispondente.
- La maggior parte delle schede madri moderne ha uno o più slot PCIe, generalmente con almeno uno slot PCIe x16 destinato alla GPU.



2.9 Modello di Memoria CUDA

- I kernel CUDA operano sulla memoria del device.
- CUDA Runtime fornisce funzioni per:
 - Allocare memoria sul device.
 - Rilasciare memoria sul device quando non più necessaria.
 - Trasferire dati bidirezionalmente tra la memoria dell'host e quella del device.

Standard C	CUDA C	Funzione
malloc	cudaMalloc	Alloca memoria dinamica
memcpy	cudaMemcpy	Copia dati tra aree di memoria
memset	cudaMemset	Inizializza memoria a un valore specifico
free	cudaFree	Libera memoria allocata dinamicamente

Nota Importante: è responsabilità del programmatore gestire correttamente l'allocazione, il trasferimento e la deallocazione della memoria per ottimizzare le prestazioni.

2.10 Gerarchia di Memoria

In CUDA, esistono diversi tipi di memoria, ciascuno con caratteristiche specifiche in termini di accesso, velocità, e visibilità. Per ora, ci concentriamo su due delle più importanti:

Global Memory

- Accessibile da tutti i thread su tutti i blocchi
- Più grande ma più lenta rispetto alla shared memory
- Persiste per tutta la durata del programma CUDA
- È adatta per memorizzare dati grandi e persistenti

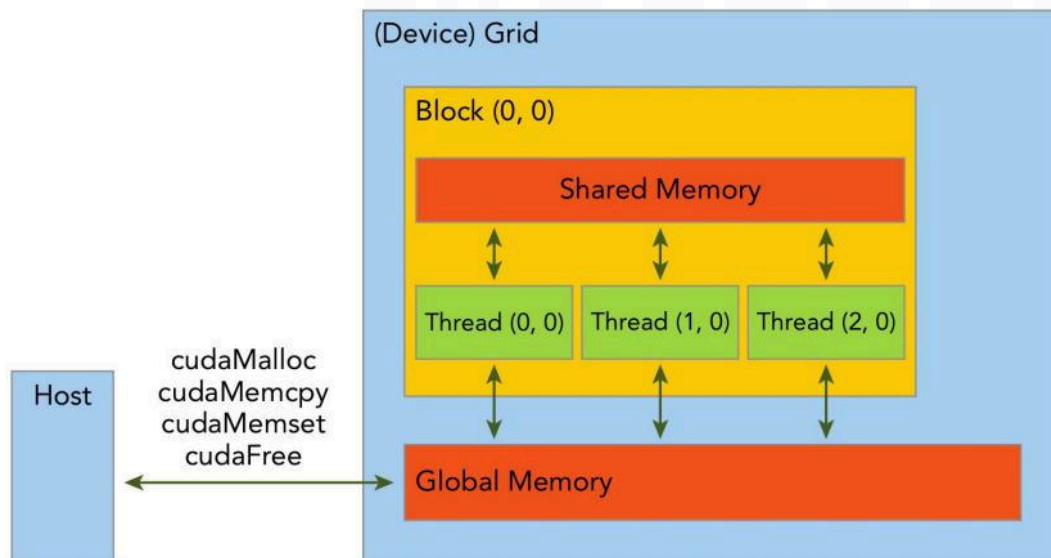
Shared Memory

- Condivisa tra i thread all'interno di un singolo blocco
- Più veloce, ma limitata in dimensioni
- Esiste solo per la durata del blocco di thread
- Utilizzata per dati temporanei e intermedi

Funzioni

- `cudaMalloc`: Alloca memoria sulla GPU.
- `cudaMemcpy`: Trasferisce dati tra host e device.
- `cudaMemset`: Inizializza la memoria del device.
- `cudaFree`: Libera la memoria allocata sul device.

Nota: Queste funzioni operano principalmente sulla Global Memory.



2.11 Allocazione della Memoria sul Device

`cudaMalloc` è una funzione CUDA utilizzata per allocare memoria sulla GPU (device).

Firma della Funzione ([Documentazione Online](#))

```
1 cudaError_t cudaMalloc(void devPtr, size_t size)
```

CUDA

Parametri

- `devPtr`: Puntatore doppio che conterrà l'indirizzo della memoria allocata sulla GPU.
- `size`: Dimensione in byte della memoria da allocare.

Valore di Ritorno

`cudaError_t`: codice di errore (`cudaSuccess` se l'allocazione ha successo).

Note Importanti

- **Allocazione:** Riserva memoria lineare contigua sulla GPU a runtime.
- **Puntatore:** Aggiorna puntatore CPU con indirizzo memoria GPU.
- **Stato iniziale:** La memoria allocata non è inizializzata.

`cudaMemset` è una funzione CUDA utilizzata per impostare un valore specifico in un blocco di memoria allocato sulla GPU (device).

Firma della Funzione ([Documentazione Online](#))

```
1 cudaError_t cudaMemset(void devPtr, int value, size_t count)
```

CUDA

Parametri

- **devPtr:** Puntatore alla memoria allocata sulla GPU.
- **value:** Valore da impostare in ogni byte della memoria.
- **count:** Numero di byte della memoria da impostare al valore specificato.

Valore di Ritorno

cudaError_t: Codice di errore (cudaSuccess se l'inizializzazione ha successo).

Note Importanti

- **Utilizzo:** Comunemente utilizzata per azzerare la memoria (impostando **value** a 0).
- **Gestione:** L'inizializzazione deve avvenire dopo l'allocazione della memoria tramite `cudaMalloc`.
- **Efficienza:** È preferibile usare `cudaMemset` per grandi blocchi di memoria per ridurre l'overhead.

2.11.1 Esempio di Allocazione di Memoria sulla GPU

Mostra come allocare memoria sulla GPU utilizzando `cudaMalloc`.

```
1 float d_array; // Dichiarazione di un puntatore per la memoria sul device
   (GPU)
2 size_t size = 10 * sizeof(float); // Calcola la dimensione della memoria da
   allocare (10 float)
3 // Allocazione della memoria sul device
4 cudaError_t err = cudaMalloc((void*)&d_array, size);
5 // Controlla se l'allocazione della memoria ha avuto successo
6 if (err != cudaSuccess) {
7     // Se c'è un errore, stampa un messaggio di errore con la descrizione
   dell'errore
8     printf("Errore nell'allocazione della memoria: %s\n",
   cudaGetErrorString(err));
9 } else {
10    // Se l'allocazione ha successo, stampa un messaggio di conferma
11    printf("Memoria allocata con successo sulla GPU.\n");}
```

CUDA

2.12 Trasferimento Dati

`cudaMemcpy` è una funzione CUDA per il trasferimento di dati tra la memoria dell'host e del device, o all'interno dello stesso tipo di memoria.

Firma della Funzione ([Documentazione Online](#))

```
1 cudaError_t cudaMemcpy(void dst, const void src, size_t count,
   cudaMemcpyKind kind)
```

CUDA

Parametri

- **dst:** Puntatore alla memoria di destinazione.
- **src:** Puntatore alla memoria sorgente.

- `count`: Numero di byte da copiare.
- `kind`: Direzione della copia (`cudaMemcpyKind`).

Tipi di Trasferimento (`kind`)

- `cudaMemcpyHostToHost`: Da host a host
- `cudaMemcpyHostToDevice`: Da host a device
- `cudaMemcpyDeviceToHost`: Da device a host
- `cudaMemcpyDeviceToDevice`: Da device a device

Valore di Ritorno

`cudaError_t`: Codice di errore (`cudaSuccess` se il trasferimento ha successo).

Note importanti

- Funzione sincrona: blocca l'host fino al completamento del trasferimento.
- Per prestazioni ottimali, minimizzare i trasferimenti tra host e device.

2.12.1 Spazi di Memoria Differenti

Attenzione: I puntatori del device non devono essere dereferenziati nel codice host (spazi di memoria CPU e GPU differenti).

Esempio: assegnazione errata come

```
1 host_array = dev_ptr
```

CUDA

invece di

```
1 cudaMemcpy(host_array, dev_ptr, nBytes, cudaMemcpyDeviceToHost)
```

CUDA

Conseguenza dell'errore: Accesso a indirizzi non validi → possibile blocco o crash dell'applicazione.

Soluzione: La Unified Memory, introdotta in CUDA 6 e oggi ottimizzata, consente di usare un unico puntatore valido per CPU e GPU, con gestione automatica della migrazione dati (vedremo in seguito).

2.13 Deallocazione della Memoria sul Device

`cudaFree` è una funzione CUDA utilizzata per liberare la memoria precedentemente allocata sulla GPU (device).

Firma della Funzione ([Documentazione Online](#))

```
1 cudaError_t cudaFree(void devPtr)
```

CUDA

Parametri

`devPtr`: Puntatore alla memoria sul device che deve essere liberata. Questo puntatore deve essere stato precedentemente restituito tramite la chiamata `cudaMalloc`.

Valore di Ritorno

`cudaError_t`: Codice di errore (`cudaSuccess` se la deallocazione ha successo).

Note Importanti

- **Gestione:** È responsabilità del programmatore assicurarsi che ogni blocco di memoria allocato con `cudaMalloc` sia liberato per evitare perdite di memoria (memory leaks) sulla GPU.
- **Efficienza:** La deallocazione è sincrona e può avere overhead significativo; è consigliato minimizzare il numero di chiamate.

2.13.1 Esempio di Allocazione e Trasferimento Dati

Mostra come allocare e trasferire dati dalla memoria host alla memoria device.

```

1  size_t size = 10 * sizeof(float); // Calcola la dimensione della memoria da
    allocare (10 float)
2  float h_data = (float)malloc(size); // Alloca memoria sull'host (CPU) per
    memorizzare i dati
3  for (int i = 0; i < 10; ++i) h_data[i] = (float)i; // Inizializza ogni elemento
    di h_data
4  float d_data; // Dichiarazione di un puntatore per la memoria sulla GPU (device)
5  cudaMalloc((void*)&d_data, size); // Allocazione della memoria sulla GPU
6  // Copia dei dati dalla memoria dell'host (CPU) alla memoria del device (GPU)
7  cudaError_t err = cudaMemcpy(d_data, h_data, size, cudaMemcpyHostToDevice);
8  // Controlla se la copia è avvenuta con successo
9  if (err != cudaSuccess) {
10     // Se c'è un errore, stampa un messaggio di errore e termina il programma
11     fprintf(stderr, "Errore nella copia H2D: %s\n", cudaGetErrorString(err));
12     exit(EXIT_FAILURE);
13 }
14
15 // Esegui operazioni sulla memoria della GPU (d_data)
16 // (Le operazioni specifiche da eseguire non sono mostrate in questo esempio)
17 // Copia dei risultati dalla memoria della GPU (device) alla memoria dell'host
    (CPU)
18 err = cudaMemcpy(h_data, d_data, size, cudaMemcpyDeviceToHost);
19 // Controlla se la copia è avvenuta con successo
20 if (err != cudaSuccess) {
21     fprintf(stderr, "Errore nella copia D2H: %s\n", cudaGetErrorString(err));
22     exit(EXIT_FAILURE);
23 }
24 free(h_data); // Libera la memoria allocata sull'host
25 cudaFree(d_data); // Libera la memoria allocata sulla GPU

```

2.14 Organizzazione dei Thread in CUDA

CUDA adotta una gerarchia a due livelli per organizzare i thread basata su blocchi di thread e griglie di blocchi.

2.14.1 Struttura Gerarchica

1. Grid (Griglia)

- Array di thread blocks.
- È organizzata in una struttura 1D, 2D o 3D.
- Rappresenta l'intera computazione di un kernel.

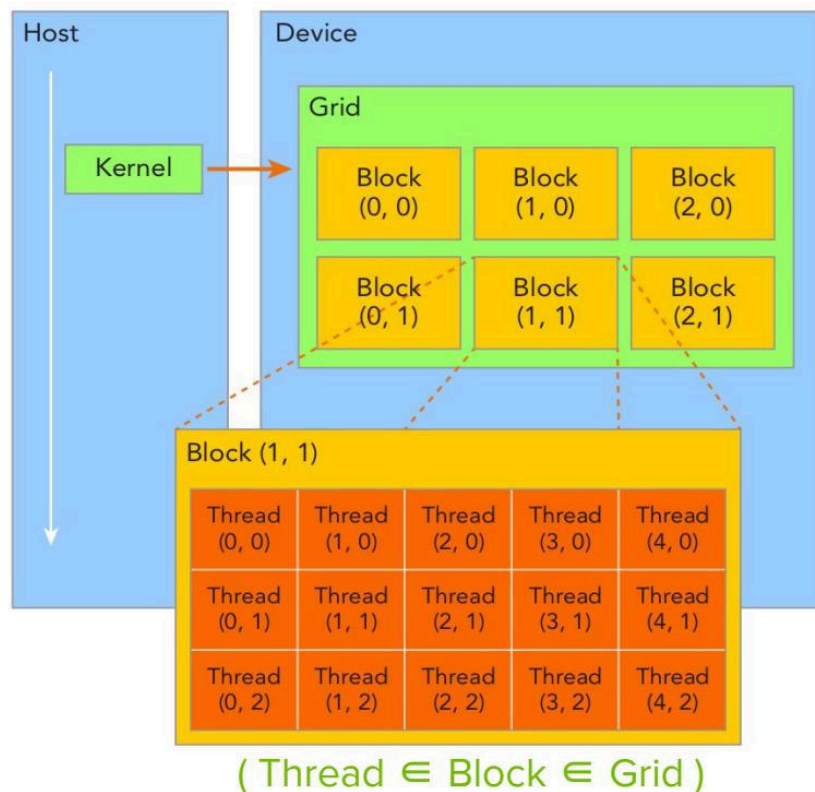
- Contiene tutti i thread che eseguono il singolo kernel.
- Condivide lo stesso spazio di memoria globale.

2. Block (Blocco)

- Un thread block è un gruppo di thread eseguiti logicamente in parallelo.
- Ha un ID univoco all'interno della sua griglia.
- I blocchi sono organizzati in una struttura 1D, 2D o 3D.
- I thread di un blocco possono sincronizzarsi (non automaticamente) e condividere memoria.
- I thread di blocchi diversi non possono sincronizzarsi direttamente (solo tramite memoria globale o kernel successivi)

3. Thread

- Ha un proprio ID univoco all'interno del suo blocco.
- Ha accesso alla propria memoria privata (registri).



2.14.2 Perché una Gerarchia di Thread?

Mappatura Intuitiva

La gerarchia di thread (grid, blocchi, thread) permette di scomporre problemi complessi in unità di lavoro parallele più piccole e gestibili, rispecchiando spesso la struttura intrinseca del problema stesso.

Organizzazione e Ottimizzazione

Il programmatore può definire le dimensioni dei blocchi e della griglia per adattare l'esecuzione alle caratteristiche specifiche dell'hardware e del problema, ottimizzando l'utilizzo delle risorse.

Efficienza nella Memoria

I thread in un blocco possono condividere dati tramite memoria on-chip veloce (es. shared memory), riducendo gli accessi alla memoria globale più lenta, migliorando dunque significativamente le prestazioni.

Scalabilità e Portabilità

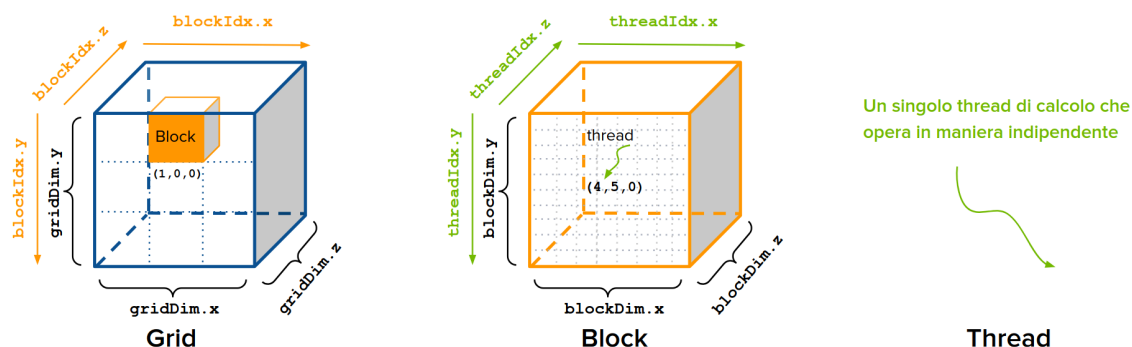
La gerarchia è scalabile e permette di adattare l'esecuzione a GPU con diverse capacità e numero di core. Il codice CUDA, quindi, risulta più portabile e può essere eseguito su diverse architetture GPU.

Sincronizzazione Granulare

I thread possono essere sincronizzati solo all'interno del proprio blocco, evitando costose sincronizzazioni globali che possono creare colli di bottiglia.

2.15 Identificazione dei Thread in CUDA

Ogni thread ha un'identità unica definita da coordinate specifiche nella gerarchia grid-bloc. Tali coordinate, diverse per ogni thread, sono essenziali per calcolare indici di lavoro e accedere correttamente ai dati.



Variabili di Identificazione (Coordinate)

1. `blockIdx` (indice del blocco all'interno della griglia)
 - Componenti: `blockIdx.x`, `blockIdx.y`, `blockIdx.z`
2. `threadIdx` (indice del thread all'interno del blocco)
 - Componenti: `threadIdx.x`, `threadIdx.y`, `threadIdx.z`

Entrambe sono variabili built-in di tipo `uint3` pre-inizializzate dal CUDA Runtime e accessibili solo all'interno del kernel.

Variabili di Dimensioni

1. `blockDim` (dimensione del blocco in termini di thread)
 - Tipo: `dim3` (lato host), `uint3` (lato device, built-in)
 - Componenti: `blockDim.x`, `blockDim.y`, `blockDim.z`
2. `gridDim` (dimensione della griglia in termini di blocchi)
 - Tipo: `dim3` (lato host), `uint3` (lato device, built-in)
 - Componenti: `gridDim.x`, `gridDim.y`, `gridDim.z`

`uint3` è un built-in vector type di CUDA con tre campi (x,y,z) ognuno di tipo unsigned int

2.15.1 Dimensione delle Griglie e dei Blocchi

- La scelta delle dimensioni ottimali dipende dalla struttura dati del problema e dalle capacità hardware/risorse della GPU.

- Le dimensioni di griglia e blocchi vengono definite nel codice host prima del lancio del kernel.
- Sia le griglie che i blocchi utilizzano il tipo `dim3` (lato host) con tre campi `unsigned int`. I campi non utilizzati vengono inizializzati a 1 e ignorati.
- 9 possibili configurazioni (1D, 2D, 3D per griglia e blocco) in tutto anche se in genere si usa la stessa per entrambi.

2.16 Struttura `dim3`

Definizione

- `dim3` è una struttura definita in `vector_types.h` usata per specificare le dimensioni di griglia e blocchi.
- Supporta le dimensioni 1, 2 e 3:

Esempi

CUDA

```
1 dim3 gridDim(256); // Definisce una griglia di 256x1x1 blocchi.
2 dim3 blockDim(512, 512);` // Definisce un blocco di 512x512x1 threads.
```

Utilizzato per specificare le dimensioni di griglia e blocchi quando si lancia un kernel dal lato host:

```
1 kernel_name<<<gridDim, blockDim>>>(...);
```

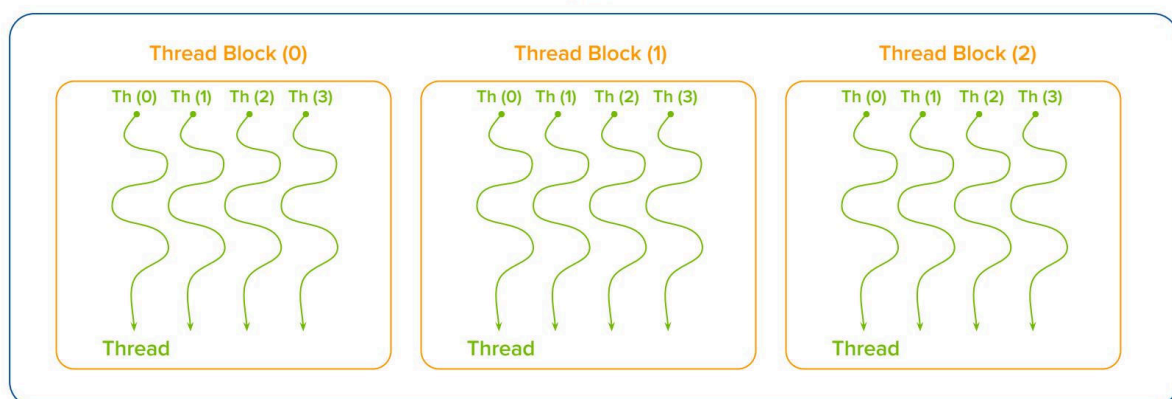
CUDA

Codice Originale:

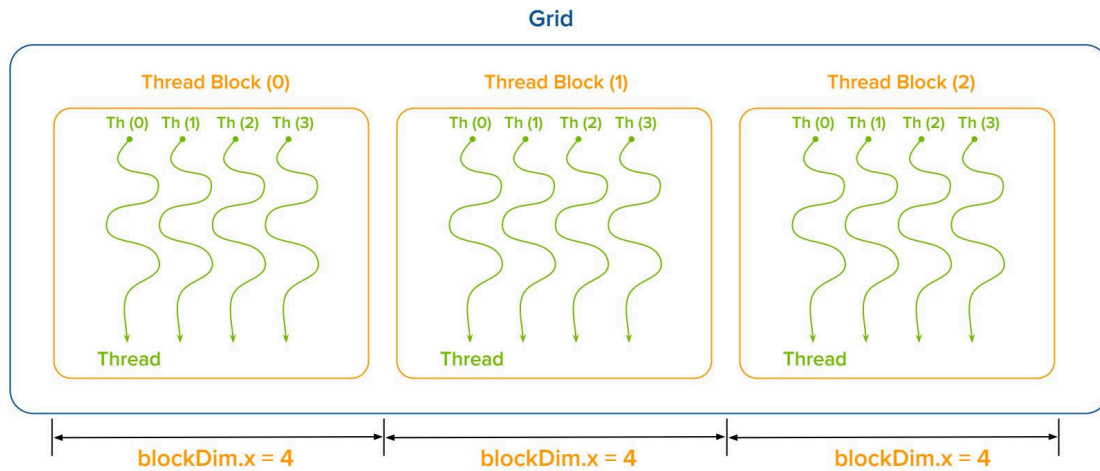
```
1 struct __device_builtin__ dim3
2 {
3     unsigned int x, y, z;
4     #if defined(__cplusplus)
5     __host__ __device__ dim3(unsigned int vx = 1, unsigned int vy = 1, unsigned int
6     vz = 1) : x(vx), y(vy), z(vz) {}
7     __host__ __device__ dim3(uint3 v) : x(v.x), y(v.y), z(v.z) {}
8     __host__ __device__ operator uint3(void) { uint3 t; t.x = x; t.y = y; t.z = z;
9     return t; }
10    #endif / __cplusplus /
11 };
```

CUDA

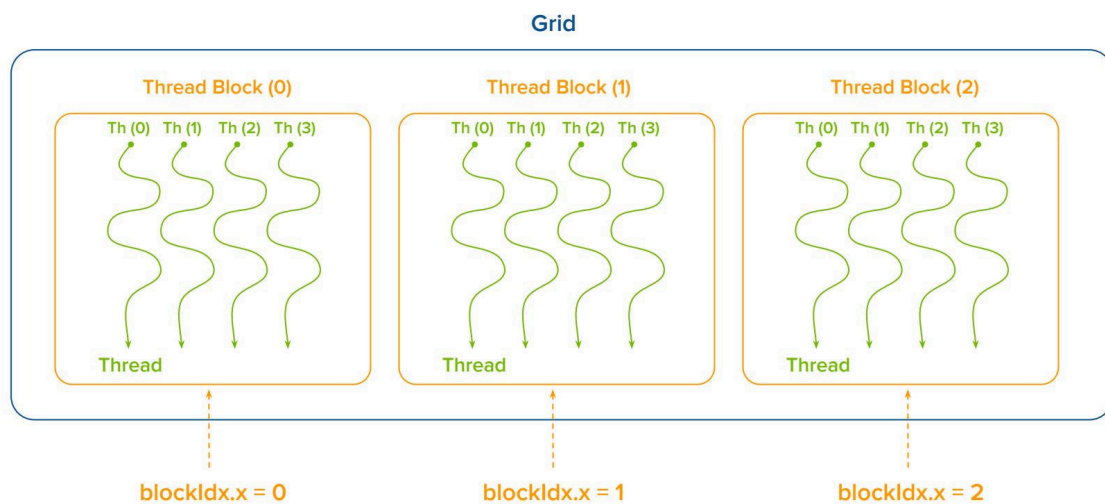
`gridDim.x`: Numero di blocchi nella griglia, in questo caso 3.



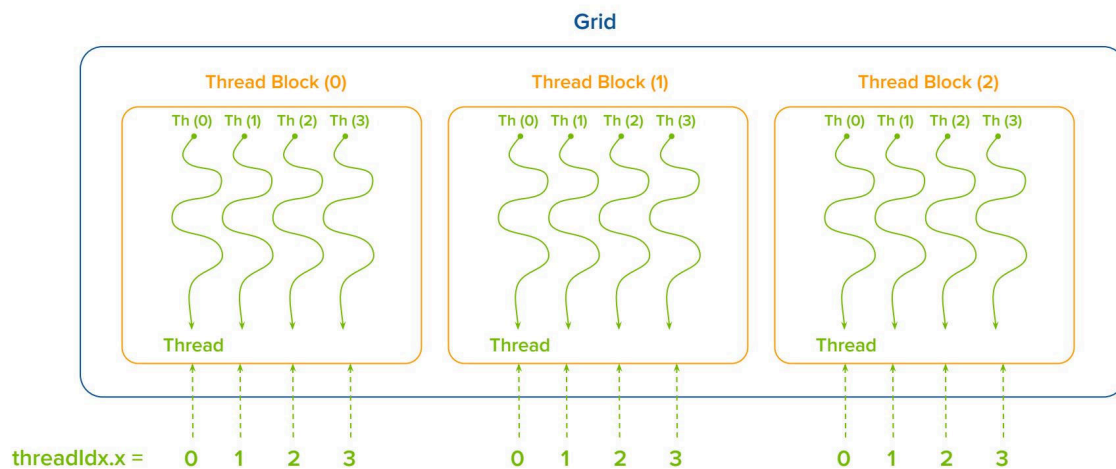
`blockDim.x`: Numero di thread per blocco, in questo caso 4.



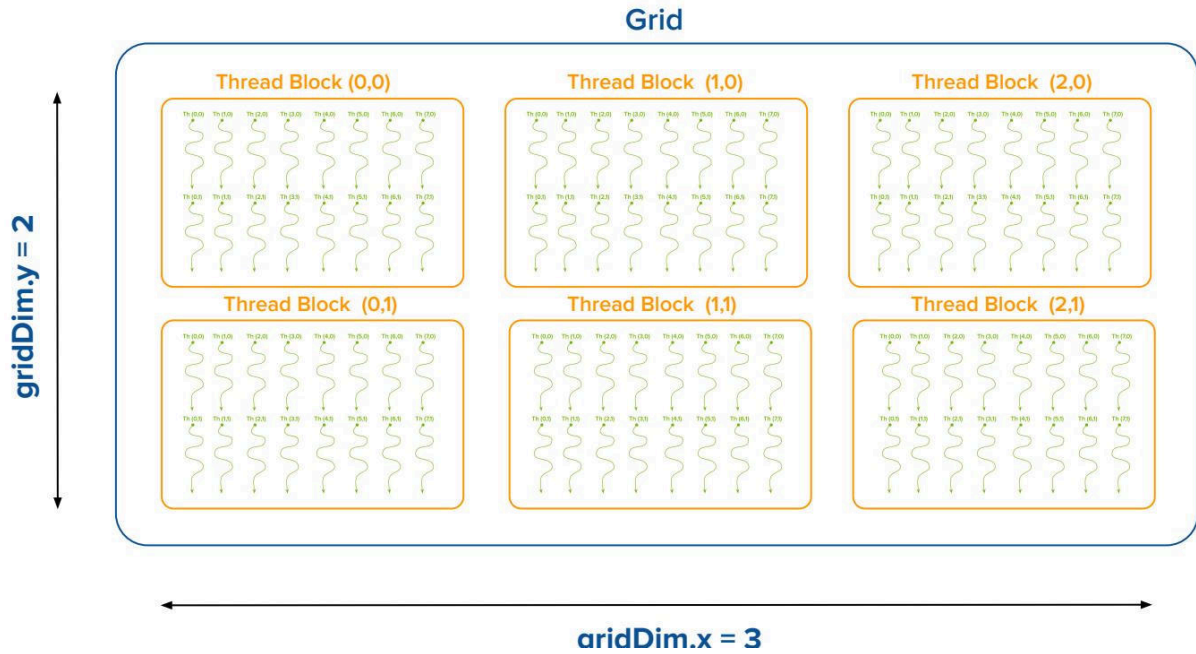
`blockIdx.x`: Indice di un blocco nella griglia.



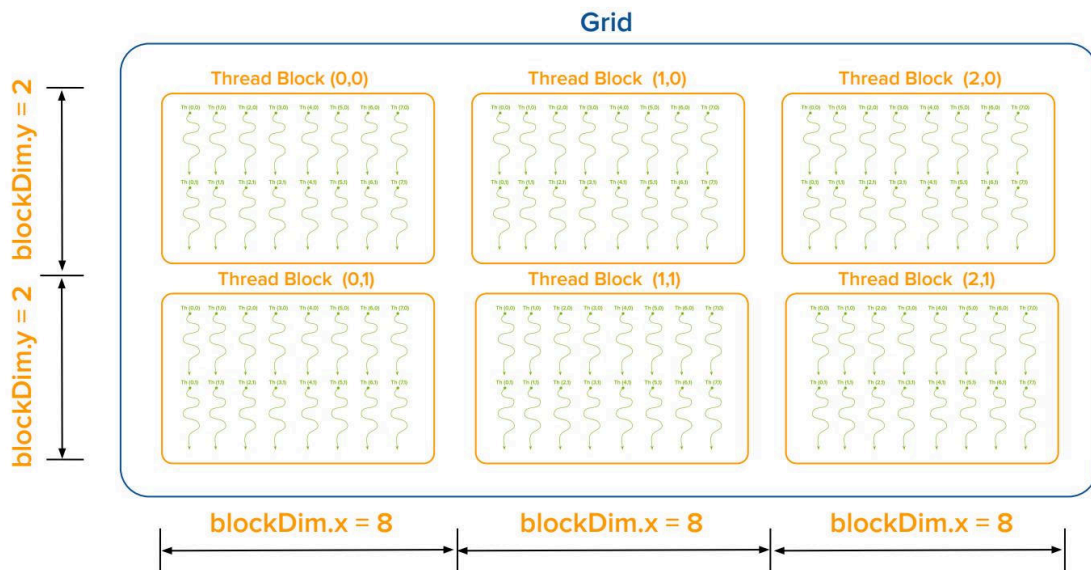
`threadIdx.x`: Indice del thread all'interno del blocco.



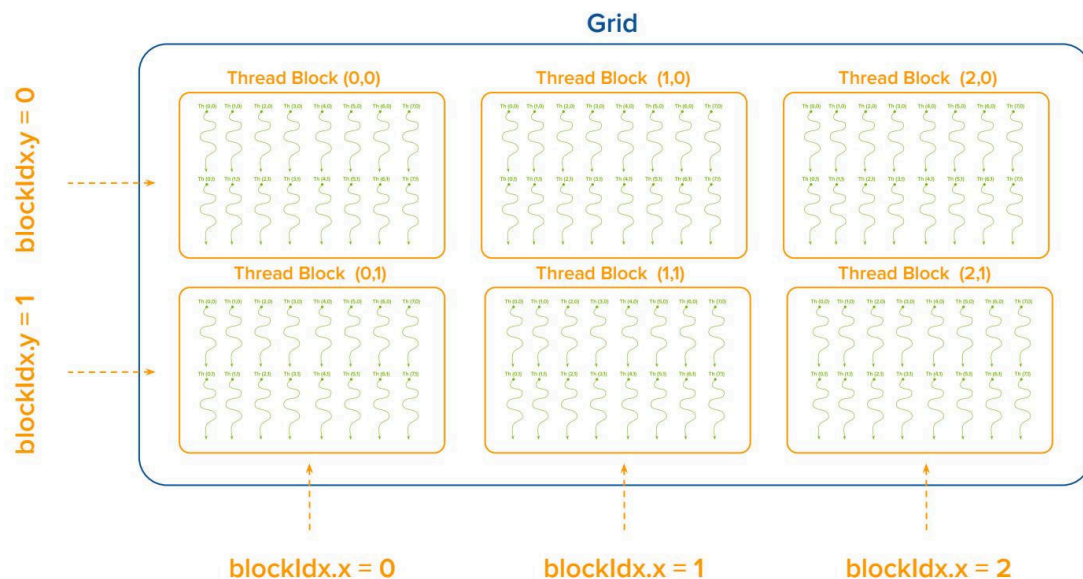
`gridDim.x`, `gridDim.y`: Numero di blocchi nella griglia lungo le dimensioni x e y.



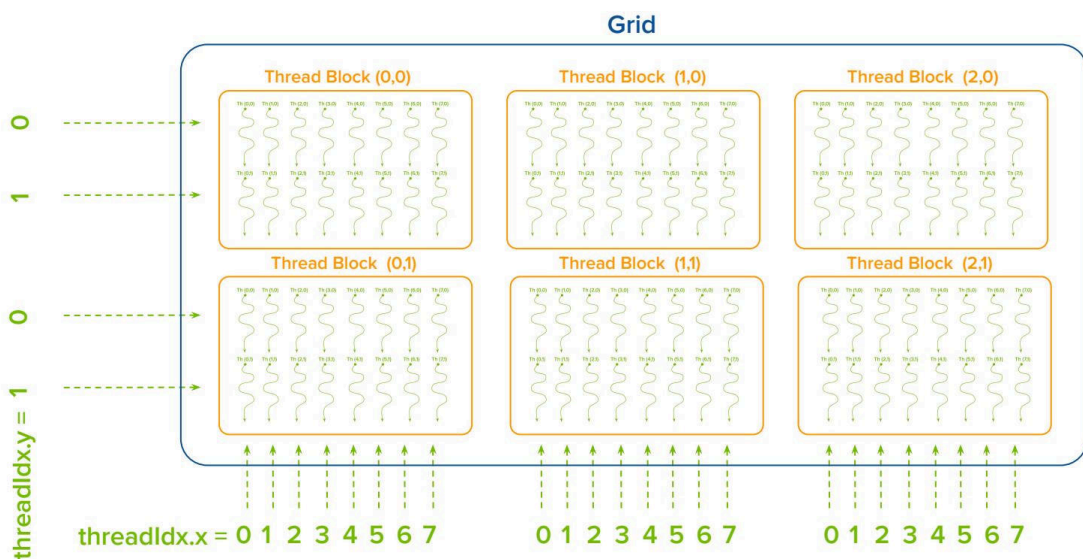
blockDim.x, blockDim.y: Numero di thread per blocco lungo le dimensioni x e y.



blockIdx.x, blockIdx.y: Indici del blocco lungo le dimensioni x e y della griglia.



threadIdx.x, threadIdx.y: Indici x e y del thread nel blocco 2D.



2.17 Esecuzione di un Kernel CUDA

2.17.1 Cos'è un Kernel CUDA?

- Un kernel CUDA è una funzione che viene eseguita in parallelo sulla GPU da migliaia o milioni di thread.
- Rappresenta il nucleo computazionale di un programma CUDA.
- Nei kernel viene definita la logica di calcolo per un singolo thread e l'accesso ai dati associati a quel thread.
- Ogni thread esegue lo stesso codice kernel, ma opera su diversi elementi dei dati.

Sintassi della chiamata Kernel CUDA

```
1 kernel_name <<<gridSize,blockSize>>>(argument list);
```

CUDA

- `gridSize`: Dimensione della griglia (num. di blocchi).
- `blockSize`: Dimensione del blocco (num. di thread per blocco).
- `argument list`: Argomenti passati al kernel.

Sintassi Standard C

```
1 function_name (argument list);
```

C

Con `gridSize` e `blockSize` si definisce:

- Numero totale di thread per un kernel.
- Il layout dei thread che si vuole utilizzare.

Come Eseguiamo il Codice in Parallelo sul Dispositivo?

- Sequenziale (non ottimale): `kernel_name<<<1, 1>>>(args);` // 1 blocco, 1 thread per blocco
- Parallelo: `kernel_name<<<256, 64>>>(args);` // 256 blocchi, 64 thread per blocco

2.18 Qualificatori di Funzione in CUDA

I qualificatori di funzione in CUDA sono essenziali per specificare dove una funzione verrà eseguita e da dove può essere chiamata.

Qualificatore	Esecuzione	Chiamata	Note
<code>__global__</code>	Sul Device	Dall'Host	Deve avere tipo di ritorno <code>void</code>
<code>__device__</code>	Sul Device	Solo dal Device	
<code>__host__</code>	Sull'Host	Solo dall'Host	Può essere omesso

```
1 __global__ void kernelFunction(int data, int size);
```

CUDA

- Funzione kernel (eseguita sulla GPU, chiamabile solo dalla CPU).

```
1 __device__ int deviceHelper(int x);
```

CUDA

- Funzione device (eseguita sulla GPU, chiamabile solo dalla GPU).

```
1 __host__ int hostFunction(int x);
```

CUDA

- Funzione host (eseguibile su CPU).

2.18.1 Combinazione dei qualificatori host e device

In CUDA, combinando `__host__` e `__device__`, una funzione può essere eseguita sia sulla CPU che sulla GPU.

```
1 __host__ __device__ int hostDeviceFunction(int x);
```

CUDA

Permette di scrivere una sola volta funzioni che possono essere utilizzate in entrambi i contesti.

2.19 Kernel CUDA: Regole e Comportamento

1. **Esclusivamente Memoria Device** (`__global__` e `__device__`)
 - Accesso consentito solo alla memoria della GPU. Niente puntatori verso la memoria host.
2. **Ritorno void** (`__global__`)
 - I kernel non restituiscono valori direttamente. La comunicazione con l'host avviene tramite la memoria.

3. **Nessun supporto per argomenti variabili** (`__global__` e `__device__`)
 - I kernel non possono avere un numero variabile di argomenti.
4. **Nessun supporto per variabili statiche locali** (`__global__` e `__device__`)
 - Tutte le variabili devono essere passate come argomenti o allocate dinamicamente.
5. **Nessun supporto per puntatori a funzione** (`__global__` e `__device__`)
 - Non è possibile utilizzare puntatori a funzione all'interno di un kernel.
6. **Comportamento asincrono** (`__global__`)
 - I kernel vengono lanciati in modo asincrono rispetto al codice host, salvo sincronizzazioni esplicite.

2.20 Configurazioni di un Kernel CUDA

2.20.1 Combinazioni di Griglia 1D (Esempi)

La configurazione di griglia e blocchi può essere 1D, 2D o 3D (9 combinazioni in totale), permettendo una mappatura efficiente (ed intuitiva) su array, matrici o dati volumetrici.

```

1 // 1D Grid, 1D Block
2 dim3 gridSize(4);
3 dim3 blockSize(8);
4 kernel_name<<<gridSize, blockSize>>>(args);
5 // 1D Grid, 2D Block
6 dim3 gridSize(4);
7 dim3 blockSize(8, 4);
8 kernel_name<<<gridSize, blockSize>>>(args);
9 // 1D Grid, 3D Block
10 dim3 gridSize(4);
11 dim3 blockSize(8, 4, 2);
12 kernel_name<<<gridSize, blockSize>>>(args);

```

CUDA

Ottimale per problemi con dati strutturati linearmente, come l'elaborazione di **vettori** o **stringhe**, dove ogni thread può lavorare su una porzione contigua dei dati.

Nota: L'efficienza di una configurazione dipende da vari fattori come la dimensione dei dati, l'architettura della GPU e la natura del problema.

2.20.2 Combinazioni di Griglia 2D (Esempi)

La configurazione di griglia e blocchi può essere 1D, 2D o 3D (9 combinazioni in totale), permettendo una mappatura efficiente (ed intuitiva) su array, matrici o dati volumetrici.

```

1 // 2D Grid, 1D Block
2 dim3 gridSize(4, 2);
3 dim3 blockSize(8);
4 kernel_name<<<gridSize, blockSize>>>(args);
5 // 2D Grid, 2D Block
6 dim3 gridSize(4, 2);
7 dim3 blockSize(8, 4);
8 kernel_name<<<gridSize, blockSize>>>(args);
9 // 2D Grid, 3D Block
10 dim3 gridSize(4, 2);

```

CUDA


```

11 dim3 blockSize(8, 4, 2);
12 kernel_name<<<gridSize, blockSize>>>(args);

```

Ideale per problemi con dati strutturati in matrici o immagini, dove ogni thread può gestire un pixel o un elemento della matrice, sfruttando la vicinanza spaziale dei dati.

Nota: L'efficienza di una configurazione dipende da vari fattori come la dimensione dei dati, l'architettura della GPU e la natura del problema.

2.20.3 Combinazioni di Griglia 3D (Esempi)

La configurazione di griglia e blocchi può essere 1D, 2D o 3D (9 combinazioni in totale), permettendo una mappatura efficiente (ed intuitiva) su array, matrici o dati volumetrici.

```

1 // 3D Grid, 1D Block
2 dim3 gridSize(4, 2, 2);
3 dim3 blockSize(8);
4 kernel_name<<<gridSize, blockSize>>>(args);
5 // 3D Grid, 2D Block
6 dim3 gridSize(4, 2, 2);
7 dim3 blockSize(8, 4);
8 kernel_name<<<gridSize, blockSize>>>(args);
9 // 3D Grid, 3D Block
10 dim3 gridSize(4, 2, 2);
11 dim3 blockSize(8, 4, 2);
12 kernel_name<<<gridSize, blockSize>>>(args);

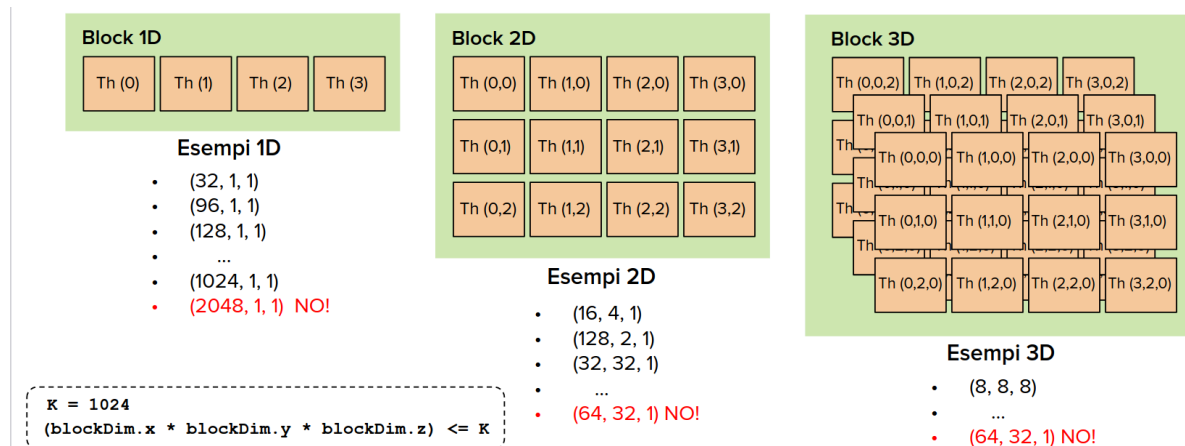
```

CUDA

Ottimale per problemi con **dati volumetrici**, come simulazioni fisiche o rendering 3D, dove ogni thread può operare su un voxel o una porzione dello spazio 3D.

2.21 Numero di Thread per Blocco

- Il **numero massimo** totale di thread per blocco è 1024 per la maggior parte delle GPU (compute capability $\geq 2.x$).
- Un blocco può essere organizzato in 1, 2 o 3 dimensioni, ma ci sono limiti per ciascuna dimensione. Esempio:
 - $x : 1024, y : 1024, z : 64$
- Il prodotto delle dimensioni x, y e z **non** può superare 1024 (queste limitazioni potrebbero cambiare in futuro).



2.22 Compute Capability (CC) - Limiti SM

- La **Compute Capability (CC)** di NVIDIA è un numero che identifica le caratteristiche e le capacità di una GPU NVIDIA in termini di funzionalità supportate e limiti hardware.
- È composta da **due numeri**: il numero principale indica la **generazione** dell'architettura, mentre il numero secondario indica **revisioni** e **miglioramenti** all'interno di quella generazione.

2.23 Identificazione dei Thread in CUDA

2.23.1 Esempio Codice CUDA

```

1  #include <cuda_runtime.h>
2  // Kernel
3  __global__ void kernel_name() {
4      // Accesso alle variabili built-in
5      int blockIdx_x = blockIdx.x, blockIdx_y = blockIdx.y, blockIdx_z = blockIdx.z;
6      int threadIdx_x = threadIdx.x, threadIdx_y = threadIdx.y, threadIdx_z =
7      threadIdx.z;
8      int totalThreads_x = blockDim.x, totalThreads_y = blockDim.y, totalThreads_z =
9      blockDim.z;
10     int totalBlocks_x = gridDim.x, totalBlocks_y = gridDim.y, totalBlocks_z =
11     gridDim.z;
12     // Logica del kernel...
13 }
14
15 int main() {
16     // Definizione delle dimensioni della griglia e del blocco (Caso 3D)
17     dim3 gridDim(4, 4, 2); // 4x4x2 blocchi
18     dim3 blockDim(8, 8, 4); // 8x8x4 thread per blocco
19     // Lancio del kernel
20     kernel_name<<<gridDim, blockDim>>>();
21     // Resto del Programma
22 }

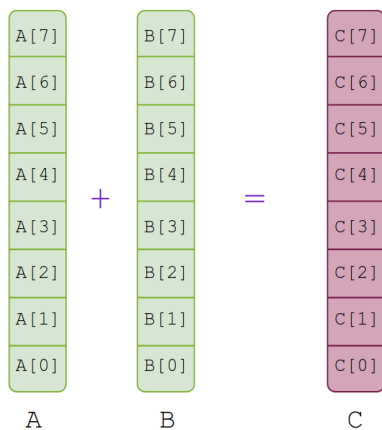
```

CUDA

2.24 Tecniche di Mapping e Dimensionamento

2.24.1 Somma di Array in CUDA

Il Problema: Vogliamo sommare due array elemento per elemento in parallelo utilizzando CUDA.

**Approccio Tradizionale (CPU)**

- Gli elementi degli array vengono elaborati in sequenza, **uno alla volta**.
- Questo approccio è **inefficiente** per array di grandi dimensioni.
- Utilizza solo **un core** della CPU, rallentando il processo.

Approccio CUDA (GPU)

- Gli elementi degli array vengono **elaborati in parallelo**.
- La GPU è progettata per eseguire calcoli **paralleli** su larga scala.
- **Migliaia di core** della GPU lavorano insieme, accelerando enormemente il calcolo.

Approccio Tradizionale (CPU)

- Gli elementi degli array vengono elaborati in sequenza, uno alla volta.
- Questo approccio è inefficiente per array di grandi dimensioni.
- Utilizza solo un core della CPU, rallentando il processo.

Approccio CUDA (GPU)

- Gli elementi degli array vengono elaborati in parallelo.
- La GPU è progettata per eseguire calcoli paralleli su larga scala.
- Migliaia di core della GPU lavorano insieme,

2.25 Confronto: Somma di Vettori in C vs CUDA C**Codice C Standard**

```

1 void sumArraysOnHost(float A, float B,
2 float C, int N) {
3     for (int idx = 0; idx < N; idx++)
4         C[idx] = A[idx] + B[idx];
5 }
6 // Chiamata della funzione
7 sumArraysOnHost(A, B, C, N);

```

C**Caratteristiche**

- **Esecuzione:** Sequenziale
- **Iterazione:** Loop Esplicito
- **Indice:** Variabile di Loop (idx)
- **Scalabilità:** Limitata dalla CPU

Vantaggi

- Portabilità su qualsiasi sistema
- Facilità di debugging

Codice CUDA C

CUDA

```

1 __global__ void sumArraysOnGPU(float A, float B,
2 float C, int N) {
3     int idx = ? // Come accedere ai dati?
4     if (idx < N) C[idx] = A[idx] + B[idx];
5 }
6 // Chiamata del kernel
7 sumArraysOnGPU<<<gridDim,blockDim>>>(A, B, C, N);

```

Caratteristiche

- **Esecuzione:** Parallela
- **Iterazione:** Implicita (un thread per elemento)
- **Indice:** ?
- **Scalabilità:** Elevata (sfrutta molti core GPU)

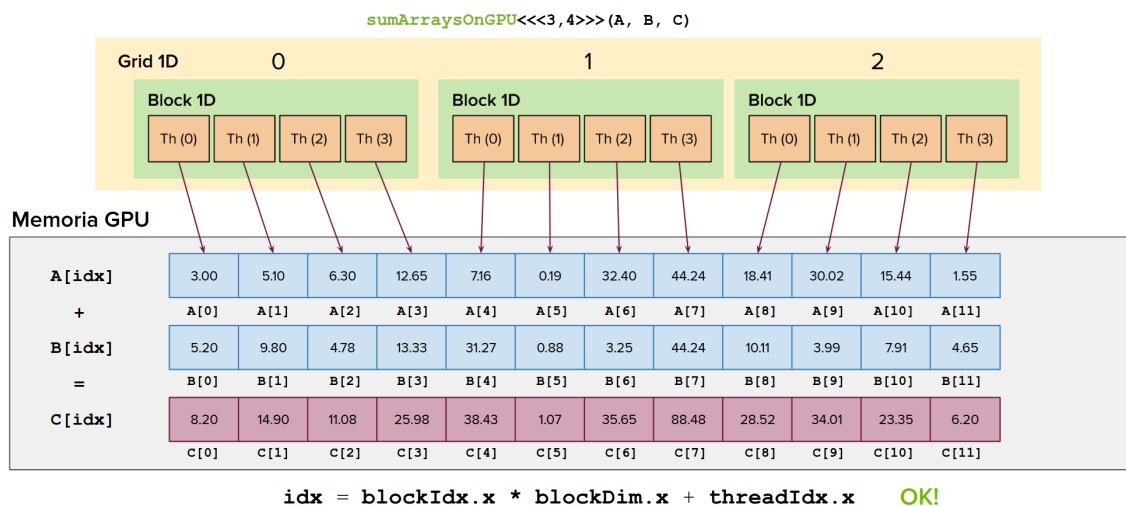
Vantaggi

- Altamente parallelo
- Eccellenti prestazioni su grandi dataset
- Sfrutta la potenza di calcolo delle GPU

Come mappare gli indici dei thread agli elementi dell'array?

idx = threadIdx.x OK! Ma..

Come mappare gli indici dei thread agli elementi dell'array?



Proprietà chiave di questo approccio

- **Copertura completa:** Tutti i 12 thread (3 blocchi x 4 thread per blocco) sono utilizzati per elaborare i 12 elementi degli array.
- **Mapping corretto:** Ogni thread è associato a un unico elemento degli array A, B e C.
- **Nessuna ripetizione:** L'indice idx, univoco per ogni thread, assicura che ogni elemento dell'array venga elaborato esattamente una volta, evitando ridondanze.
- **Parallelismo massimizzato:** La formula idx permette di sfruttare appieno il parallelismo della GPU, assegnando un compito specifico ad ogni thread disponibile.

- **Scalabilità:** Questa formula si adatta bene a dimensioni di array diverse, purché si adegui il numero di blocchi.
- **Bilanciamento del carico:** Il lavoro è distribuito uniformemente tra tutti i thread, garantendo un utilizzo efficiente delle risorse.
- **Accessi coalescenti:** I thread adiacenti in un blocco accedono a elementi di memoria adiacenti, favorendo accessi coalescenti e migliorando l'efficienza della memoria.

Quindi il codice CUDA sarà il seguente:

Codice CUDA C

```

1 __global__ void sumArraysOnGPU(float A, float B,
2 float C, int N) {
3     int idx = blockDim.x*blockIdx.x + threadIdx.x;
4     if (idx < N) C[idx] = A[idx] + B[idx];
5 }
6 // Chiamata del kernel
7 sumArraysOnGPU<<<gridDim,blockDim>>>(A, B, C, N);

```

CUDA

2.26 Identificazione dei Thread e Mapping dei Dati in CUDA

Le variabili di identificazione sono accessibili solo all'interno del kernel e permettono ai thread di conoscere la propria posizione all'interno della gerarchia e di adattare il proprio comportamento di conseguenza.

Perché Identificare i Thread?

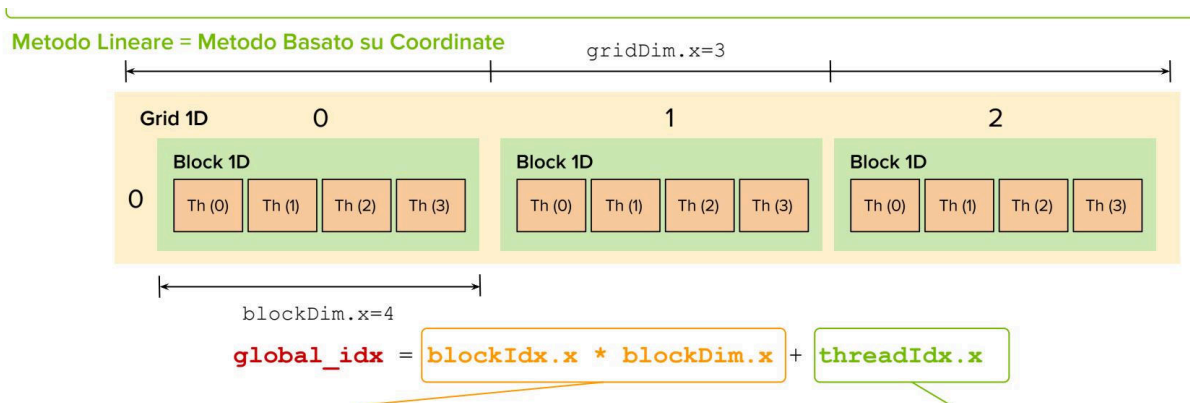
- L'indice globale di un thread identifica univocamente quale porzione di dati deve elaborare.
- Essenziale per gestire correttamente l'accesso alla memoria e coordinare l'esecuzione di algoritmi complessi.

2.26.1 Struttura dei Dati e Calcolo dell'Indice Globale

- Anche le strutture più complesse, come matrici (2D) o array tridimensionali (3D), vengono memorizzate come una sequenza di elementi contigui in memoria nella GPU, tipicamente organizzati in array lineari.
- Ogni thread elabora uno o più elementi in base al proprio indice globale.
- Esistono diversi metodi per calcolare l'indice globale di un thread (es. Metodo Lineare, Coordinate-based).
- Metodi diversi possono produrre mappature diverse tra thread e dati, influenzando prestazioni (come la coalescenza degli accessi in memoria) e la leggibilità del codice.

2.26.2 Calcolo dell'Indice Globale del Thread - Grid 1D, Block 1D

In CUDA, ogni thread ha un indice globale (`global_idx`) che lo identifica nell'esecuzione del kernel. Il programmatore lo calcola usando l'indice del thread nel blocco e l'indice del blocco nella griglia.



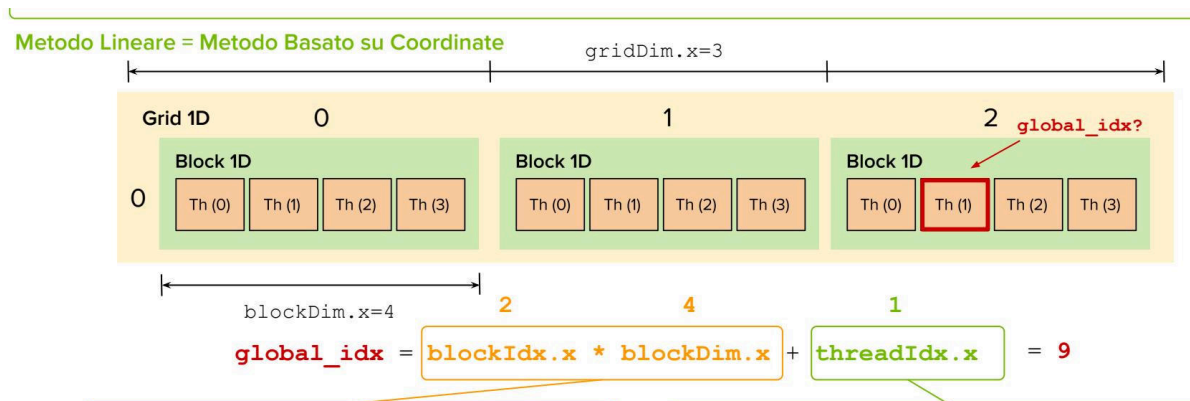
$\text{blockIdx.x} * \text{blockDim.x}$

- Calcola l'offset di tutti i thread nei blocchi precedenti al blocco corrente.
- Moltiplicando blockIdx.x per blockDim.x , otteniamo il numero totale di thread che si trovano nei blocchi precedenti.

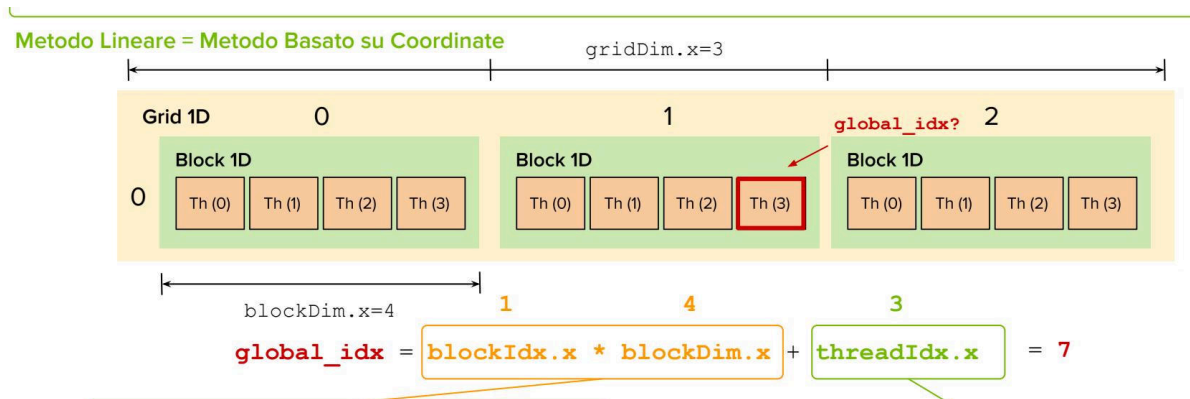
threadIdx.x

- Identifica la posizione del thread all'interno del blocco corrente.
- È l'indice del thread all'interno del blocco corrente, da 0 a $\text{blockDim.x} - 1$.

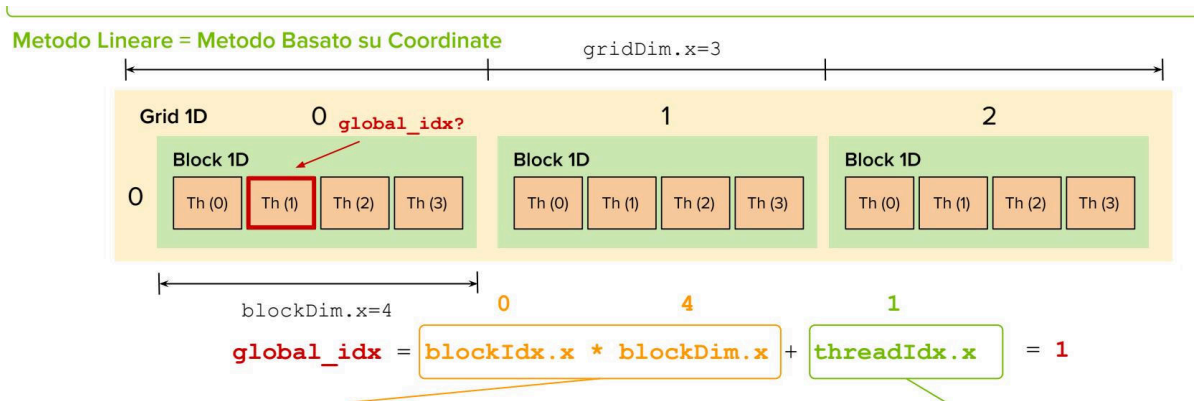
In questo esempio viene mostrato come calcolare l'indice globale per il thread Th(1) appartenente al blocco unidimensionale con indice $\text{blockIdx.x} = 2$



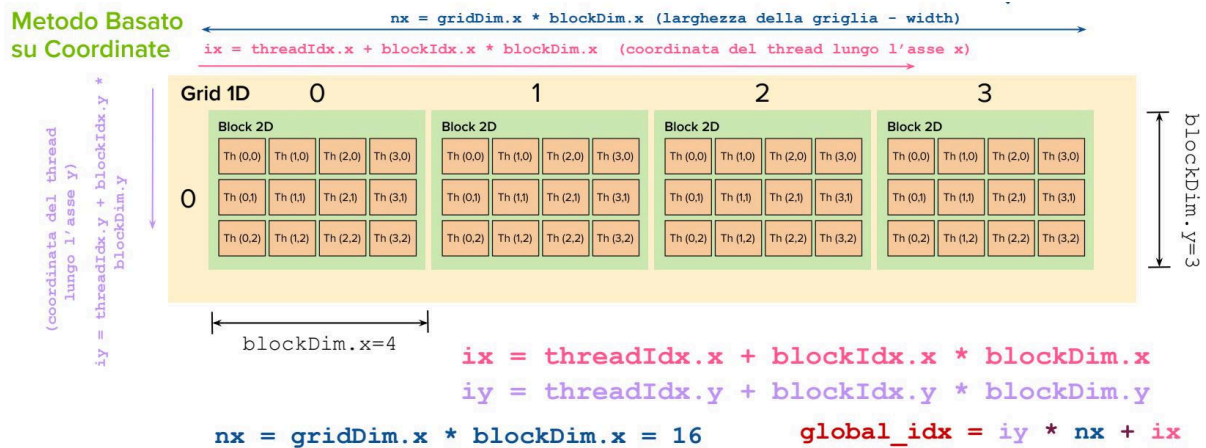
In questo esempio viene mostrato come calcolare l'indice globale per il thread Th(3) appartenente al blocco unidimensionale con indice $\text{blockIdx.x} = 1$



In questo esempio viene mostrato come calcolare l'indice globale per il thread Th(1) appartenente al blocco unidimensionale con indice `blockIdx.x = 0`



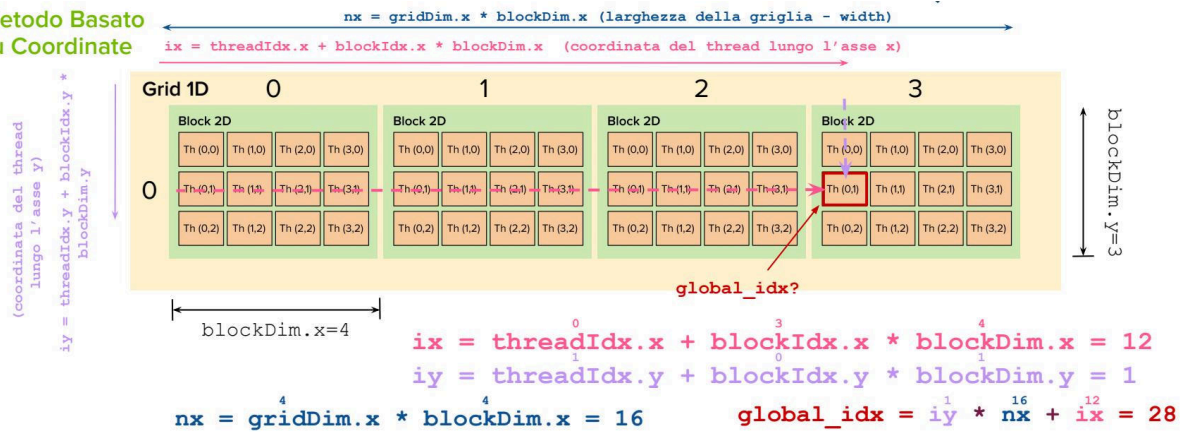
2.27 Calcolo dell'Indice Globale del Thread - Grid 1D, Block 2D



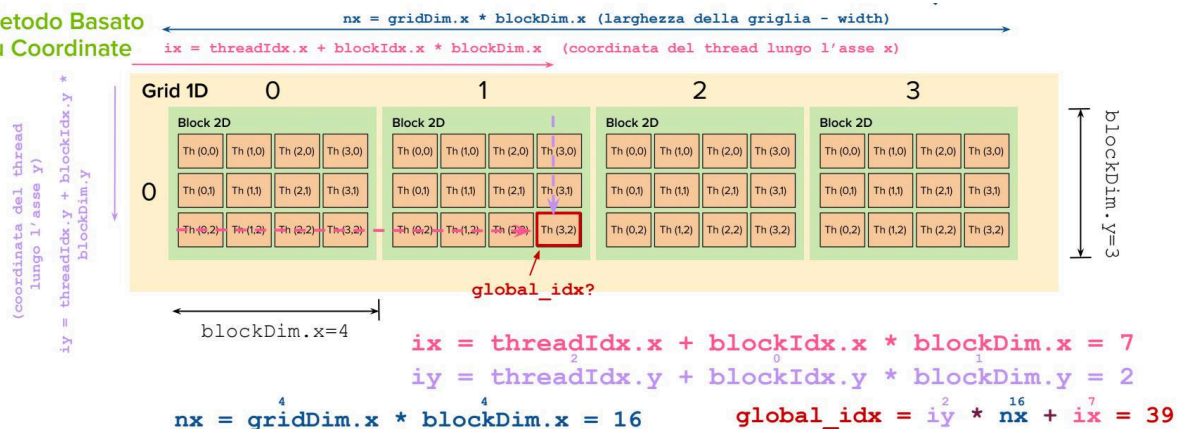
- $ix = \text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x}$: determina l'indice del thread lungo l'asse x, prendendo in considerazione la posizione nel blocco (`threadIdx.x`) e il numero di blocchi precedenti (`blockIdx.x * blockDim.x`).
- $iy = \text{threadIdx.y} + \text{blockIdx.y} * \text{blockDim.y}$: determina l'indice del thread lungo l'asse y, considerando sia la posizione locale (`threadIdx.y`) che i blocchi precedenti lungo y (`blockIdx.y * blockDim.y`).
- $\text{global_idx} = iy * nx + ix$: calcola l'indice globale sommando `ix` all'indice globale lungo y, dove `nx` rappresenta il numero di thread per riga (in questo caso, $nx = \text{gridDim.x} * \text{blockDim.x}$).

Seguono alcuni esempi.

Metodo Basato su Coordinate

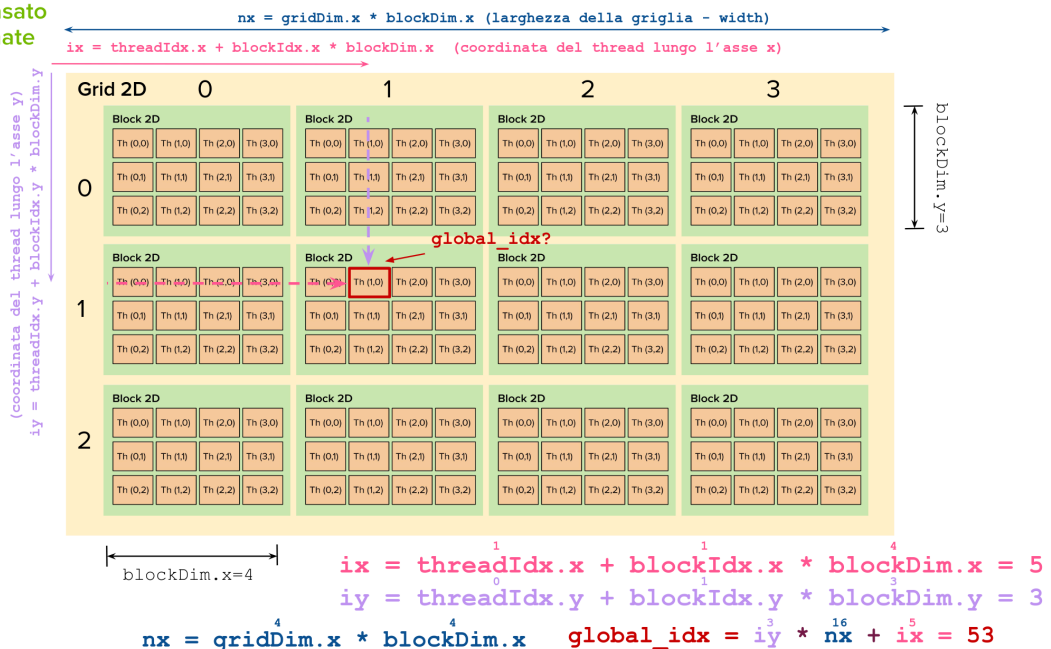


Metodo Basato su Coordinate



2.28 Calcolo dell'Indice Globale del Thread - Grid 2D, Block 2D

Metodo Basato su Coordinate



2.29 Metodo Basato su Coordinate per Indici Globali in CUDA

Caratteristiche del Metodo Basato su Coordinate

- Calcola indici separati per ogni dimensione della griglia e dei blocchi.
- Riflette naturalmente la disposizione multidimensionale dei dati.

- Facilita la comprensione della posizione del thread nello spazio
- Richiede un passaggio aggiuntivo per combinare gli indici in un indice globale.

2.29.1 Calcolo degli Indici Coordinati

Calcolo degli Indici Coordinati

Caso 1D) $x = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x} \rightarrow \text{idx} = x$ (equivalente al caso lineare)

Caso 2D) $x = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$
 $y = \text{blockIdx.y} * \text{blockDim.y} + \text{threadIdx.y} \rightarrow \text{idx} = y * \overset{(nx)}{\text{width}} + x$

Caso 3D) $x = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$
 $y = \text{blockIdx.y} * \text{blockDim.y} + \text{threadIdx.y}$
 $z = \text{blockIdx.z} * \text{blockDim.z} + \text{threadIdx.z} \rightarrow \text{idx} = z * \overset{(ny)}{(\text{height} * \overset{(nx)}{\text{width}})} + y * \text{width} + x$

Calcolo dell'Indice Globale

2.29.2 Esempio di Utilizzo (Caso 2D)

```
1 __global__ void kernel2D(float data, int width, int height) {
2     int x = blockIdx.x * blockDim.x + threadIdx.x;
3     int y = blockIdx.y * blockDim.y + threadIdx.y;
4     if (x < width && y < height) { // width e height si riferiscono alle dimensioni
        dell'array dati
5         int idx = y * width + x;
6         // Operazioni su data[global_idx]
7     }
8 }
```

CUDA

2.29.3 Come Calcolare la Dimensione della Griglia e del Blocco

Approccio Generale

- Definire manualmente prima la dimensione del blocco (numero di thread per blocco).
- Poi, calcolare automaticamente la dimensione della griglia in base ai dati e alla dimensione del blocco.

Motivazioni

- La **dimensione del blocco** è legata alle **caratteristiche hardware** della GPU e la natura del problema.
- La **dimensione della griglia** si adatta alla **dimensione del blocco** e al **volume dei dati** da processare.

Calcolo delle Dimensioni (Caso 1D)

```
1 int blockSize = 256; int dataSize = 1024; // Dimensione del blocco e dei dati
2 dim3 blockDim(blockSize); dim3 gridDim((dataSize + blockSize - 1) / blockSize);
3 kernel_name<<<gridDim, blockDim>>>(args); // Lancio del kernel
```

CUDA

Spiegazione del Calcolo

La formula $(\text{dataSize} + \text{blockSize} - 1) / \text{blockSize}$ assicura un numero sufficiente di blocchi per coprire tutti i dati, anche se `dataSize` non è un multiplo esatto di `blockSize`.

- **Divisione semplice:** $\text{dataSize} / \text{blockSize}$ fornisce il numero di blocchi completamente pieni.
- Se ci sono dati residui che non riempiono un intero blocco, la divisione semplice li ignorerebbe.
- Aggiungere $\text{blockSize} - 1$ a `dataSize` “compensa” questi dati residui, includendo l’ultimo blocco parziale. Equivalente a calcolare la ceil della divisione.

2.29.4 Esempio 1 (Dati Residui): dataSize = 1030, blockSize = 256**Calcolo delle Dimensioni (Caso 1D)**

CUDA

```

1 int blockSize = 256; int dataSize = 1030; // Dimensione del blocco e dei dati
2 dim3 blockDim(blockSize); dim3 gridDim((dataSize + blockSize - 1) / blockSize);
3 kernel_name<<<gridDim, blockDim>>>(args); // Lancio del kernel

```

- **Divisione semplice:** $1030/256 = 4$ blocchi; **ignorerebbe** l'ultimo blocco parziale perché

$$256 * 4 = 1024$$

$$1030 - 1024 = 6$$

quindi 6 elementi residui ❌

- Con la formula $(1030 + 256 - 1)/256 = 1285/256 = 5$ blocchi; nessun elemento residuo ✅
- In questo caso, la divisione semplice avrebbe dato 4 blocchi, ma c'è un residuo di 6 elementi ($1030 \bmod 256 = 6$); la formula include anche il blocco parziale, quindi otteniamo 5 blocchi.

2.29.5 Esempio 2 (Multiplo Perfetto): dataSize = 1024, blockSize = 256

- **Divisione semplice:** $1024/256 = 4 \rightarrow$ copre esattamente 1024 elementi ✅
- Con la formula $(1024 + 256 - 1)/256 = 1279/256 = 4$ blocchi ✅
- **Spiegazione:** $1279/256 = 4.996$, ma essendo divisione intera $\rightarrow 4$ blocchi; aggiungere 255 non basta per raggiungere 1280 (soglia del 5° blocco), quindi non si arrotonda per eccesso.
- **Importante:** La formula funziona correttamente sia con residui che senza, garantendo sempre il numero esatto di blocchi necessari senza sprechi

2.29.6 Calcolo delle Dimensioni (Caso 2D)

```

1 int blockSizeX = 16, blockSizeY = 16; // Dimensione del blocco
2 int dataSizeX = 1024, dataSizeY = 512; // Dimensione dei dati
3 dim3 blockDim(blockSizeX, blockSizeY); // Definizione del blocco 2D
4 dim3 gridDim( // Calcolo della griglia 2D
5   (dataSizeX + blockSizeX - 1) / blockSizeX, // Numero di blocchi in X
6   (dataSizeY + blockSizeY - 1) / blockSizeY // Numero di blocchi in Y
7 );
8 kernel_name<<<gridDim, blockDim>>>(args); // Lancio del kernel

```

CUDA

2.29.7 Calcolo delle Dimensioni (Caso Generale 3D)

```

1 int blockSizeX = 16, blockSizeY = 16, blockSizeZ = 16; // Dimensione del
  blocco
2 int dataSizeX = 1024, dataSizeY = 512, dataSizeZ = 256; // Dimensione dei dati
3 dim3 blockDim(blockSizeX, blockSizeY, blockSizeZ); // Definizione del blocco 3D
4 dim3 gridDim( // Calcolo della griglia 3D
5   (dataSizeX + blockSizeX - 1) / blockSizeX, // Numero di blocchi in X
6   (dataSizeY + blockSizeY - 1) / blockSizeY, // Numero di blocchi in Y
7   (dataSizeZ + blockSizeZ - 1) / blockSizeZ // Numero di blocchi in Z
8 );
9 kernel_name<<<gridDim, blockDim>>>(args); // Lancio del kernel

```

CUDA

2.30 Analisi delle Prestazioni

2.30.1 Verifica del Kernel CUDA (Somma di Array)

Il controllo dei kernel CUDA mira a confermare l'affidabilità dei calcoli eseguiti sulla GPU.

```

1 void checkResult(float *hostRef, float *gpuRef, const int N) {
2     double epsilon = 1.0E-8;
3     int match = 1;
4     for (int i = 0; i < N; i++) {
5         if (fabsf(hostRef[i] - gpuRef[i]) > epsilon)
6         {
7             match = 0;
8             printf("Arrays do not match!\n");
9             printf("host %5.2f gpu %5.2f at current %d\n", hostRef[i], gpuRef[i], i);
10            break;
11        }
12    }
13    if (match) printf("Arrays match.\n\n");
14 }

```

- `hostRef`: risultati attesi dalla somma
- `gpuRef`: risultati calcolati dal kernel
- `fabsf`: funzione della libreria C che calcola il valore assoluto di un numero in virgola mobile a precisione singola

Suggerimenti per la Verifica (basic)

- Verifica ogni elemento degli array per assicurarsi che i risultati del kernel corrispondano ai valori attesi.
- Usa una piccola tolleranza (`epsilon`) per confronti in virgola mobile, in quanto ci possono essere errori di arrotondamento legate alla natura delle rappresentazioni numeriche nei computer.
- **(Alternativa) Configurazione <<< 1, 1>>>:**
 - Forza l'esecuzione del kernel con un solo blocco e un thread.
 - Emula un'implementazione sequenziale.

2.30.2 Gestione degli Errori in CUDA

Problema

- **Asincronicità:** molte chiamate CUDA sono asincrone, rendendo difficile associare un errore alla specifica chiamata che lo ha causato.
- **Complessità di Debugging:** l'errore può emergere in una parte del programma diversa e lontana dal punto in cui è stato generato, rendendo l'individuazione della causa complicata.
- **Gestione Manuale:** controllare ogni chiamata CUDA manualmente è tedioso e soggetto a errori.

Macro CHECK

CUDA

```

1 // Fornisce file, riga, codice e descrizione dell'errore.
2 #define CHECK(call){
3     const cudaError_t error = call;
4     if (error != cudaSuccess){
5         printf("Error: %s:%d, ", __FILE__, __LINE__);
6         printf("code:%d, reason: %s\n", error, cudaGetErrorString(error));
7         exit(1);
8     }
9 }

```

Esempi di Utilizzo

CUDA

```

1 CHECK(cudaMalloc(&d_input, size)); // Allocazione
2 CHECK(cudaMemcpy(d_input, h_input, size, cudaMemcpyHostToDevice));
3 kernel_function <<<numBlocks, blockSize >>>(argument list); // Lancia il kernel
4 CHECK(cudaGetLastError()); // Primo controllo: errori di lancio del kernel
5 CHECK(cudaDeviceSynchronize()); // Secondo controllo: errori durante l'esecuzione
  del kernel. Usare solo in DEBUG (Overhead di performance!)

```

2.30.3 Profiling delle Prestazioni in CUDA**Introduzione al Profiling**

- Misurare e ottimizzare le prestazioni è fondamentale per garantire l'efficienza del codice.
- Il profiling permette di misurare le prestazioni, analizzare l'uso delle risorse e individuare possibili ottimizzazioni.

Importanza della Misurazione del Tempo

- **Identificazione dei Colli di Bottiglia:** Individuare le sezioni di codice che limitano le prestazioni (un'implementazione naïve raramente è ottimale)
- **Analisi degli Effetti delle Modifiche:** valutare l'impatto delle modifiche sul tempo di esecuzione.
- **Confronto tra Implementazioni:** valutare le prestazioni tra diverse strategie di implementazione.
- **Analisi del Bilanciamento Carico/Calcolo:** verificare se il carico di lavoro è distribuito in modo efficiente tra thread, blocchi e host-device.

2.30.4 Metodi Principali

1. **Timer CPU:** Semplice e diretto, utilizza funzioni di sistema per ottenere il tempo di esecuzione.
2. **NVIDIA Profiler** (deprecato): Strumento da riga di comando per analizzare attività di CPU e GPU.
3. **NVIDIA Nsight Systems e Nsight Compute:** Strumenti moderni avanzati per analisi approfondita e ottimizzazione a livello di sistema e kernel.

2.30.4.1 Timer CPU

- Timer eseguito dall'host, misura il tempo **wall-clock** visto dalla CPU.
- Soluzione semplice e pratica basata su funzioni di sistema standard.
- **Può misurare qualsiasi operazione:** kernel GPU, trasferimenti memoria, codice CPU.

- Il tempo include anche gli overhead: lancio dei kernel, sincronizzazione, latenze di comunicazione.
- Per operazioni GPU asincrone (kernel), richiede sincronizzazione esplicita.

Funzione del Timer della CPU

CUDA

```

1 #include <time.h>
2 double cpuSecond() {
3     struct timespec ts;
4     timespec_get(&ts, TIME_UTC);
5     return ((double)ts.tv_sec + (double)ts.tv_nsec * 1.e-9);
6 }
```

- La funzione utilizza `timespec_get()` per ottenere il tempo corrente del sistema.
- Restituisce il tempo in secondi, combinando secondi e nanosecondi.
- Precisione teorica fino al **nanosecondo**.

Utilizzo Per Misurare un Kernel CUDA

CUDA

```

1 double iStart = cpuSecond(); // Registra il tempo di inizio
2 kernel_name<<<grid, block>>>(argument list); // Lancia il kernel CUDA
3 cudaDeviceSynchronize(); // Attende il completamento del kernel
4 double iElaps = cpuSecond() - iStart; // Calcola il tempo trascorso
```

- La chiamata a `cudaDeviceSynchronize()` è cruciale per assicurare che tutto il lavoro sulla GPU sia completato prima di misurare il tempo finale. Questo è necessario poiché le chiamate ai kernel CUDA sono **asincrone** rispetto all'host (senza riflettere solo il tempo di lancio del kernel).
- Il tempo misurato include sia l'**esecuzione** sia l'**overhead** di lancio e sincronizzazione.

Pro

- Facile da implementare e utilizzare.
- Non richiede librerie CUDA **specifiche** per il timing.
- Funziona su **qualsiasi sistema** con supporto CUDA.
- Efficace per **kernel lunghi** e **misure approssimative**.

Contro

- Impreciso per kernel molto brevi ($< 1ms$).
- Include **overhead** non relativo all'esecuzione del kernel (es., sistema operativo, utilizzo CPU, etc.).
- Non fornisce dettagli sulle **fasi interne** del kernel.
- Precisione influenzata dal **carico dell'host**.

Somma di due array

CUDA

```

1  int main(int argc, char **argv)
2  {
3      printf("%s Starting...\n", argv[0]);
4
5      // Configurazione del device
6      int dev = 0;
7      cudaDeviceProp deviceProp;
8      CHECK(cudaGetDeviceProperties(&deviceProp, dev)); // Ottiene le proprietà del
device
9      printf("Using Device %d: %s\n", dev, deviceProp.name);
10     CHECK(cudaSetDevice(dev)); // Seleziona il device CUDA da utilizzare
11
12     // Dimensione dei vettori
13     int nElem = 1 << 24; // 2^24 elementi (16M) - bit shifting
14     printf("Vector size %d\n", nElem);
15
16     // Allocazione della memoria host
17     size_t nBytes = nElem * sizeof(float);
18     float *h_A, *h_B, *hostRef, *gpuRef;
19     h_A = (float *)malloc(nBytes); // Alloca memoria per il vettore A su host
20     h_B = (float *)malloc(nBytes); // Alloca memoria per il vettore B su host
21     hostRef = (float *)malloc(nBytes); // Alloca memoria per il risultato
calcolato su host
22     gpuRef = (float *)malloc(nBytes); // Alloca memoria per il risultato calcolato
su GPU
23
24     // Inizializzazione dei dati su host
25     double iStart, iElaps;
26     iStart = cpuSecond();
27     initData(h_A, nElem); // Inizializza il vettore A
28     initData(h_B, nElem); // Inizializza il vettore B
29     iElaps = cpuSecond() - iStart;
30     printf("Data initialization time: %f sec\n", iElaps);
31     memset(hostRef, 0, nBytes); // Inizializza a zero il vettore risultato su host
32     memset(gpuRef, 0, nBytes); // Inizializza a zero il vettore risultato della
GPU
33
34     // Somma dei vettori su host per verifica
35     iStart = cpuSecond();
36     sumArraysOnHost(h_A, h_B, hostRef, nElem); // Calcola la somma su CPU per
confronto
37     iElaps = cpuSecond() - iStart;
38     printf("sumArraysOnHost elapsed %f sec\n", iElaps);
39     // Allocazione della memoria su device
40     float *d_A, *d_B, *d_C;

```

```

1  CHECK(cudaMalloc((float**)&d_A, nBytes)); // Alloca memoria per A su GPU
2  CHECK(cudaMalloc((float**)&d_B, nBytes)); // Alloca memoria per B su GPU
3  CHECK(cudaMalloc((float**)&d_C, nBytes)); // Alloca memoria per il risultato
   su GPU
4  // Copia dei dati dall'host al device
5  CHECK(cudaMemcpy(d_A, h_A, nBytes, cudaMemcpyHostToDevice)); // Copia A su GPU
6  CHECK(cudaMemcpy(d_B, h_B, nBytes, cudaMemcpyHostToDevice)); // Copia B su GPU
7
8  // Configurazione del kernel
9  dim3 block(1024); // Dimensione del blocco: 1024 threads
10 dim3 grid((nElem + block.x - 1) / block.x); // Calcola il numero di blocchi
   necessari
11
12 // Esecuzione del kernel su device
13 iStart = cpuSecond();
14 sumArraysOnGPU<<<grid, block>>>(d_A, d_B, d_C, nElem); // Lancia il kernel
   CUDA
15 CHECK(cudaDeviceSynchronize()); // Attende il completamento del kernel
16 iElaps = cpuSecond() - iStart;
17 printf("sumArraysOnGPU <<<%d, %d>>> elapsed %f sec\n", grid.x, block.x,
   iElaps);
18
19 // Copia dei risultati dal device all'host
20 CHECK(cudaMemcpy(gpuRef, d_C, nBytes, cudaMemcpyDeviceToHost)); // Copia dalla
   GPU
21
22 // Verifica dei risultati
23 checkResult(hostRef, gpuRef, nElem); // Confronta i risultati di CPU e GPU
24
25 // Liberazione della memoria su device
26 CHECK(cudaFree(d_A)); // Libera la memoria di A su GPU
27 CHECK(cudaFree(d_B)); // Libera la memoria di B su GPU
28 CHECK(cudaFree(d_C)); // Libera la memoria del risultato su GPU
29
30 // Liberazione della memoria su host
31 free(h_A); // Libera la memoria di A su host
32 free(h_B); // Libera la memoria di B su host
33 free(hostRef); // Libera la memoria del risultato CPU su host
34 free(gpuRef); // Libera la memoria del risultato GPU su host
35 return 0;}
```

```
int dev = 0
```

- Rappresenta l'indice della GPU NVIDIA che si intende utilizzare.
- 0 solitamente si riferisce al primo dispositivo CUDA disponibile nel sistema

Alternative e Pratiche Comuni

1. Selezione del dispositivo tramite argomenti:

```
if (argc > 1) dev = atoi(argv[1]);
```

2. Utilizzo di variabili d'ambiente:

```
1 char* deviceIndex = getenv("CUDA_VISIBLE_DEVICES");
2 if (deviceIndex) dev=atoi(deviceIndex);
```

CUDA

```
cudaSetDevice(dev)
```

- Imposta il dispositivo CUDA attivo per le operazioni successive
- Assicura che tutte le allocazioni e le operazioni CUDA successive utilizzino questo dispositivo specifico.

```
initialData(h_A, nElem); // Inizializza il vettore A
initialData(h_B, nElem); // Inizializza il vettore B
```

- Inizializza un array di float con valori casuali (compresi fra 0 e 25.5)

```
1 void initialData(float *ip, int size){
2     // Genera dati casuali
3     time_t t;
4     srand((unsigned int) time(&t));
5     for (int i = 0; i < size; i++){
6         ip[i] = (float)(rand() & 0xFF) / 10.0f;
7     }
8 }
```

CUDA

APPENDICES